

Time Series Forecasting-Sparkling Wine

Problem Statement:

ABC Estate Wines has been a leader in the wine industry for many years, offering high-quality wines to consumers all around the world. As the company continues to expand its reach and grow its customer base, it is essential to analyze market trends and forecast future sales to ensure continued success.

In this report, we will focus on analyzing the sales data for sparkling wine in the 20th century. As an analyst for ABC Estate Wines, I have been tasked with reviewing this data to identify patterns, trends, and opportunities for growth in the sparkling wine market. This knowledge will help us to make informed decisions about how to position our products in the market, optimize our sales strategies, and forecast future sales trends.

Overall, this report aims to provide valuable insights into the sparkling wine market and how ABC Estate Wines can continue to succeed in this highly competitive industry.

1. Read the data as an appropriate Time Series data and plot the data.

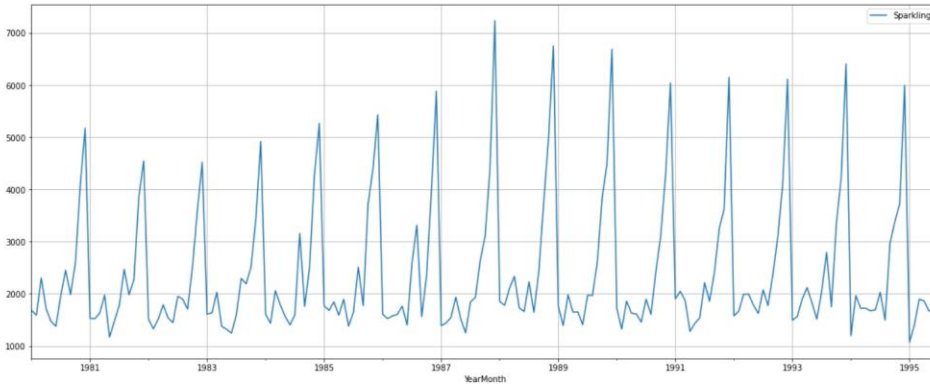
Sparkling		Sparkling	
YearMonth		YearMonth	
1980-01-01	1686	1995-03-01	1897
1980-02-01	1591	1995-04-01	1862
1980-03-01	2304	1995-05-01	1670
1980-04-01	1712	1995-06-01	1688
1980-05-01	1471	1995-07-01	2031

We can see data is from 1980 to 1995.

The number of Rows and Columns of the Dataset:

The dataset has 187 rows and 1 column.

Plot of the dataset:



Post Ingestion of Dataset:

We have divided the dataset further by extraction month and year columns from the YearMonth column and renamed the sparkling column name to Sales for better analysis of the dataset. The new dataset has 187 rows and 3 columns.

Rows of the new data set;

Sparkling		Year	Month				
YearMonth				Sales	Year	Month	
1980-01-01	1686	1980	1	YearMonth			
1980-02-01	1591	1980	2	1995-03-01	1897	1995	3
1980-03-01	2304	1980	3	1995-04-01	1862	1995	4
1980-04-01	1712	1980	4	1995-05-01	1670	1995	5
1980-05-01	1471	1980	5	1995-06-01	1688	1995	6
				1995-07-01	2031	1995	7

2. Perform appropriate Exploratory Data Analysis to understand the data and also perform decomposition.

Data Type;

Index: DateTime

Sales: integer

Month: integer Year: integer

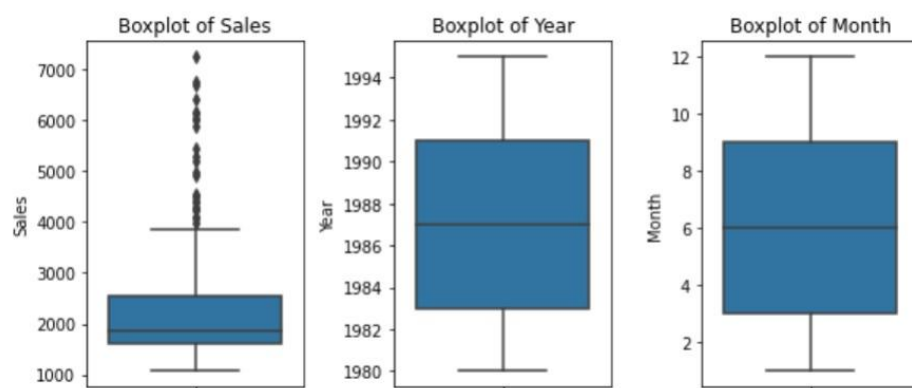
Statistical summary:

	count	mean	std	min	25%	50%	75%	max
Sales	187.0	2402.0	1295.0	1070.0	1605.0	1874.0	2549.0	7242.0
Year	187.0	1987.0	5.0	1980.0	1983.0	1987.0	1991.0	1995.0
Month	187.0	6.0	3.0	1.0	3.0	6.0	9.0	12.0

Null Value:

There are no null values present in the dataset. So we can do further analysis smoothly.

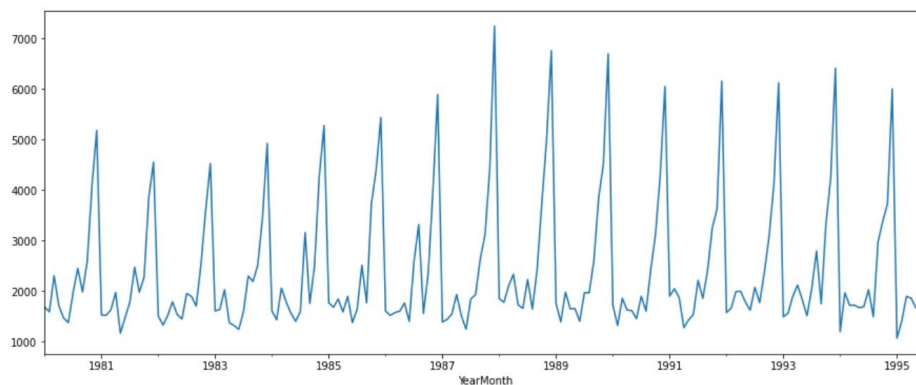
Boxplot of dataset:



The box plot shows:

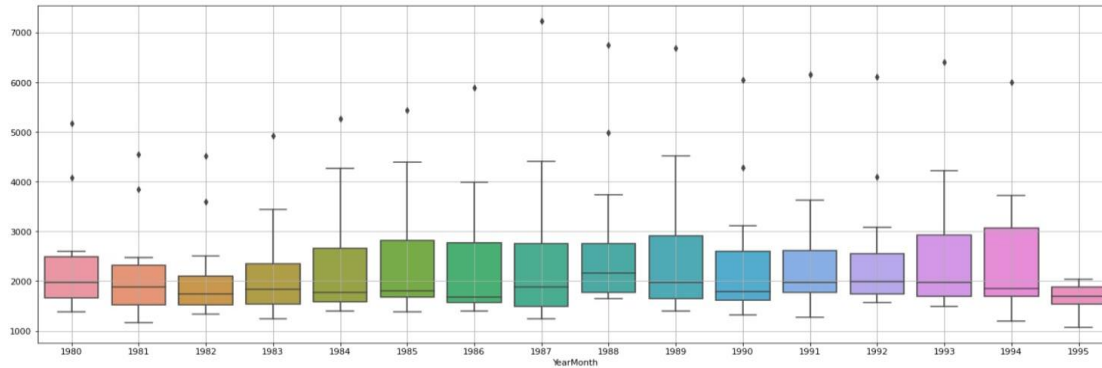
- Sales boxplot has outliers we can treat them but we are choosing not to treat them as they do not give much effect on the time series model.

Line plot of sales:



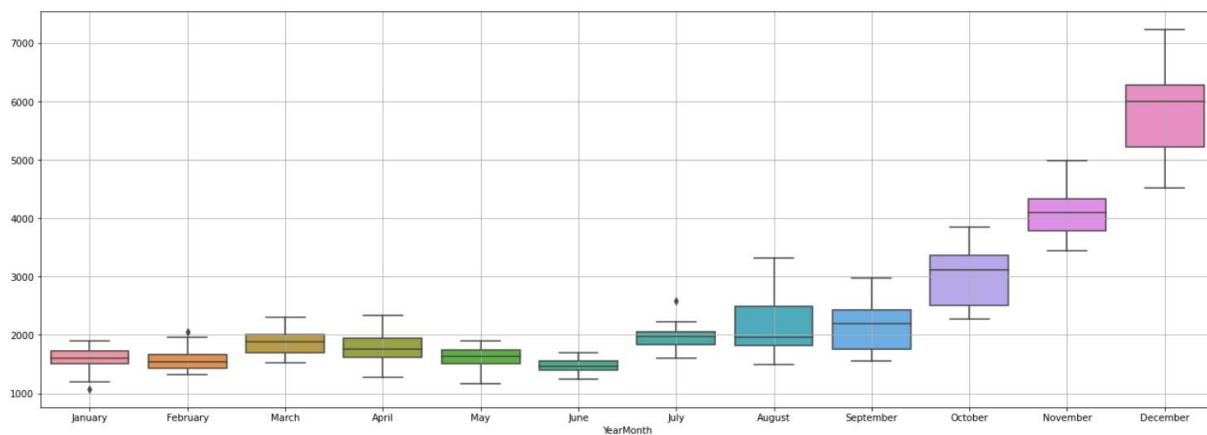
The line plot shows the patterns of trend and seasonality and also shows that there was a peak in the year 1988.

Boxplot Yearly:



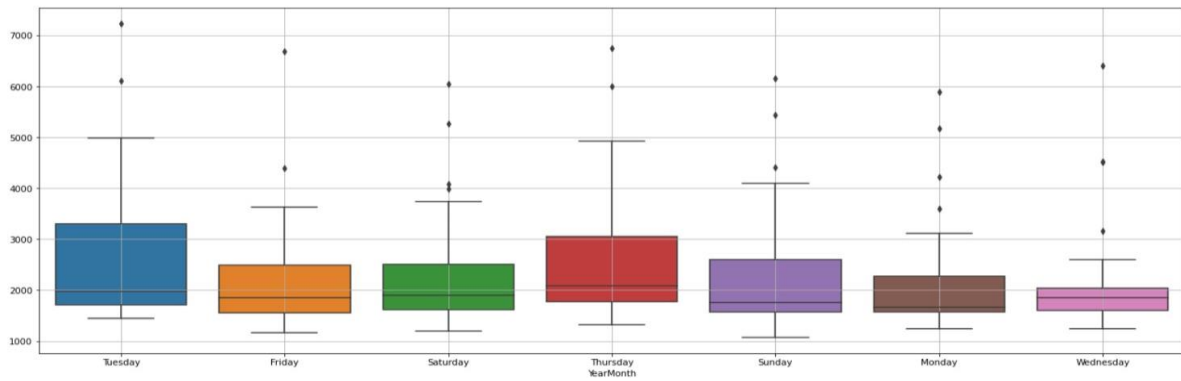
This yearly box plot shows there is consistency over the years and there was a peak in 1988-1989. Outliers are present in all years.

Boxplot Monthly:



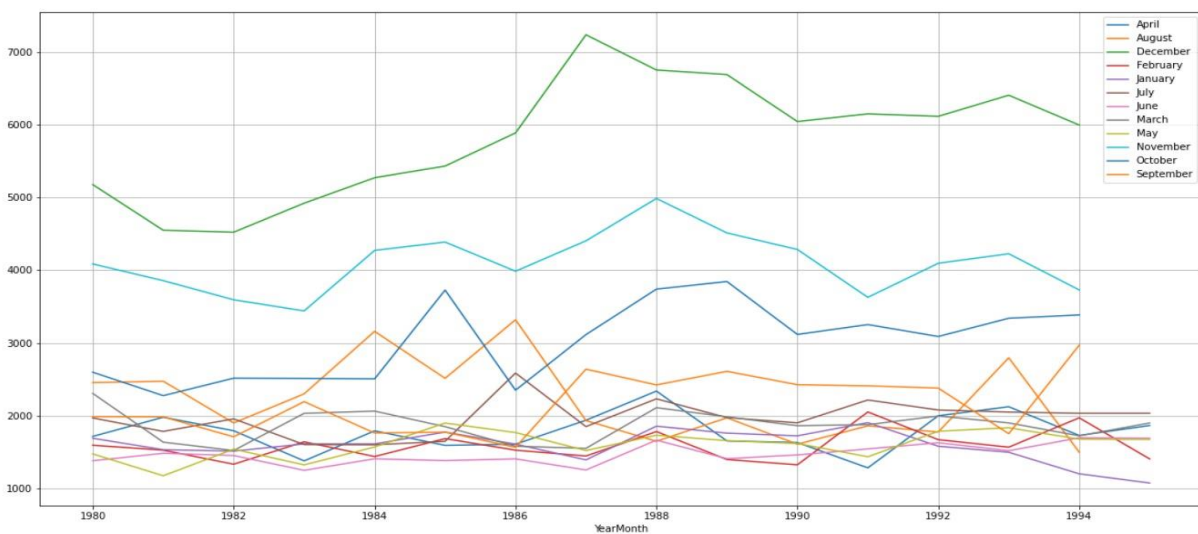
The plot shows that sales are highest in the month of December and lowest in the month of January. Sales are consistent from January to July then from August the sales start to increase. Outliers are present in January, February and July.

Boxplot Weekday wise:



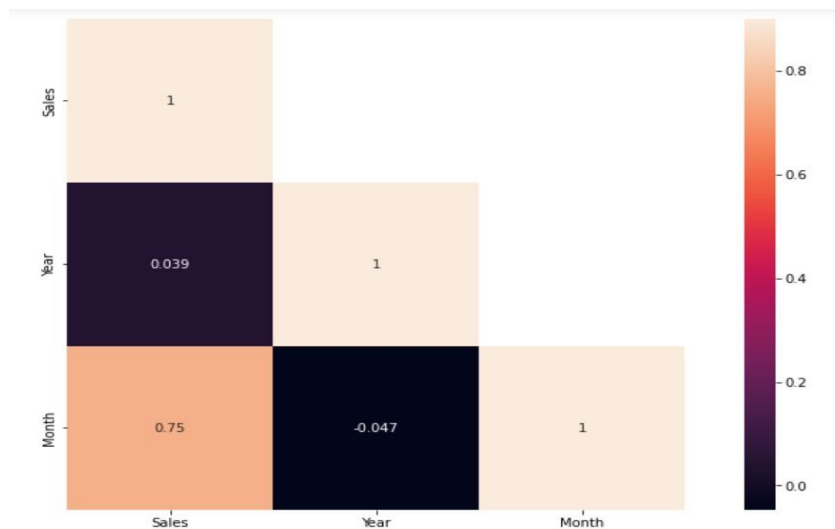
Tuesday has more sales than other days and Wednesday has the lowest sales of the week. Outliers are present on all days which is understandable.

Graph of Monthly Sales over the years:



This plot shows that December has the highest sales over the years and the year 1988 was the year with the highest number of sales.

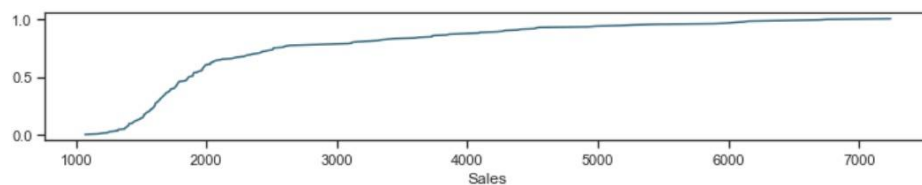
Correlation plot



This heat map shows that there is a low correlation between sales and year. there is a more correlation between month and sales. It indicated seasonal patterns in sales

Plot ECDF: Empirical Cumulative Distribution Function

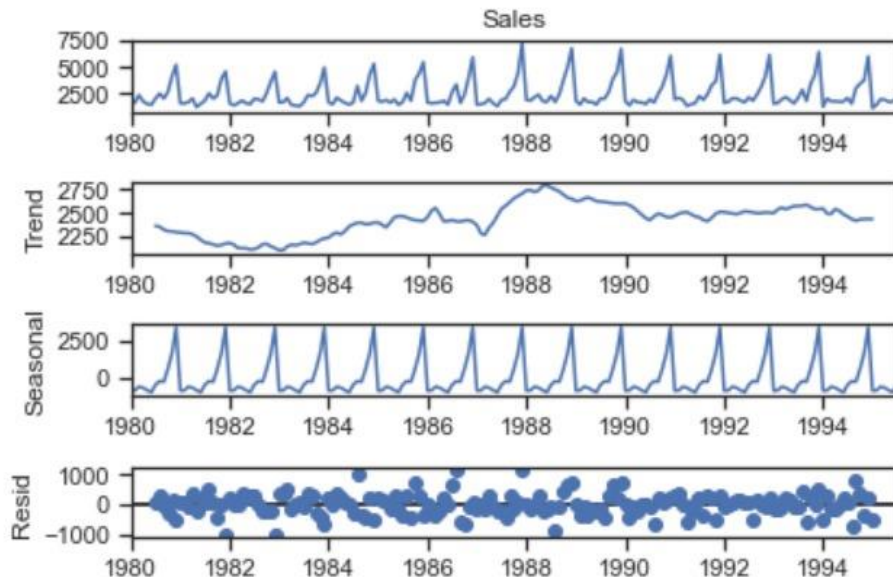
This graph shows the distribution of data.



This plot shows:

- More than 50% of sales have been less than 2000
- Highest values is 7000
- Approx 80% of sales have been less than 3000

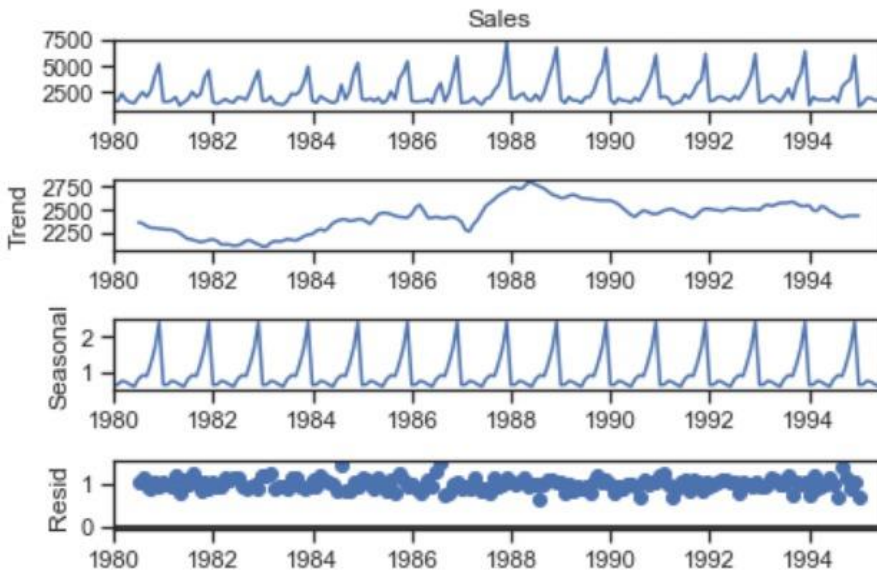
Decomposition –Additive



The plots show:

- Peak year 1988-1989
- It also shows that the trend has declined over the year after 1988-1989.
- Residue is spread and is not in a straight line.
- Both trend and seasonality are present.

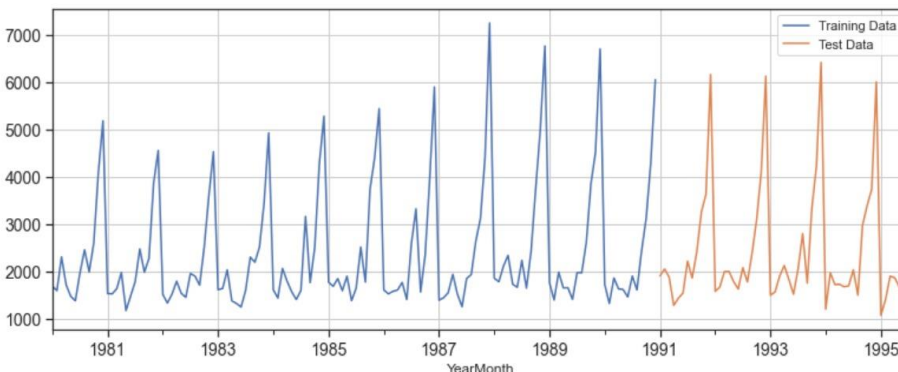
Decomposition-Multiplicative



The plots show

- Peak year 1988-1989
- It also shows that the trend has declined over the year after 1988-1989.
- Residue is spread and is in approx a straight line.
- Both trend and seasonality are present.
- Residue is 0 to 1, while additive is 0 to 1000.
- So multiplicative model is selected owing to a more stable residual plot and lower range of residuals.

3. Split the data into training and test. The test data should start in 1991.



As per the instructions given in the project we have split the data, around 1991. With training data from 1980 to 1990 December. Test data starts from the first month of January 1991 till the end.

Rows and Columns:

train dataset has 132 rows and 3 columns. test dataset has 55 and 3 columns.

Few Rows of datasets:

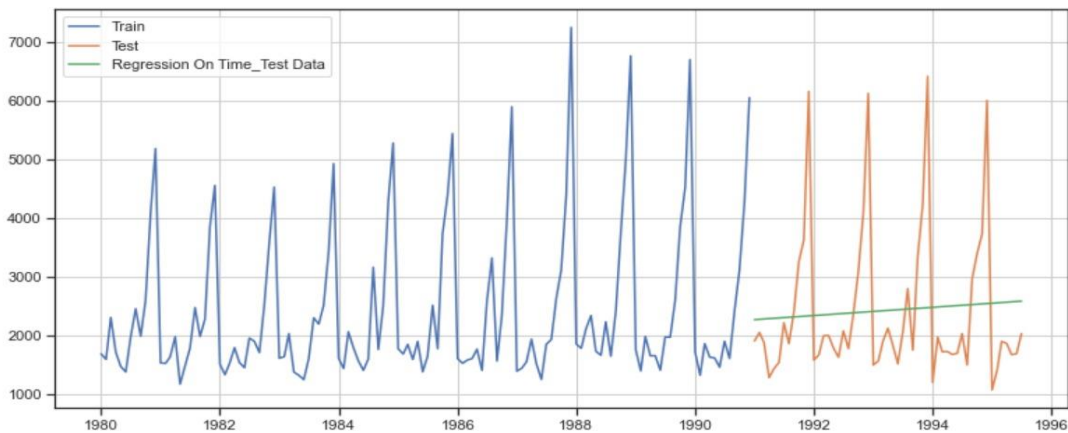
Rows of dataset:				First few rows of Test Data			
First few rows of Training Data							
YearMonth	Sales	Year	Month	YearMonth	Sales	Year	Month
1980-01-01	1686	1980	1	1991-01-01	1902	1991	1
1980-02-01	1591	1980	2	1991-02-01	2049	1991	2
1980-03-01	2304	1980	3	1991-03-01	1874	1991	3
1980-04-01	1712	1980	4	1991-04-01	1279	1991	4
1980-05-01	1471	1980	5	1991-05-01	1432	1991	5
Last few rows of Training Data				Last few rows of Test Data			
YearMonth	Sales	Year	Month	YearMonth	Sales	Year	Month
1990-08-01	1605	1990	8	1995-03-01	1897	1995	3
1990-09-01	2424	1990	9	1995-04-01	1862	1995	4
1990-10-01	3116	1990	10	1995-05-01	1670	1995	5
1990-11-01	4286	1990	11	1995-06-01	1688	1995	6
1990-12-01	6047	1990	12	1995-07-01	2031	1995	7

4. Build all the 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, and simple average models. should also be built on the training data and check the performance on the test data using RMSE.

- Model 1: Linear Regression
- Model 2: Naive Approach
- Model 3: Simple Average
- Model 4: Moving Average(MA)
- Model 5: Simple Exponential Smoothing

- Model 6: Double Exponential Smoothing (Holt's Model)
- Model 7: Triple Exponential Smoothing (Holt - Winter's Model)

Model 1: Linear Regression

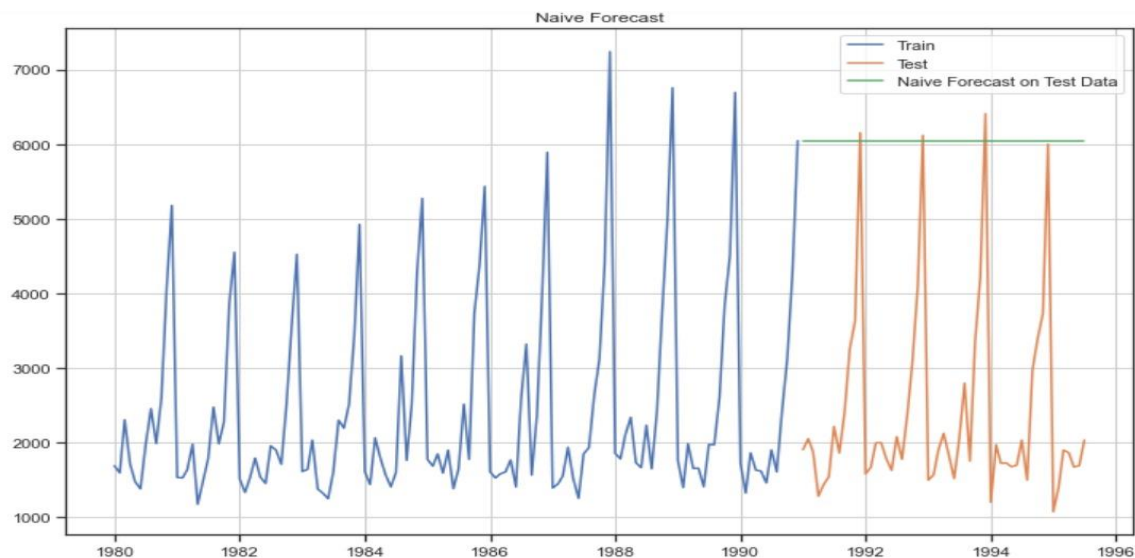


The green line indicates the predictions made by the model, while the orange values are the actual test values. It is clear the predicted values are very far off from the actual values

Model was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

Linear Regression 1275.867052

Model 2: Naive Approach:

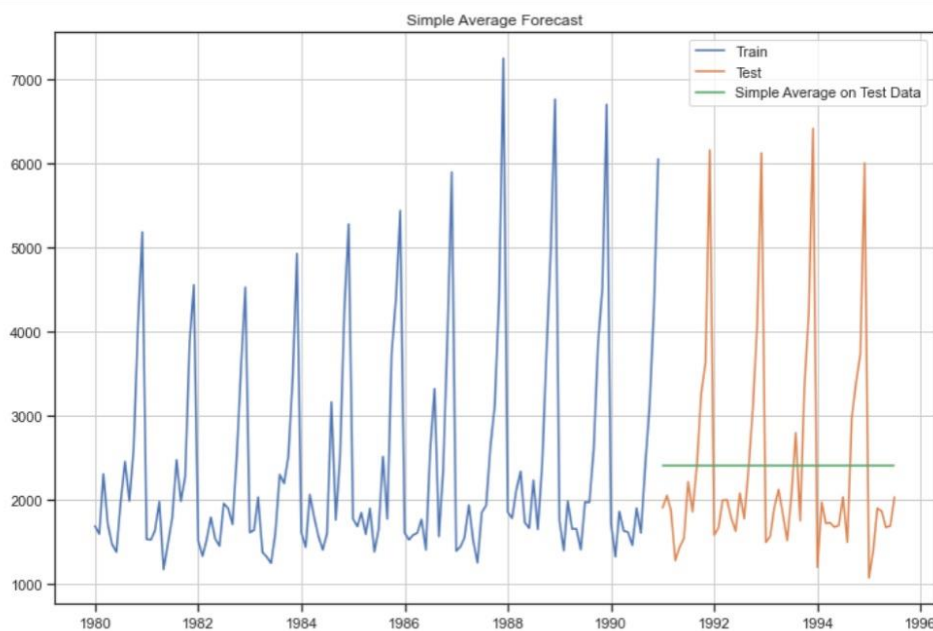


The green line indicates the predictions made by the model, while the orange values are the actual test values. It is clear the predicted values are very far off from the actual values

Model was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

Naive Model 3864.279352

Method 3: Simple Average

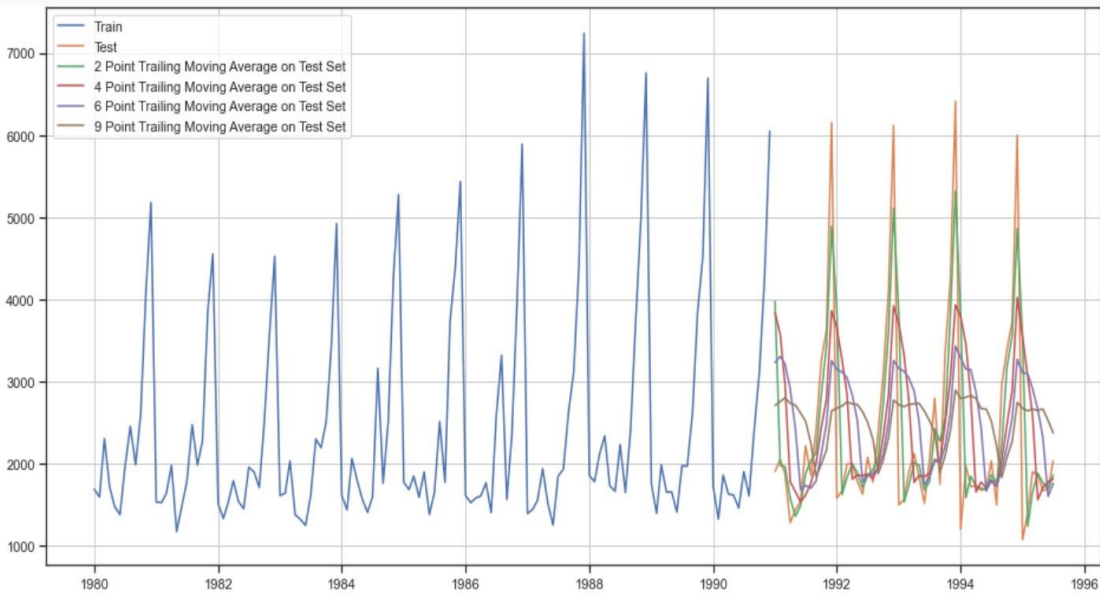


The green line indicates the predictions made by the model, while the orange values are the actual test values. It is clear the predicted values are very far off from the actual values

Model was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

Simple Average Model 1275.081804

Method 4: Moving Average(MA)



Model was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

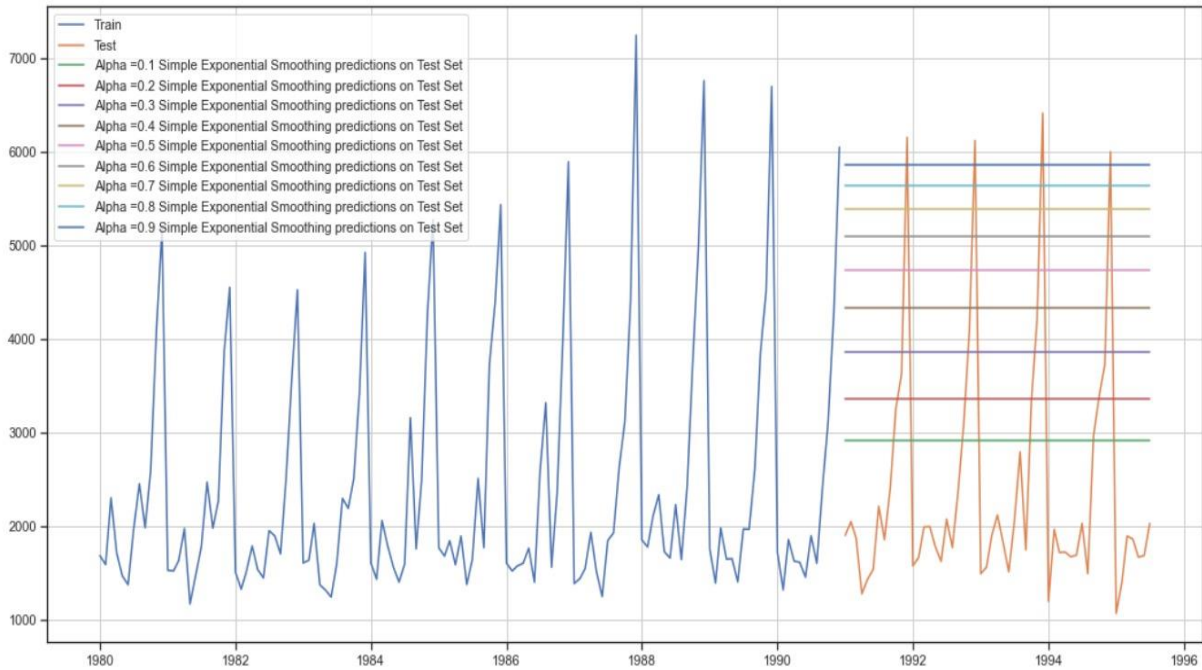
2pointTrailingMovingAverage 813.400684

4pointTrailingMovingAverage 1156.589694

6pointTrailingMovingAverage 1283.927428 9pointTrailingMovingAverage 1346.278315

We have made multiple moving average models with rolling windows varying from 2 to 9. Rolling average is a better method than simple average as it takes into account only the previous n values to make the prediction, where n is the rolling window defined. This takes into account the recent trends and is in general more accurate. The higher the rolling window, the smoother will be its curve, since more values are being taken into account.

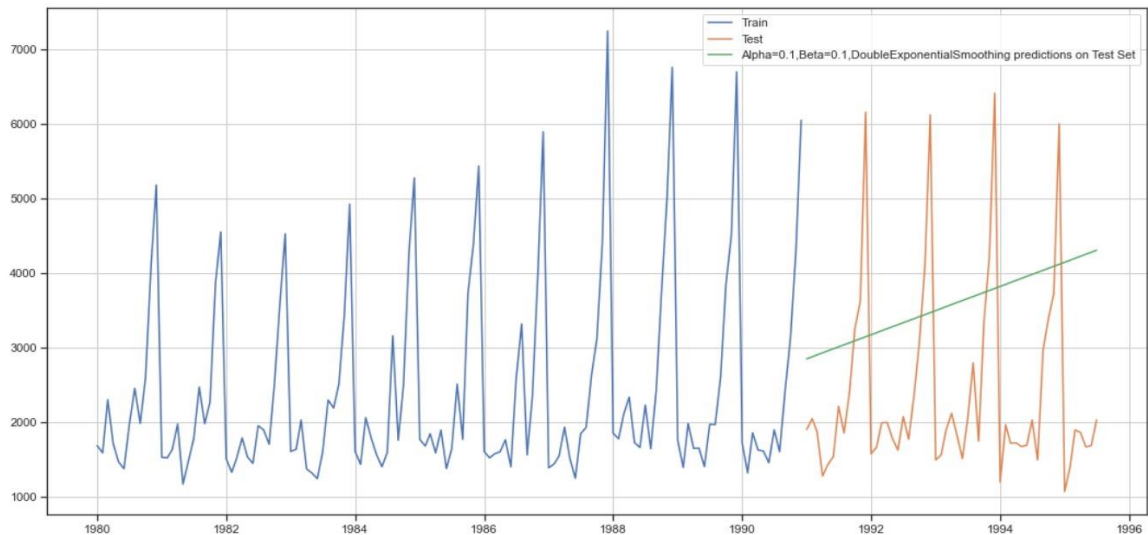
Method 5: Simple Exponential Smoothing



Model was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

	Alpha Values	Train RMSE	Test RMSE
0	0.1	1333.873836	1375.393398
1	0.2	1356.042987	1595.206839
2	0.3	1359.511747	1935.507132
3	0.4	1352.588879	2311.919615
4	0.5	1344.004369	2666.351413
5	0.6	1338.805381	2979.204388
6	0.7	1338.844308	3249.944092
7	0.8	1344.462091	3483.801006
8	0.9	1355.723518	3686.794285

Method 6: Double Exponential Smoothing (Holt's Model)

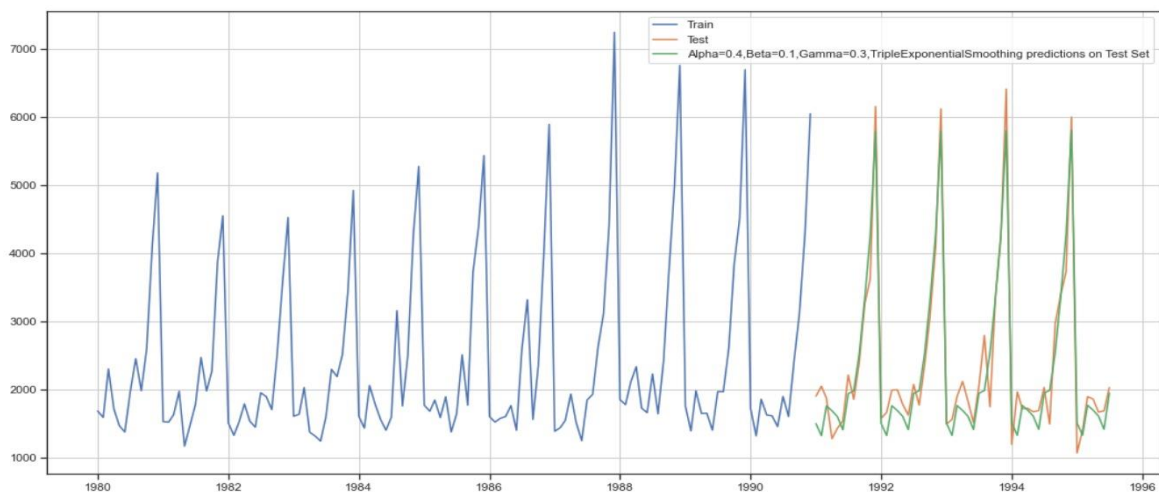


The green line indicates the predictions made by the model, while the orange values are the actual test values. It is clear the predicted values are very far off from the actual values

Model was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

Alpha Value = 0.1, beta value = 0.1, Double Exponential Smoothing = 1778.564670

Method 7: Triple Exponential Smoothing (Holt - Winter's Model)



Output for a best alpha, beta, and gamma values are shown by the green color line in the above plot. The best model had both a multiplicative trends, as well as a seasonality Model, which was evaluated using the RMSE metric. Below is the RMSE calculated for this model.

Alpha=0.4, Beta=0.1, Gamma=0.3, TripleExponentialSmoothing 317.434302

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. Note: Stationarity should be checked at $\alpha = 0.05$.

Check for stationarity of the whole Time Series data.

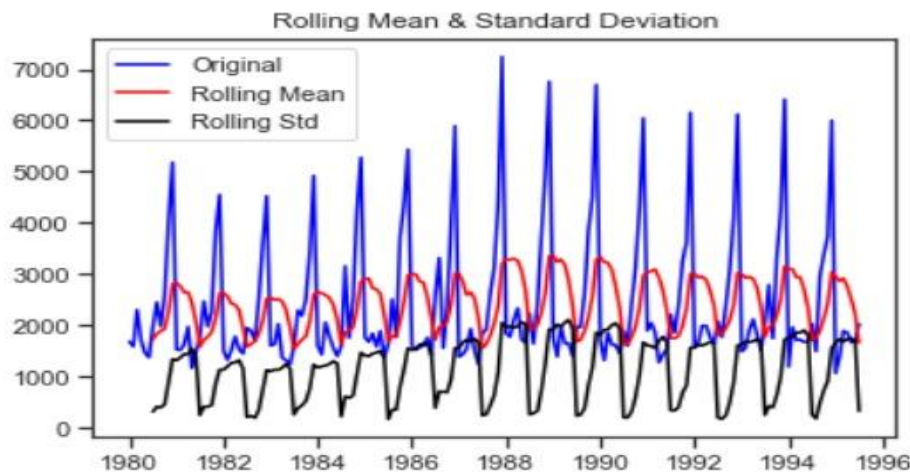
The Augmented Dickey-Fuller test is a unit root test which determines whether there is a unit root and subsequently whether the series is non-stationary.

The hypothesis in a simple form for the ADF test is:

- H_0 : The Time Series has a unit root and is thus non-stationary.
- H_1 : The Time Series does not have a unit root and is thus stationary.

We would want the series to be stationary for building ARIMA models and thus we would want the p-value of this test to be less than the α value.

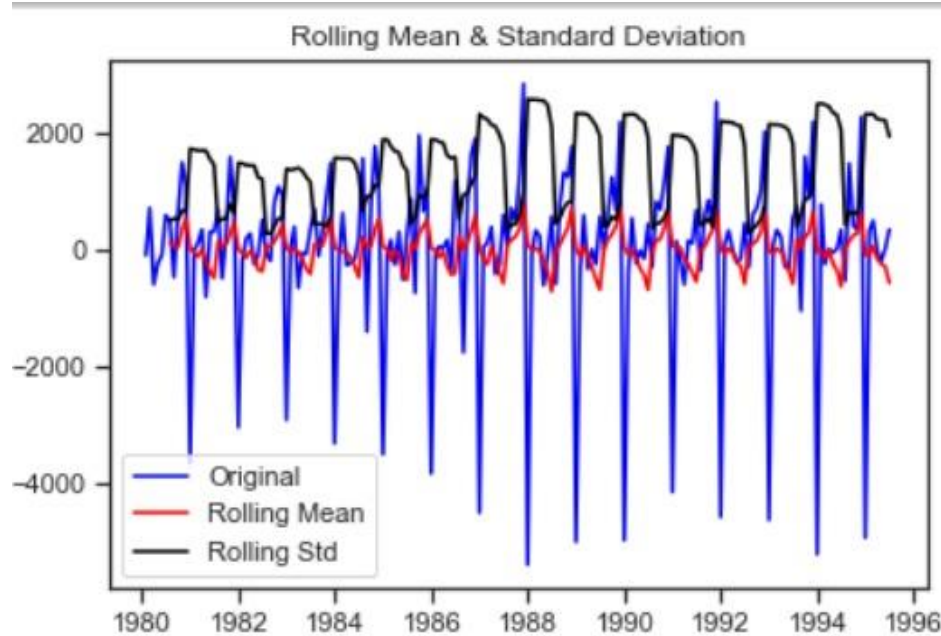
We see that at 5% significant level the Time Series is non-stationary.



Results of Dickey-Fuller Test:

p-value 0.601061

In order to try and make the series stationary we used the differencing approach. We used `.diff()` function on the existing series without any argument, implying the default diff value of 1 and also dropped the NaN values, since differencing of order 1 would generate the first value as NaN which need to be dropped



Results of Dickey-Fuller Test:

p-value 0.000000

Dickey - Fuller test was 0.000, which is obviously less than 0.05. Hence the null hypothesis that the series is not stationary at difference = 1 was rejected, which implied that the series has indeed become stationary after we performed the differencing. Null hypothesis was rejected since the p-value was less than alpha i.e. 0.05. Also the rolling mean plot was a straight line this time around. Also the series looked more or less the same from both the directions, indicating stationarity.

We could now proceed ahead with ARIMA/ SARIMA models, since we had made the series stationary.

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.

AUTO - ARIMA model

We employed a for loop for determining the optimum values of p, d, q , where p is the order of the AR (Auto-Regressive) part of the model, while q is the order of the MA (Moving Average) part of the model. d is the differencing that is required to make the series stationary. p, q values in the range of (0,4) were given to the for loop, while a fixed value of 1 was given for d , since we had already determined d to be 1, while checking for stationarity using the ADF test.

Some parameter combinations for the Model...

Model: (0, 1, 1)

Model: (0, 1, 2)

Model: (0, 1, 3)

Model: (1, 1, 0)

Model: (1, 1, 1)

Model: (1, 1, 2)

Model: (1, 1, 3)

Model: (2, 1, 0)

Model: (2, 1, 1)

Model: (2, 1, 2)

Model: (2, 1, 3)

Model: (3, 1, 0)

Model: (3, 1, 1)

Model: (3, 1, 2)

Model: (3, 1, 3)

Akaike information criterion (AIC) value was evaluated for each of these models and the model with least AIC value was selected.

	param	AIC
10	(2, 1, 2)	2213.509213
15	(3, 1, 3)	2221.458583
14	(3, 1, 2)	2230.768028
11	(2, 1, 3)	2232.885328
9	(2, 1, 1)	2233.777626
3	(0, 1, 3)	2233.994858
2	(0, 1, 2)	2234.408323
6	(1, 1, 2)	2234.5272
13	(3, 1, 1)	2235.49868
7	(1, 1, 3)	2235.607812
5	(1, 1, 1)	2235.755095
12	(3, 1, 0)	2257.723379
8	(2, 1, 0)	2260.365744
1	(0, 1, 1)	2263.060016
4	(1, 1, 0)	2266.608539
0	(0, 1, 0)	2267.663036

the summary report for the ARIMA model with values (p=2,d=1,q=2).

```

=====
SARIMAX Results
=====
Dep. Variable:      Sales      No. Observations:      132
Model:              ARIMA(2, 1, 2)  Log Likelihood          -1101.755
Date:              Fri, 17 Feb 2023  AIC                        2213.509
Time:              23:59:37      BIC                        2227.885
Sample:            01-01-1980      HQIC                       2219.351
                  - 12-01-1990
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       -7.740      0.000        -0.701        -0.418
ma.L1         -1.9917        0.109      -18.215      0.000        -2.206        -1.777
ma.L2          0.9999        0.110        9.108      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 1.8e+28. Standard errors may be unstable.

```

RMSE values are as below:

Auto_ARIMA 1299.978401

AUTO- SARIMA Model

A similar for loop like AUTO_ARIMA with below values was employed, resulting in the models shown below.

```
p = q = range(0, 4) d= range(0,2) D = range(0,2) pdq = list(itertools.product(p, d, q)) model_pdq
= [(x[0], x[1], x[2], 12) for x in list(itertools.product(p, D, q))]
```

Examples of some parameter combinations for Model...

Model: (0, 1, 1)(0, 0, 1, 12)

Model: (0, 1, 2)(0, 0, 2, 12)

Model: (0, 1, 3)(0, 0, 3, 12)

Model: (1, 1, 0)(1, 0, 0, 12)

Model: (1, 1, 1)(1, 0, 1, 12)

Model: (1, 1, 2)(1, 0, 2, 12)

Model: (1, 1, 3)(1, 0, 3, 12)

Model: (2, 1, 0)(2, 0, 0, 12)

Model: (2, 1, 1)(2, 0, 1, 12)

Model: (2, 1, 2)(2, 0, 2, 12)

Model: (2, 1, 3)(2, 0, 3, 12)

Model: (3, 1, 0)(3, 0, 0, 12)

Model: (3, 1, 1)(3, 0, 1, 12)

Model: (3, 1, 2)(3, 0, 2, 12)

Model: (3, 1, 3)(3, 0, 3, 12)

Akaike information criterion (AIC) value was evaluated for each of these models and the model with least AIC value was selected. Here only the top 5 models are shown.

	param	seasonal	AIC
91	(1, 1, 1)	(2, 0, 3, 12)	84.365834
252	(3, 1, 3)	(3, 0, 0, 12)	1387.497086
220	(3, 1, 1)	(3, 0, 0, 12)	1387.788331
237	(3, 1, 2)	(3, 0, 1, 12)	1388.602592
221	(3, 1, 1)	(3, 0, 1, 12)	1388.681484

the summary report for the best SARIMA model with values (2,1,2)(2,0,2,12)

SARIMAX Results

Dep. Variable: y No. Observations: 132

Model: SARIMAX(1, 1, 2)x(1, 0, 2, 12) Log Likelihood -770.792

Date: Sat, 18 Feb 2023 AIC 1555.584

Time: 01:27:20 BIC 1574.095

Sample: 0 HQIC 1563.083

- 132

Covariance Type: opg

	coef	std err	z	P> z	[0.025	0.975]
ar.L1	-0.6282	0.255	-2.464	0.014	-1.128	-0.128
ma.L1	-0.1040	0.225	-0.463	0.644	-0.545	0.337
ma.L2	-0.7277	0.154	-4.736	0.000	-1.029	-0.427
ar.S.L12	1.0439	0.014	72.835	0.000	1.016	1.072
ma.S.L12	-0.5550	0.098	-5.663	0.000	-0.747	-0.363
ma.S.L24	-0.1354	0.120	-1.133	0.257	-0.370	0.099
sigma2	1.506e+05	2.03e+04	7.401	0.000	1.11e+05	1.9e+05

Ljung-Box (L1) (Q): 0.04 Jarque-Bera (JB): 11.72

Prob(Q): 0.84 Prob(JB): 0.00

Heteroskedasticity (H): 1.47 Skew: 0.36

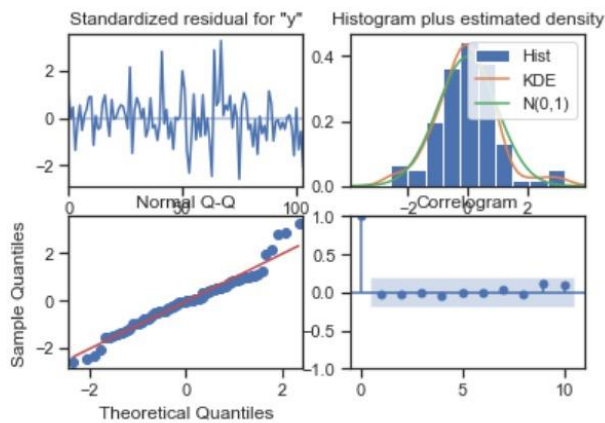
Prob(H) (two-sided): 0.26 Kurtosis: 4.48

Warnings:

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

We also plotted the graphs for the residual to determine if any further information can be extracted or all the usable information has already been extracted. Below were the plots for the best auto SARIMA model.

Plot 22: SARIMA plot



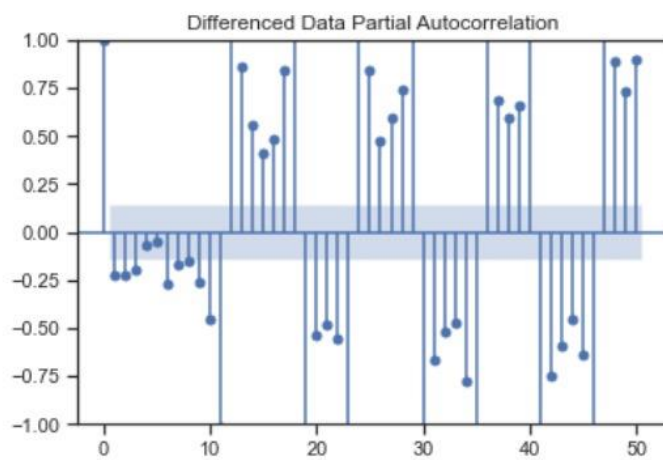
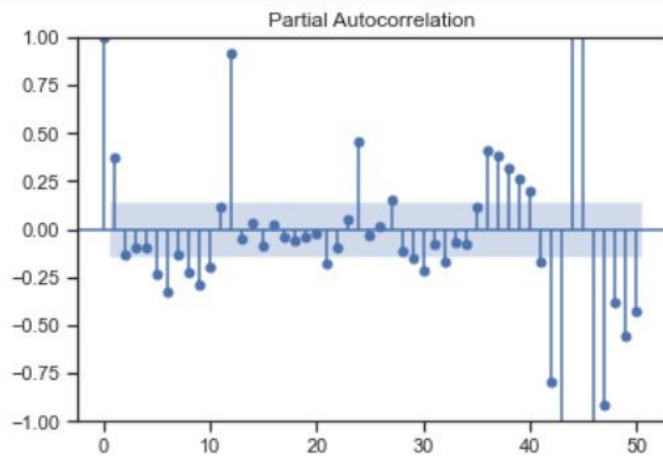
RSME of Model:

528.6069474180102

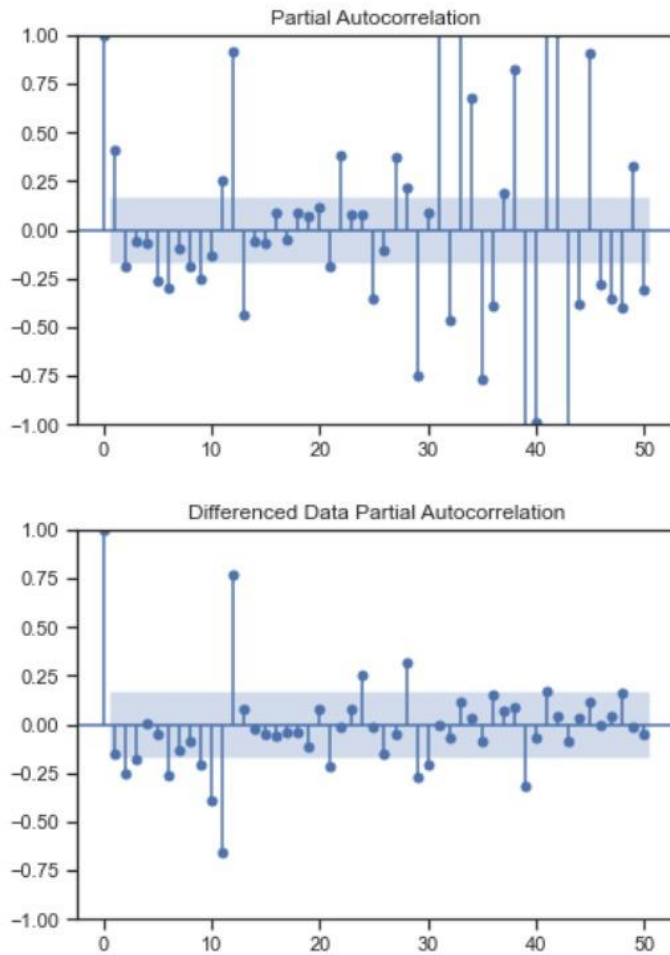
7) Build ARIMA/SARIMA models based on the cut-off points of ACF and PACF on the training data and evaluate this model on the test data using RMSE.

Manual- ARIMA Model

PFB the ACF plot on data



training data with `diff(1)`:



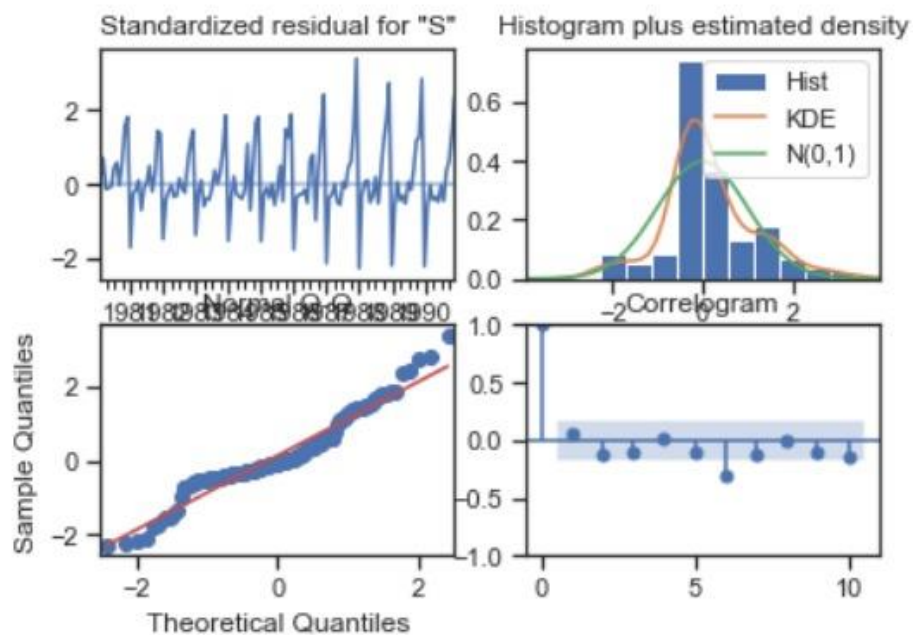
Looking at ACF plot we can see a sharp decay after lag 1 for original as well as differenced data. Hence we select the q value to be 1. i.e. $q=1$.

Looking at PACF plot we can again see significant bars till lag 1 for differenced series which is stationary in nature, post 1 the decay is large enough. Hence we choose p value to be 1. i.e. $p=1$. d values will be 1, since we had seen earlier that the series is stationary with lag 1. Hence the values selected for manual ARIMA:- $p=1$, $d=1$, $q=1$ summary from this manual ARIMA model.

```

=====
SARIMAX Results
=====
Dep. Variable:      Sales      No. Observations:      132
Model:              ARIMA(1, 1, 1)  Log Likelihood         -1114.878
Date:              Sat, 18 Feb 2023  AIC                        2235.755
Time:              01:27:35      BIC                     2244.381
Sample:            01-01-1980     HQIC                    2239.260
                  - 12-01-1990
Covariance Type:    opg
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
ar.L1          0.4494      0.043      10.366      0.000      0.364      0.534
ma.L1         -0.9996      0.102      -9.811      0.000     -1.199     -0.800
sigma2        1.401e+06    7.57e-08    1.85e+13    0.000    1.4e+06    1.4e+06
=====
Ljung-Box (L1) (Q):      0.50      Jarque-Bera (JB):      10.42
Prob(Q):                 0.48      Prob(JB):              0.01
Heteroskedasticity (H):  2.64      Skew:                  0.46
Prob(H) (two-sided):     0.00      Kurtosis:              4.03
=====

```



Model Evaluation: RSME

1319.9367298218867

Manual SARIMA Model

SARIMAX(1, 1, 1)x(1, 1, 1, 12)

Below is the summary of the manual SARIMA model

```

=====
SARIMAX Results
=====
Dep. Variable:          y          No. Observations:      132
Model:                SARIMAX(1, 1, 1)x(1, 1, 1, 12)  Log Likelihood      -882.088
Date:                 Sat, 18 Feb 2023              AIC              1774.175
Time:                 01:27:47                     BIC              1788.071
Sample:               0                          HQIC              1779.818
Covariance Type:      opg
=====

```

	coef	std err	z	P> z	[0.025	0.975]
ar.L1	0.1957	0.104	1.878	0.060	-0.009	0.400
ma.L1	-0.9404	0.053	-17.897	0.000	-1.043	-0.837
ar.S.L12	0.0711	0.242	0.294	0.769	-0.404	0.546
ma.S.L12	-0.5035	0.221	-2.277	0.023	-0.937	-0.070
sigma2	1.51e+05	1.33e+04	11.371	0.000	1.25e+05	1.77e+05

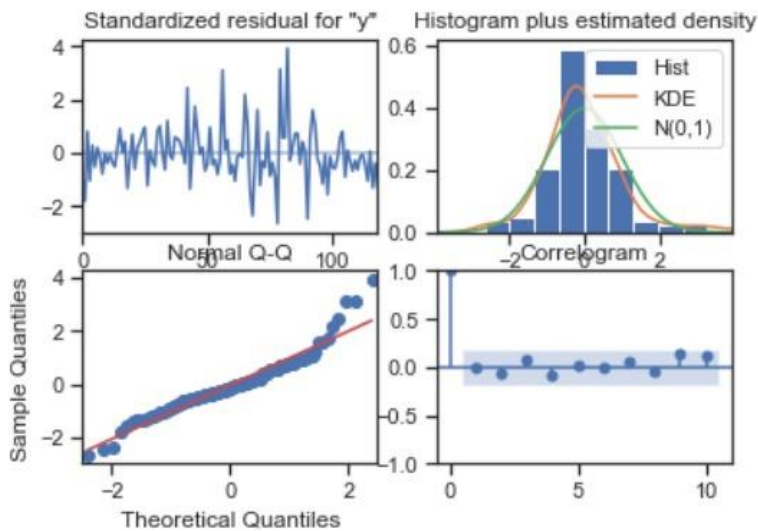
```

=====
Ljung-Box (L1) (Q):      0.01  Jarque-Bera (JB):      45.66
Prob(Q):                 0.93  Prob(JB):             0.00
Heteroskedasticity (H):   2.61  Skew:              0.82
Prob(H) (two-sided):      0.00  Kurtosis:          5.56
=====

```

Warnings:

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



Model Evaluation: RSME

359.612454

8. Build a table (create a data frame) with all the models built along with their corresponding parameters and the respective RMSE values on the test data.

	Test RMSE
Linear Regression	1275.867052
Naive Model	3864.279352
Simple Average Model	1275.081804
2pointTrailingMovingAverage	813.400684
4pointTrailingMovingAverage	1156.589694
6pointTrailingMovingAverage	1283.927428
9pointTrailingMovingAverage	1346.278315
Alpha=0.1, SimpleExponentialSmoothing	1375.393398
Alpha Value = 0.1, beta value = 0.1, DoubleExponentialSmoothing	1778.564670
Alpha=0.08621, Beta=1.3722, Gamma=0.4763, TrippleExponentialSmoothing_Auto_Fit	1316.034674
Alpha=0.4, Beta=0.1, Gamma=0.2, TripleExponentialSmoothing	317.434302
Auto_ARIMA	1299.978401
(1,1,1),(2,0,3,12),Auto_SARIMA	528.606947
ARIMA(3,1,3)	1319.936730
(1,1,1)(1,1,1,12),Manual_SARIMA	359.612454

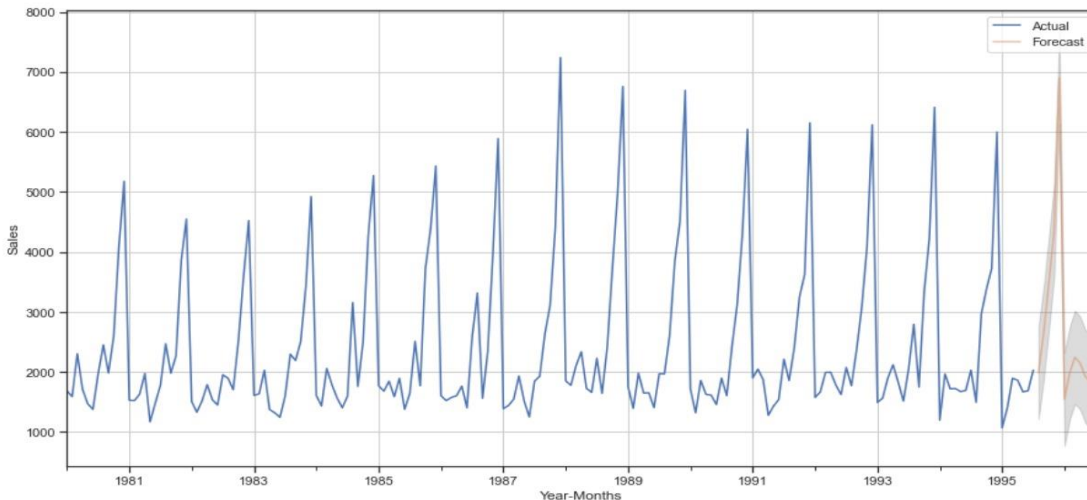
9 Based on the model-building exercise, build the most optimum model(s) on the complete data and predict 12 months into the future with appropriate confidence intervals/bands.

Based on the above comparison of all the various models that we had built, we can conclude that the triple exponential smoothing or the Holts-Winter model is giving us the lowest RMSE, hence it would be the most optimum model

sales predictions made by this best optimum model.

Sales_Predictions	
1995-08-01	1988.782193
1995-09-01	2652.762887
1995-10-01	3483.872246
1995-11-01	4354.989747
1995-12-01	6900.103171
1996-01-01	1546.800546
1996-02-01	1981.361768
1996-03-01	2245.459724
1996-04-01	2151.066942
1996-05-01	1929.355815
1996-06-01	1830.619260
1996-07-01	2272.156151

the sales prediction on the graph along with the confidence intervals. PFB the graph.



Predictions, 1 year into the future are shown in orange color, while the confidence interval has been shown in grey color.

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 sales for Sparkling wine for the company are predicted to be at least the same as last year, if not more, with peak sales for next year potentially higher than this year.
- Sparkling wine has been a consistently popular wine among customers with only a very marginal decline in sales, despite reaching its peak popularity in the late 1980s.
- Seasonality has a significant impact on the sales of Sparkling wine, with sales being slow in the first half of the year and picking up from August to December.
- It is recommended for the company to run campaigns in the first half of the year when sales are slow, particularly in the months of March to July.
- Combining promotions where Sparkling wine is paired with a less popular wine such as "Rose wine" under a special offer may encourage customers to try the underperforming wine, which could potentially boost its sales and benefit the company.

