

BUSINESS REPORT ON
TIME SERIES FORECASTING
ROSE SALES DATASET
DATA SCIENCE AND BUSINESS ANALYTICS

By: Sanjam Preet Singh Bhullar

May 2023

Table of Contents

Problem Statement.....	6
1. Read the data as an appropriate time series data and plot the data.....	7
Null values.....	8
Plotting the time series.....	9
2. Perform appropriate Exploratory Data Analysis to understand the data and also perform decomposition.	10
Yearly sales.....	10
Quarterly sales.....	11
Monthly sales.....	11
Day-wise sales.....	13
Empirical Cumulative Distribution Function (ECDF)	13
Average sales.....	14
Sales percentage change.....	14
Additive decomposition of time Series.....	15
Multiplicative decomposition of time series.....	16
3. Split the data into training and test. The test data should start in 1991.....	17
Train set.....	17
Test set.....	17
Plotting the train-test graph.....	18
4. Build all 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.....	19
Linear regression.....	19
Naïve model.....	20
Simple average.....	21

Moving average.....	22
Simple Exponential Smoothing	24
Double Exponential Smoothing.....	26
Triple Exponential Smoothing.....	28
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 $\alpha = 0.05$	31
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.....	33
ARIMA model.....	33
SARIMA model.....	35
7. Build a table (create a data frame) with all models built along with their corresponding parameters and RMSE values on the test data.....	38
8. 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.....	39
9. Comment on the model thus built and report your findings and suggest the measures that the company should be taking for future sales. Please explain and summarise the various steps performed in this project. There should be proper business interpretation and actionable insights present.....	41
Exploratory Data Analysis summary.....	41
Model-building summary.....	41
Recommendations.....	41
Appendix.....	43

List of Figures

Figure 1: Time series plot.....	9
Figure 2: Yearly sale boxplot.....	10
Figure 3: Yearly sale line plot.....	10
Figure 4: Quarterly sale boxplot.....	11
Figure 5: Monthly sale boxplot.....	11
Figure 6: Month plot.....	12
Figure 7: Monthly line plot.....	12
Figure 8: Day-wise sales boxplot.....	13
Figure 9: ECDF plot.....	13
Figure 10: Average sales plot.....	14
Figure 11: Sales percentage change plot.....	14
Figure 12: Additive decomposition.....	15
Figure 13: Multiplicative decomposition.....	16
Figure 14: Train-test time series plot.....	18
Figure 15: Regression on test set.....	19
Figure 16: Naïve forecast on test set.....	20
Figure 17: Simple average on test set.....	21
Figure 18: Moving average plot on entire dataset.....	22
Figure 19: Moving average forecast on test set.....	23
Figure 20: SES forecast ($\alpha = 0.987$)	24
Figure 21: SES forecast ($\alpha = 0.3$)	25
Figure 22: DES forecast ($\alpha = 0, \beta = 0$)	26
Figure 23: DES forecast ($\alpha = 0.1, \beta = 0.1$)	27
Figure 24: TES forecast ($\alpha = 0.06, \beta = 0.051, \gamma = 0$)	28
Figure 25: TES forecast ($\alpha = 0.088, \beta = 0.0001, \gamma = 0.0023$).....	29

Figure 26: TES forecast ($\alpha = 0.1$, $\beta = 1$, $\gamma = 0.2$).....	30
Figure 27: Non-stationary time series plot.....	31
Figure 28: Stationary time series plot ($d = 1$).....	32
Figure 29: ACF plot.....	35
Figure 30: Diagnostic plots.....	37
Figure 31: Sales forecast for 12 months.....	39
Figure 32: Sales forecast with confidence interval.....	40

List of tables

Table 1: First 5 rows of dataset.....	7
Table 2: Last 5 rows of dataset.....	7
Table 3: Null values in dataset.....	8
Table 4: Null values imputed.....	8
Table 5: First 5 rows of train set.....	17
Table 6: Last 5 rows of train set.....	17
Table 7: First 5 rows of test set.....	17
Table 8: Last 5 rows of test set.....	17
Table 9: Moving average data sample.....	22
Table 10: RMSE for different α values.....	25
Table 11: RMSE for different α and β values.....	27
Table 12: RMSE for different α , β and γ values.....	30
Table 13: ARIMA parameters.....	33
Table 14: ARIMA summary.....	34
Table 15: SARIMA parameters.....	36
Table 16: SARIMA summary.....	36
Table 17: Comparison of all models.....	38

PROBLEM STATEMENT

For this particular assignment, the data of Rose wine sales in the 20th century is to be analysed. As an analyst in the ABC Estate Wines, you are tasked to analyse and forecast wine sales in the 20th century.

1. Read the data as an appropriate time series data and plot the data.

Rose	
Time Stamp	
1980-01-31	112.0
1980-02-29	118.0
1980-03-31	129.0
1980-04-30	99.0
1980-05-31	116.0

Table 1: First 5 rows of dataset

Rose	
Time Stamp	
1995-03-31	45.0
1995-04-30	52.0
1995-05-31	28.0
1995-06-30	40.0
1995-07-31	62.0

Table 2: Last 5 rows of dataset

The Rose dataset shows the wine sales at the end of each month.

We have the data starting from January 1980 to July 1995 – a duration of 14.5 years.

```
DatetimeIndex: 187 entries, 1980-01-31 to 1995-07-31
Data columns (total 1 columns):
#   Column  Non-Null Count  Dtype
---  ---
0    Rose    185 non-null      float64
dtypes: float64(1)
```

There is **one column** in the dataset, as has been seen in tables 1 and 2.

The dataset has 187 entries, but there are **185 records for Rose wine sales**.

This means there are **two null values** in the dataset.

Null values

Rose	
Time Stamp	
1994-07-31	NaN
1994-08-31	NaN

Table 3: Null values in dataset

There are no records for two months. In a time series data, null values cannot be dropped because doing so will break the order of the data records. Null values will have to be imputed by an appropriate method.

Rose	
Time Stamp	
1994-05-31	44.000000
1994-06-30	45.000000
1994-07-31	45.386725
1994-08-31	44.635276
1994-09-30	46.000000
1994-10-31	51.000000

Table 4: Null values imputed

The two null values have been imputed in such a manner that they are close to the preceding and succeeding sale records.

Plotting the time series

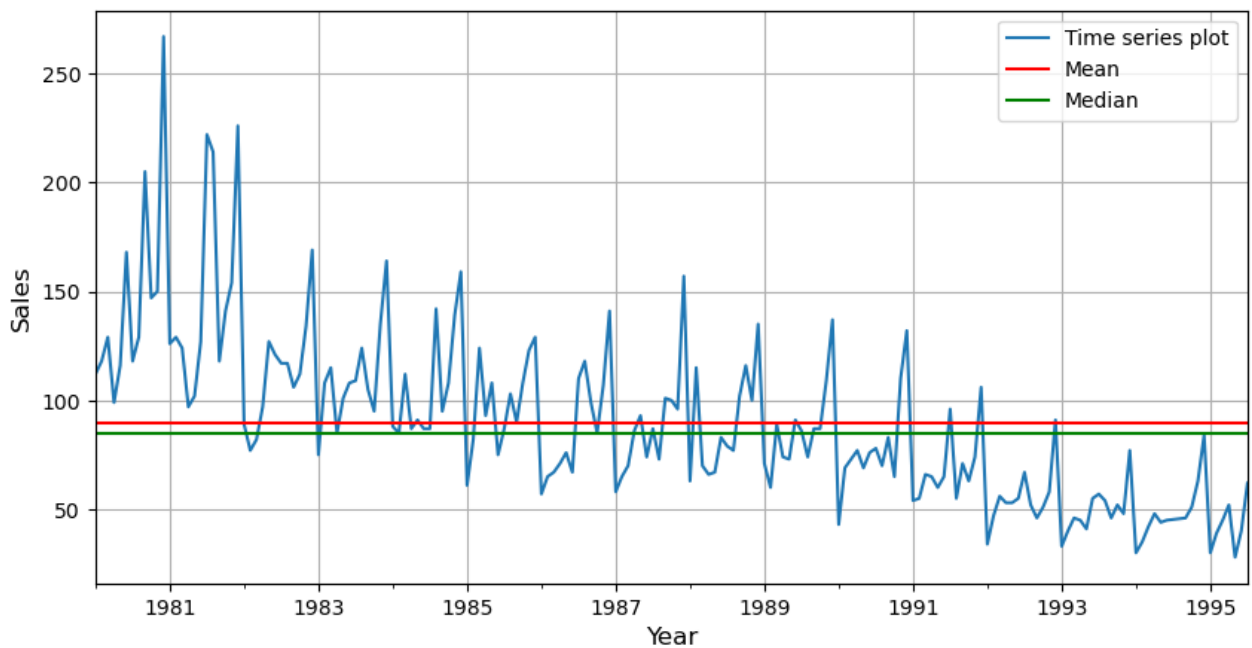


Figure 1: Time series plot

It can be seen from the plot that the Rose wine sales has a **decreasing trend over time**. The sales were high in the first part of the 1980s and, as the years rolled by, the sales fell.

The time series dataset has seasonality. It was more pronounced in the 1980s than in the 1990s.

The mean sales is more than the median sales.

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

Yearly sales

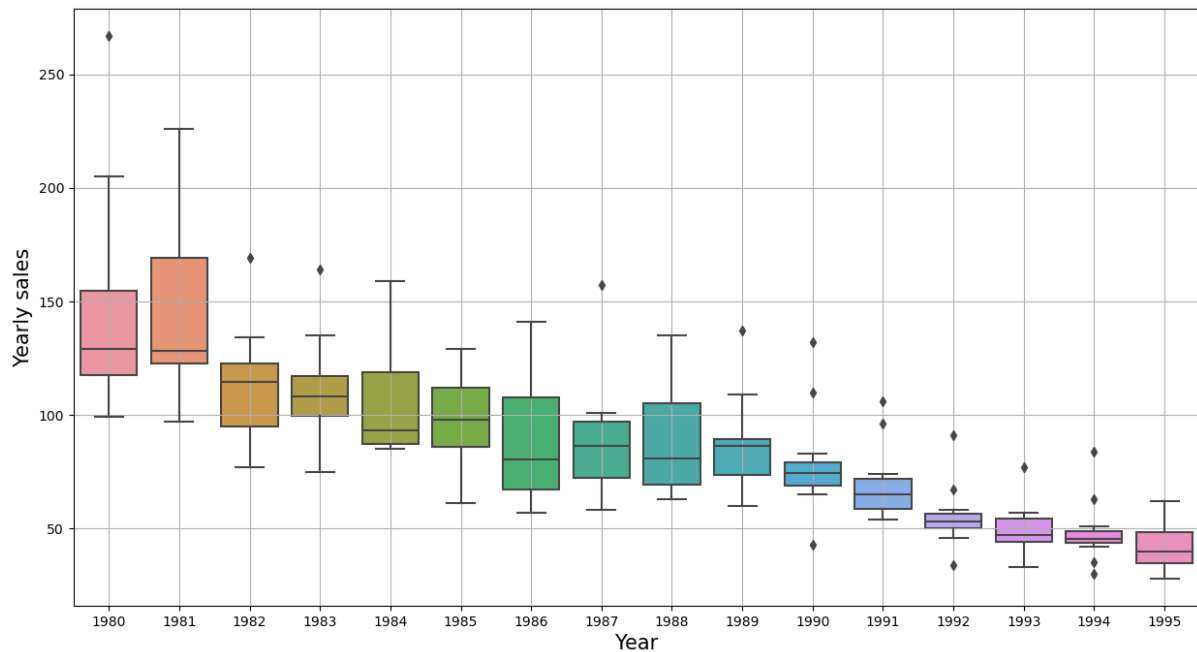


Figure 2: Yearly sale boxplot

The median yearly sales fall over time. The sales were high in the first half of the decade of 1980. In the 1990s, the fall in sales was drastic.

This plot shows the mean sales of Rose wine across years. The sales were the highest in 1981 and decreased in the subsequent years.

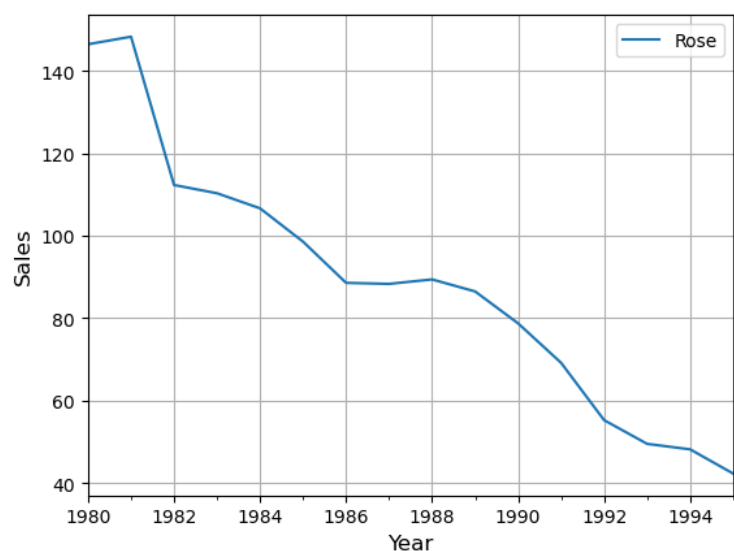


Figure 3: Yearly sale line plot

Quarterly sales

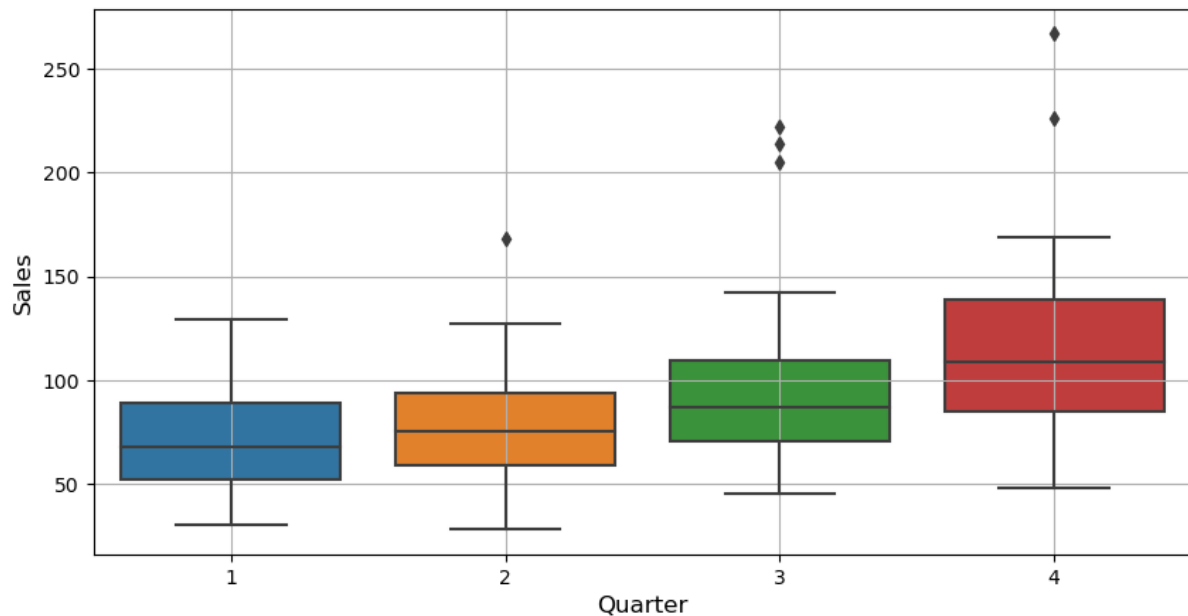


Figure 4: Quarterly sale boxplot

The quarterly median sales increase. It is the highest in quarter four and the lowest in quarter one for any given year.

Monthly sales

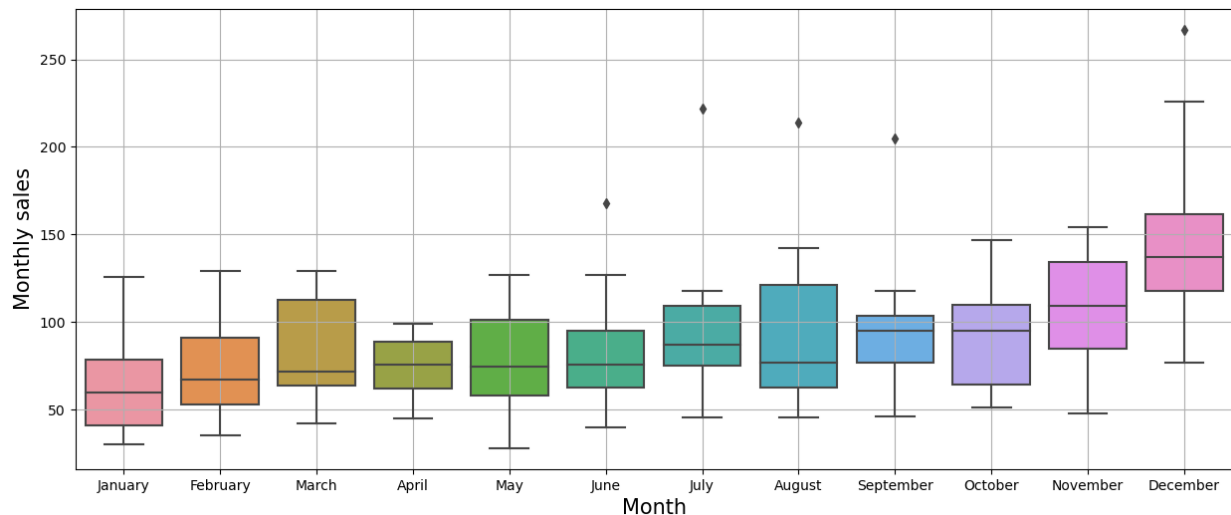


Figure 5: Monthly sale boxplot

The sales remain low at the start of the year, but pick up pace in the subsequent months. The last two months record high sales, with the highest being in December.

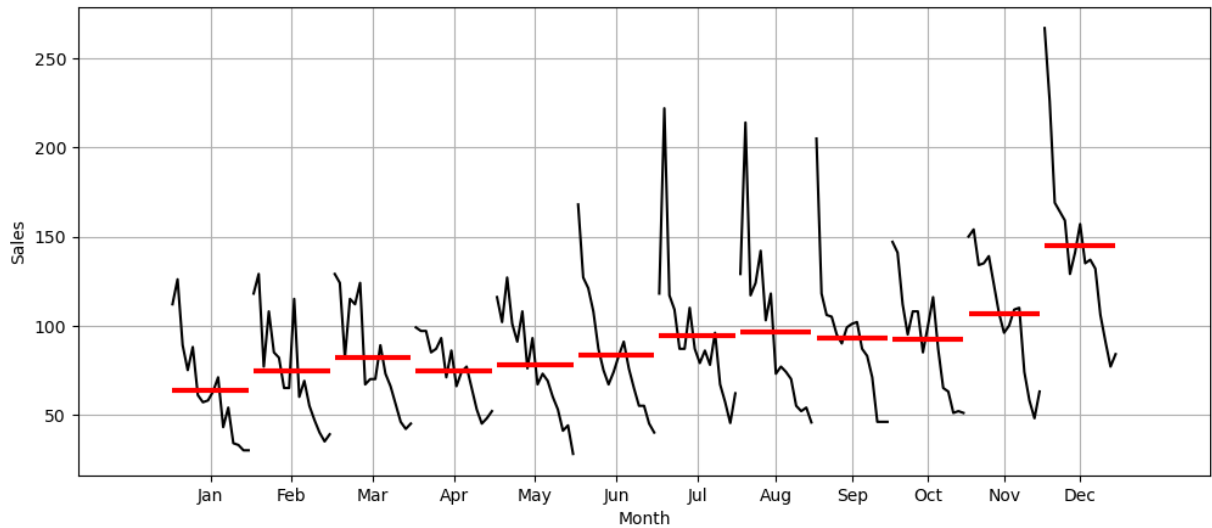


Figure 6: Month plot

This graph reiterates the observation drawn from Figure 5. **Sales are the highest in December.** The red lines denote median sales, which remain low in the initial months of a year and pick up pace in the last two months of the year.

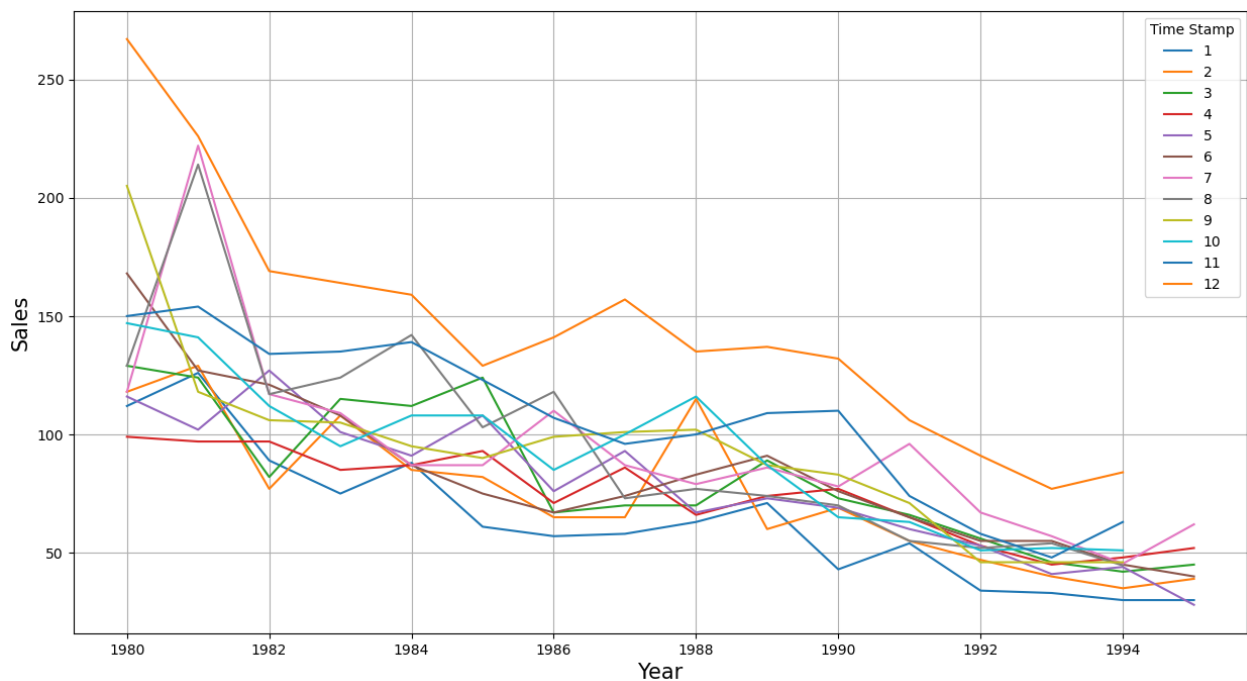


Figure 7: Monthly line plot

The line plot shows month-wise sales across years. Sales in December is the highest. Since sales has a decreasing trend, sales in December fall over time.

Day-wise sales

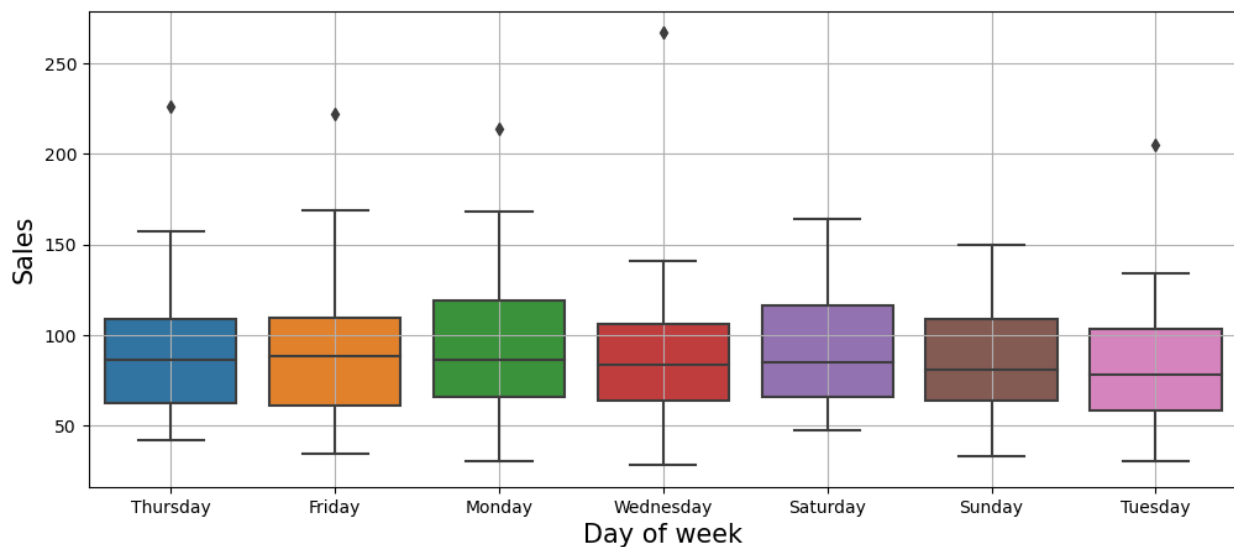


Figure 8: Day-wise sales boxplot

Median sales on all seven days of a week is more or less the same. Sales are the **lowest on Tuesday**. Surprisingly, sales on week days are more than that of Sunday's. There are a few outliers.

Empirical Cumulative Distribution Function (ECDF)

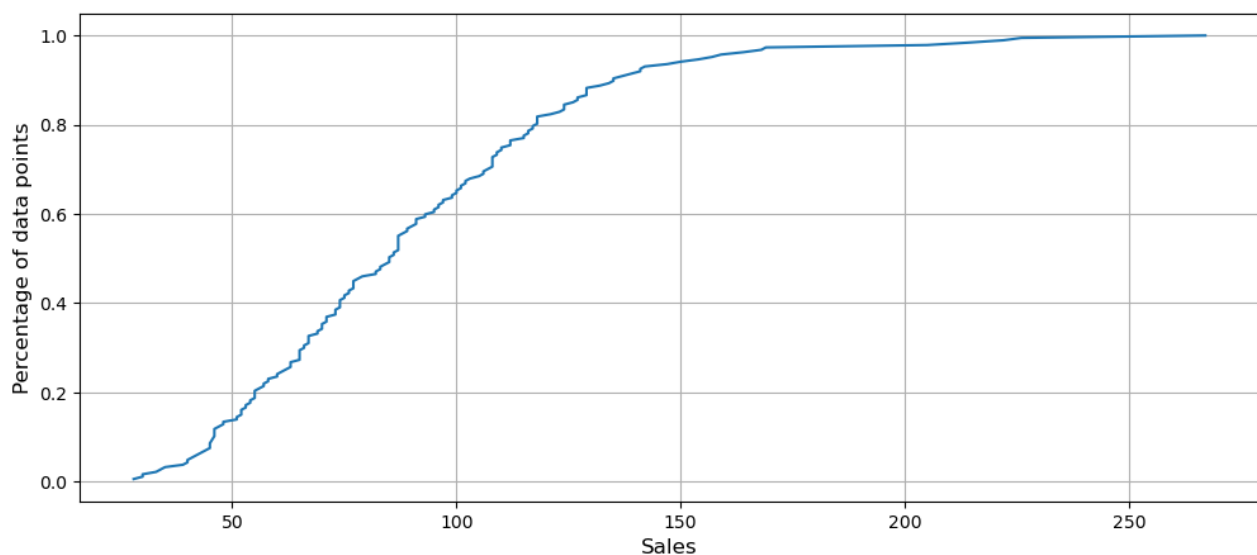


Figure 9: ECDF plot

The ECDF plot shows that 60 per cent of the data points have sale value up to 100 and over 80% of the data points have sale values up to 150.

Average sales

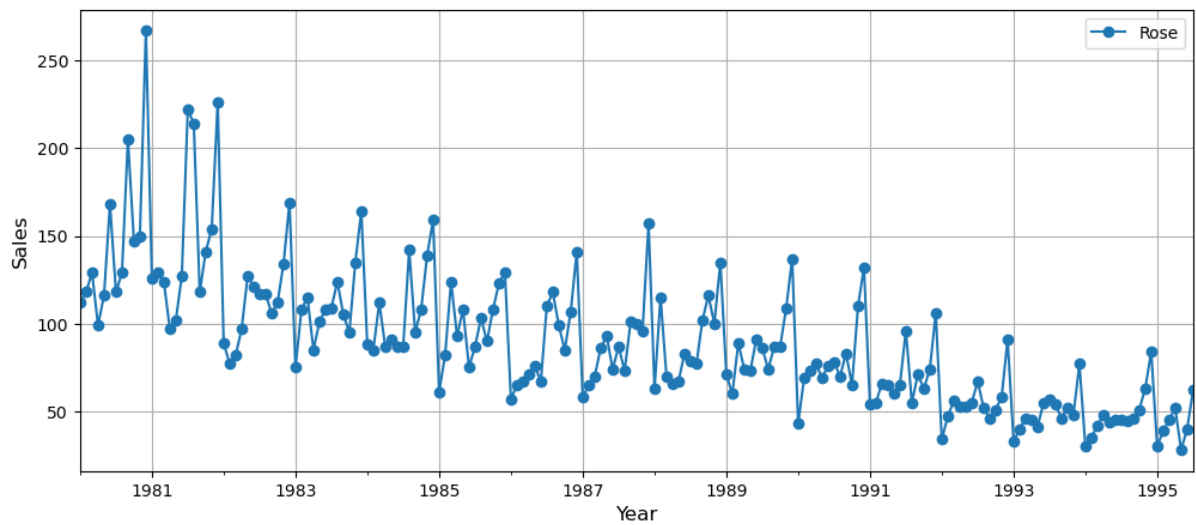


Figure 10: Average sales plot

The average sales plot reiterates the point that the time series dataset has a decreasing trend. The average sales decrease as the years roll by.

Sales percentage change

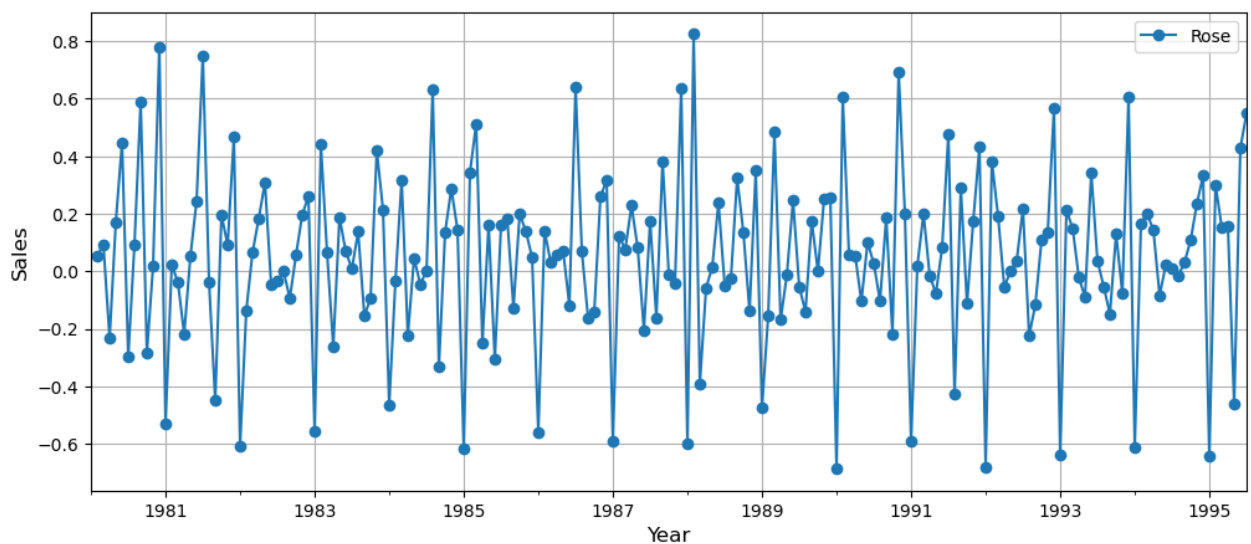


Figure 11: Sales percentage change plot

The graph shows month-on-month percentage change in sales. When percentage change in sales is taken into account, the trend goes missing from the dataset.

Additive decomposition of time series

In the additive model, the time series is the sum of three components – trend, seasonality and residual. Trend and seasonality are systemic components, while residual is an irregular component. These two components are interpretable and can be estimated.

Trend shows the long-term movement of the time series, while **seasonality denotes the intra-year fluctuations** that are repeated over the entire length of the time series.

Residual is the “noise” element – something that happened in the past and is not expected to happen in the future.

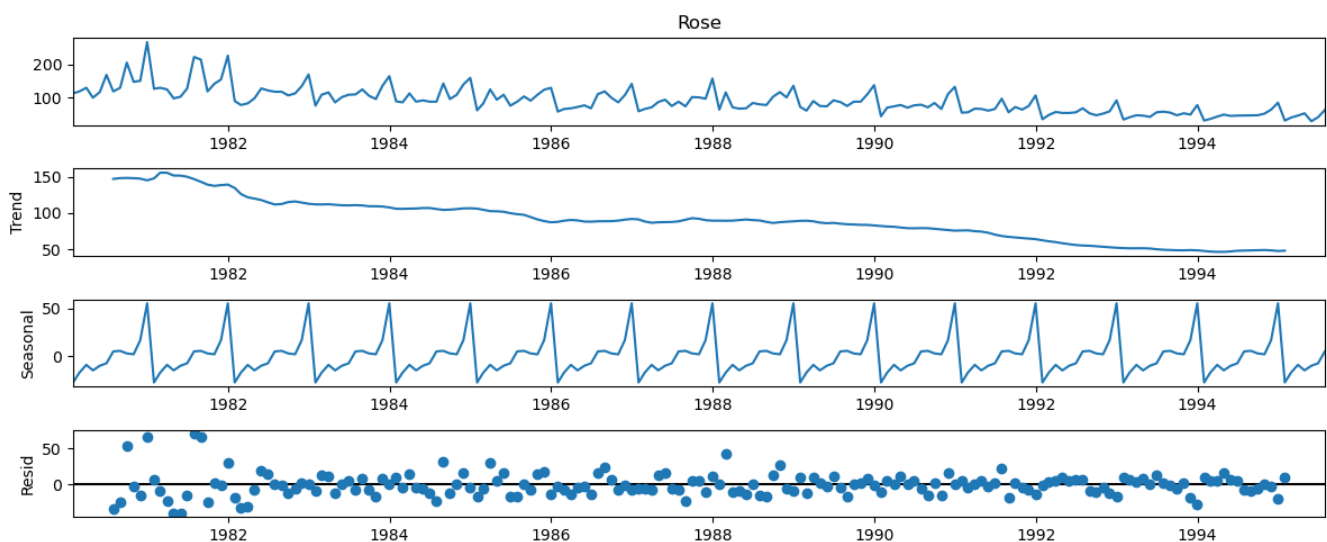


Figure 12: Additive decomposition

The first graph shows the entire time series, as was seen in Figure 1.

The second plot shows the trend over time. The time series dataset has a decreasing trend.

The third plot shows the seasonality component of the time series. The same sales pattern is repeated each year.

The “noise” element, which can be seen in the fourth graph, is the random component which cannot be explained through systematic components. Most of the residuals are centred on a single point.

Multiplicative decomposition of time series

In the multiplicative model, the time series is the multiplication of components – trend, seasonality and residual. It explains the non-linear change over time. Unlike the additive model, it tells the change in percentage terms.

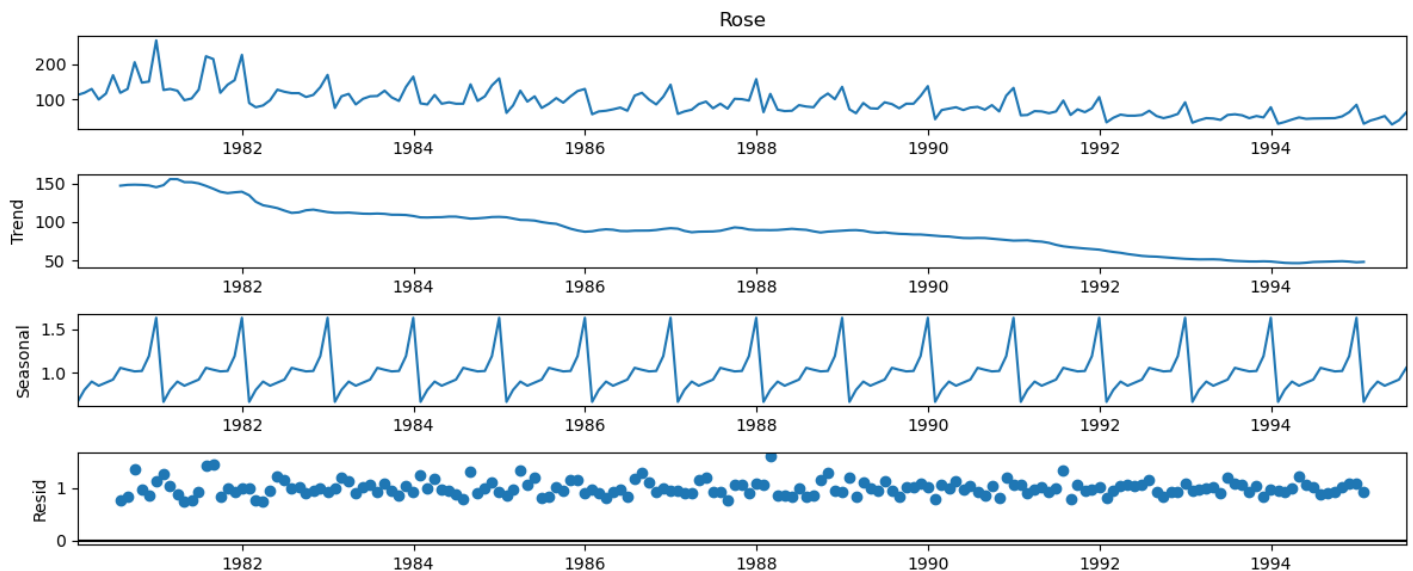


Figure 13: Multiplicative decomposition

The first three plots are the same, with the only difference being that seasonality is plotted against percentage change.

The residual points don't follow a pattern. Almost all values are centred on a single point

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

Unlike classification on regression technique, the data cannot be split randomly. The test set (unseen data) has to be the most recent one because of the ordered nature of data.

The data for Rose wine sales is split in such a manner that the train data should have records till 1990 and the test data starts from 1991.

Train set

Rose	
Time Stamp	
1980-01-31	112.0
1980-02-29	118.0
1980-03-31	129.0
1980-04-30	99.0
1980-05-31	116.0

Table 5: First 5 rows of train set

Rose	
Time Stamp	
1990-08-31	70.0
1990-09-30	83.0
1990-10-31	65.0
1990-11-30	110.0
1990-12-31	132.0

Table 6: Last 5 rows of train set

The train set contains **132 rows** and **1 column**.

Test set

Rose	
Time Stamp	
1991-01-31	54.0
1991-02-28	55.0
1991-03-31	66.0
1991-04-30	65.0
1991-05-31	60.0

Table 7: First 5 rows of test set

Rose	
Time Stamp	
1995-03-31	45.0
1995-04-30	52.0
1995-05-31	28.0
1995-06-30	40.0
1995-07-31	62.0

Table 8: Last 5 rows of test set

The test set has **55 rows** and **1 column**.

Plotting the train-test graph

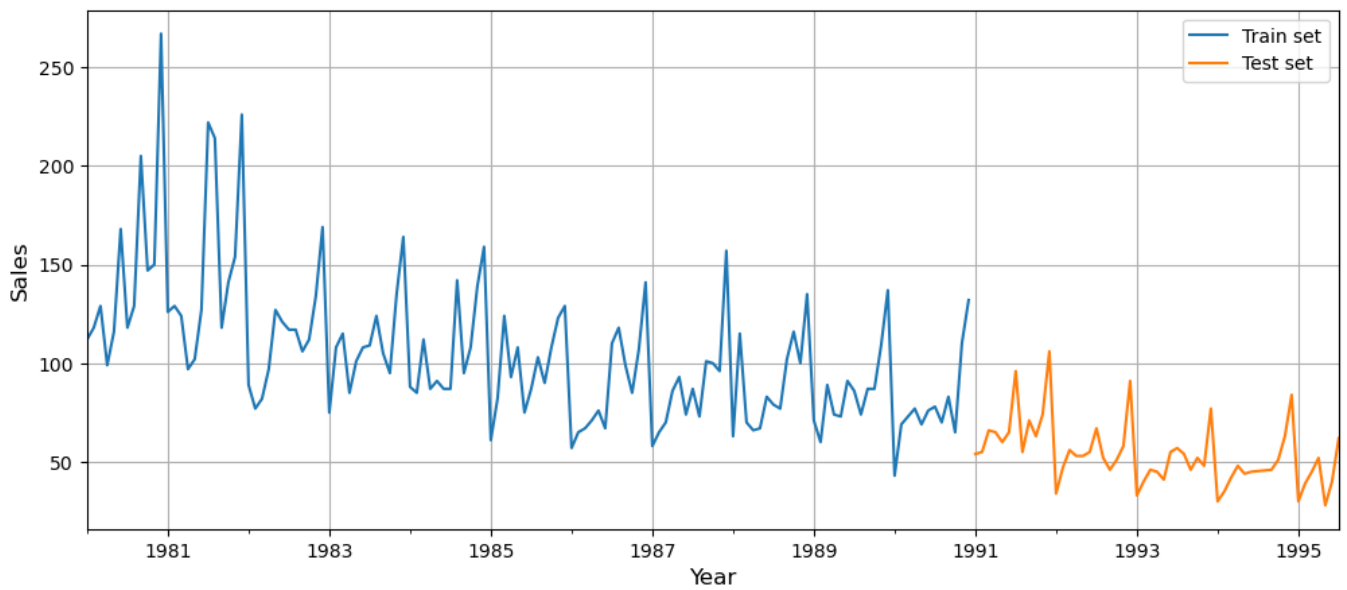


Figure 14: Train-test time series plot

The most recent sales records, starting from 1991, have been taken as the test set.

4. Build all 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.

Linear regression

Before building the regression model, we regress the wines sales variable against the order of the occurrence.

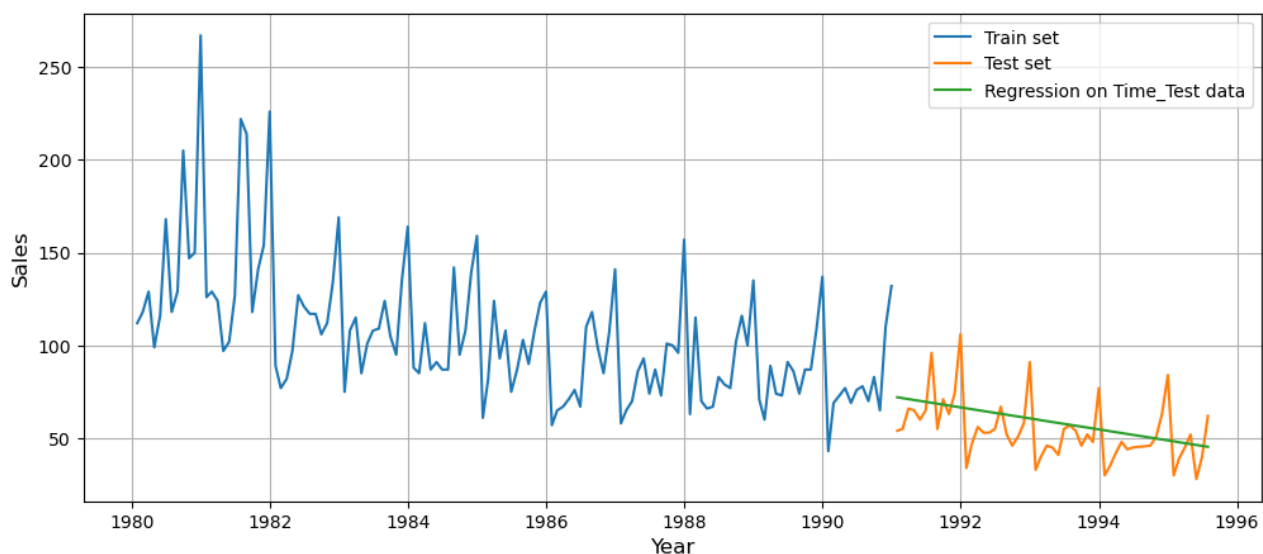


Figure 15: Regression on test set

The regression model captures the decreasing trend, but fails to factor in the seasonality component of the time series.

Root mean squared error (RMSE), which is a metric to measure the forecast accuracy, on the test set is 15.269.

Naïve model

In the Naïve approach, **the last observed value is taken to forecast for future**. In other words, the last value from the train test is used to forecast for the future, i.e., the entire test period.

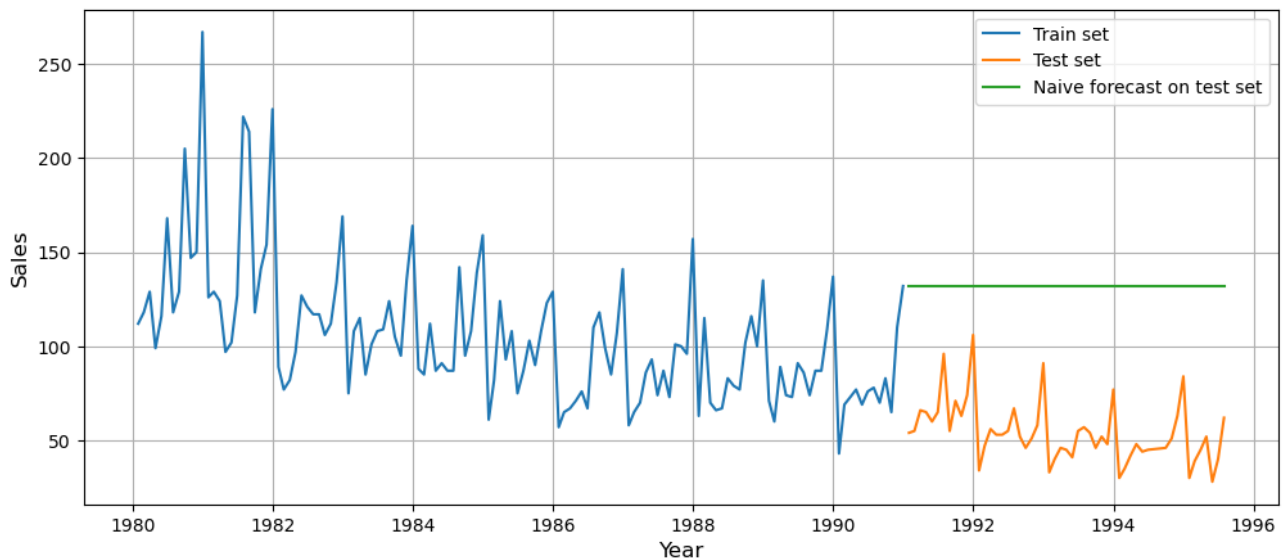


Figure 16: Naïve forecast on test set

The **forecast is constant**. It neither captures the trend nor the seasonality.

RMSE on test set is 79.719.

The RMSE value for the Naïve model is more than that of the regression model.

Simple average

For the simple average method, we will forecast by using the average of the training dataset values.

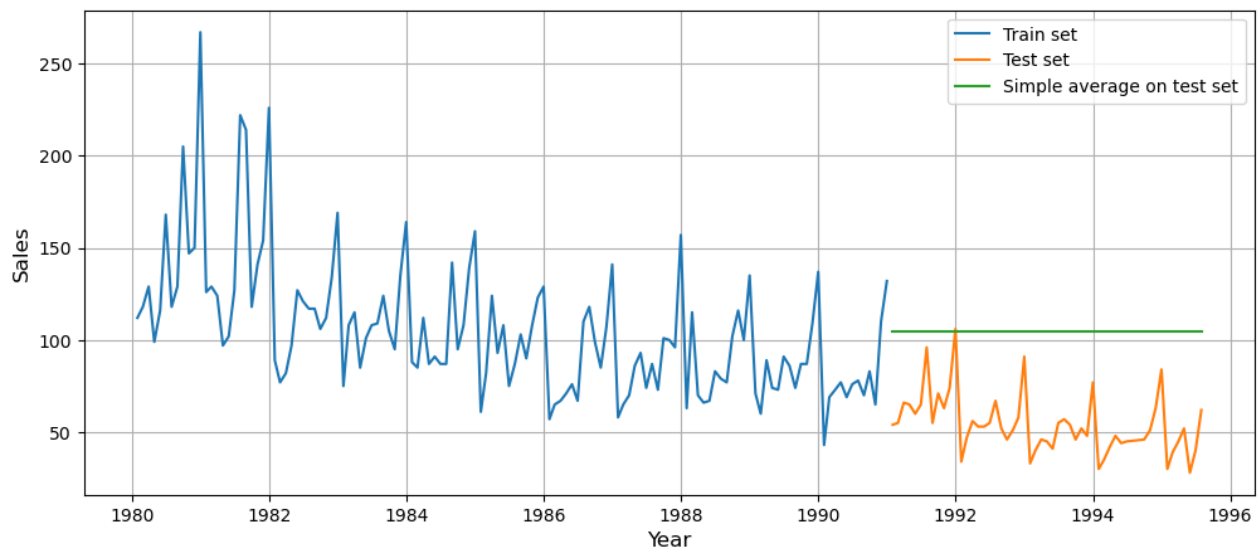


Figure 17: Simple average on test set

It can be observed that the forecast on the test set is constant. The simple average method neither factors in trend nor seasonality.

RMSE on test set is 53.4605.

Moving average

In the moving average model, we are going to calculate rolling means (or moving averages) for different intervals. The best interval can be determined by the maximum accuracy (or the minimum error).

For the Rose wine sales, we are going to take intervals of 2, 4, 6 and 9, and then average over the entire data.

	Rose	Trailing_2	Trailing_4	Trailing_6	Trailing_9
Time Stamp					
1995-03-31	45.0	42.0	49.50	52.000000	49.888889
1995-04-30	52.0	48.5	41.50	52.166667	50.629630
1995-05-31	28.0	40.0	41.00	46.333333	48.666667
1995-06-30	40.0	34.0	41.25	39.000000	48.000000
1995-07-31	62.0	51.0	45.50	44.333333	49.222222

Table 9: Moving average data sample

The table shows the last five rows of the dataset with rolling means taken at time intervals of 2, 4, 6 and 9.

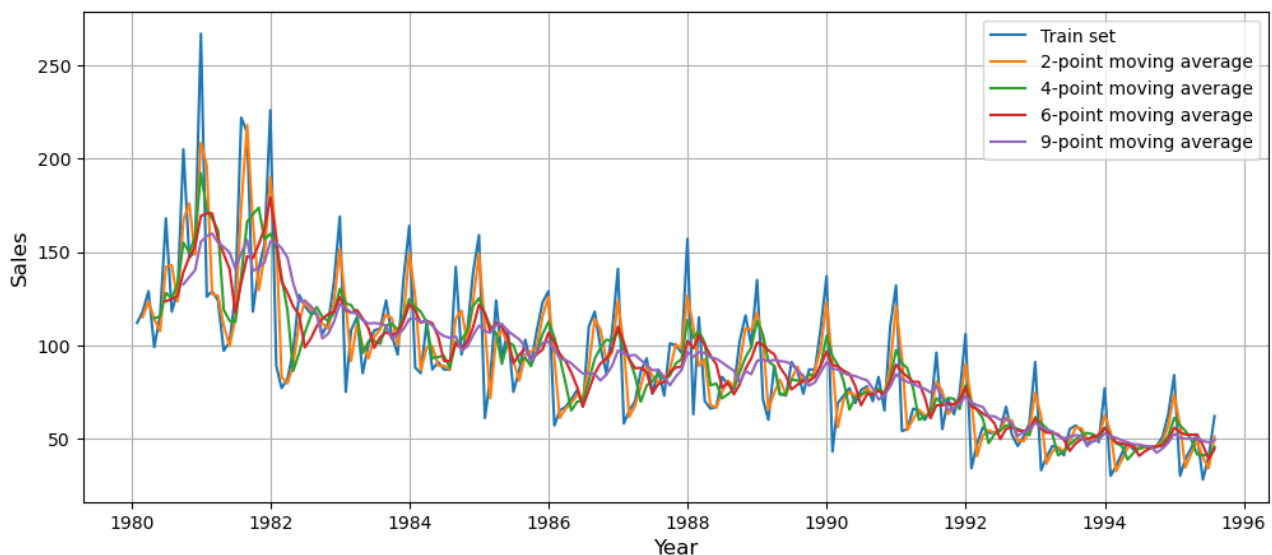


Figure 18: Moving average plot on entire dataset

The figure shows moving average on the entire time series.

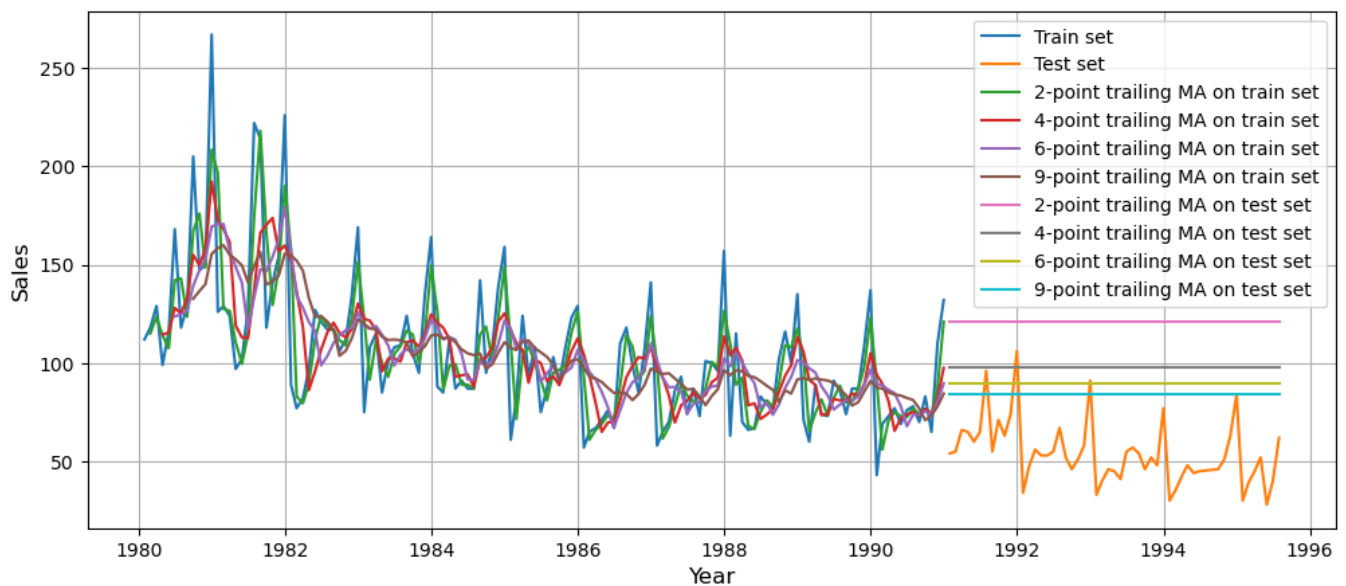


Figure 19: Moving average forecast on test set

The moving average forecast at different intervals gives a **constant prediction** for the future.

The 2-point rolling average forecasts higher sales.

The 9-point rolling average predicts the lowest sales.

RMSE for 2-point rolling moving average is 68.970.

RMSE for 4-point rolling moving average is 46.404.

RMSE for 6-point rolling moving average is 39.126.

RMSE for 9-point rolling moving average is 34.411.

RMSE is the highest for interval 2 and the lowest for interval 9.

Simple Exponential Smoothing

The Simple Exponential Smoothing (SES) method is suitable for forecasting data with no clear trend or seasonal pattern.

Parameter alpha (α) is called the smoothing constant. It corresponds to level, which is the local mean. Its value lies between 0 and 1. Since the model uses only one smoothing constant, it is called Simple Exponential Smoothing.

If α is closer to 1, forecasts follow the actual observations more closely.

If α is closer to 0, forecasts are farther from the actual observations and the prediction line is smooth.

For the Rose wine sales, **α is taken to be 0.987**. This value has been automatically generated by the model.

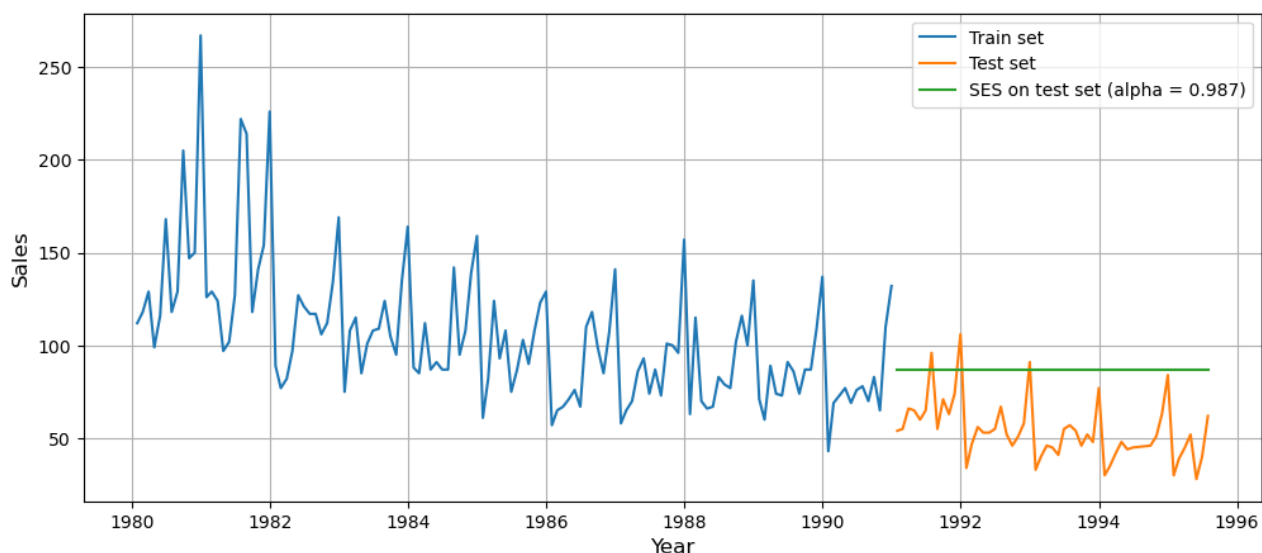


Figure 20: SES forecast ($\alpha = 0.987$)

The predictions for the future remain constant, neither capturing the trend nor seasonality.

RMSE on test set is 36.796.

We can set different values for α to check whether or not the performance of the SES model improves. The higher the alpha value, more weightage is given to the more recent observation. That means, what happened recently will happen again.

The details of how different values of α were chose are given in the appendix.

Alpha Values	Train RMSE	Test RMSE
0.3	32.470164	47.504821
0.4	33.035130	53.767406
0.5	33.682839	59.641786
0.6	34.441171	64.971288
0.7	35.323261	69.698162
0.8	36.334596	73.773992
0.9	37.482782	77.139276

Table 10: RMSE for different α values

RMSE on the test data is the least for $\alpha = 0.3$.

Therefore, another SES model will be built with $\alpha = 0.3$.

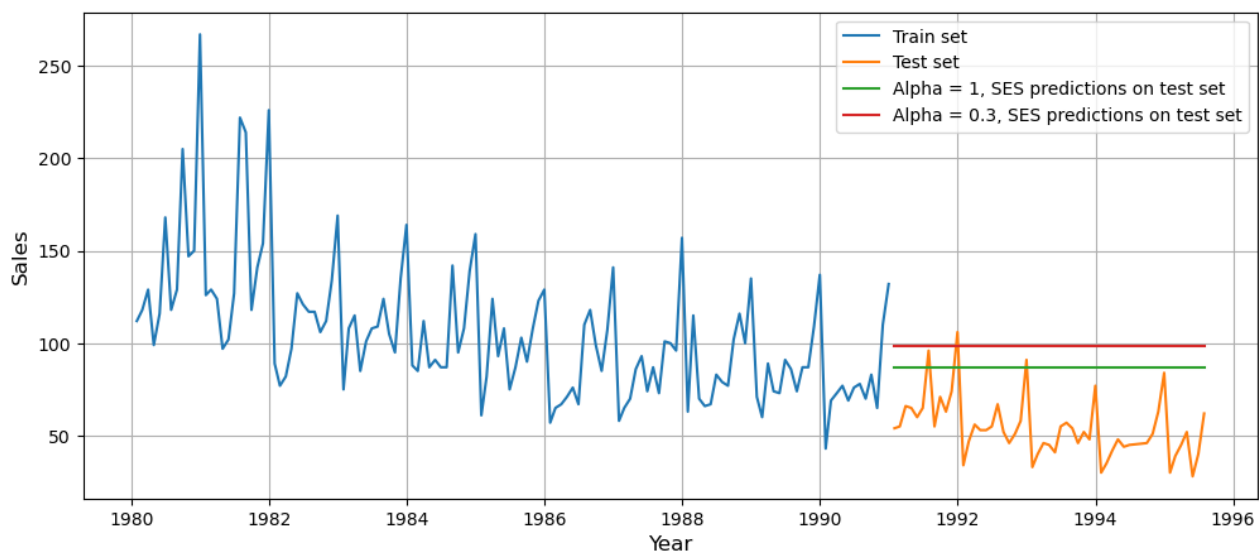


Figure 21: SES forecast ($\alpha = 0.3$)

This model could not have performed better because the value of α was less than the one in the previous model.

RMSE on test set is 47.505.

Double Exponential Smoothing

Double Exponential Smoothing (DES), also called the Holt's method, is an extension of the SES. It is applicable when data has trend, but no seasonality.

It has two smoothing parameters – level (α) and trend (β). Both α and β lie between 0 and 1.

For the DES model, the parameters are as follows:

$$\alpha = 0, \beta = 0$$

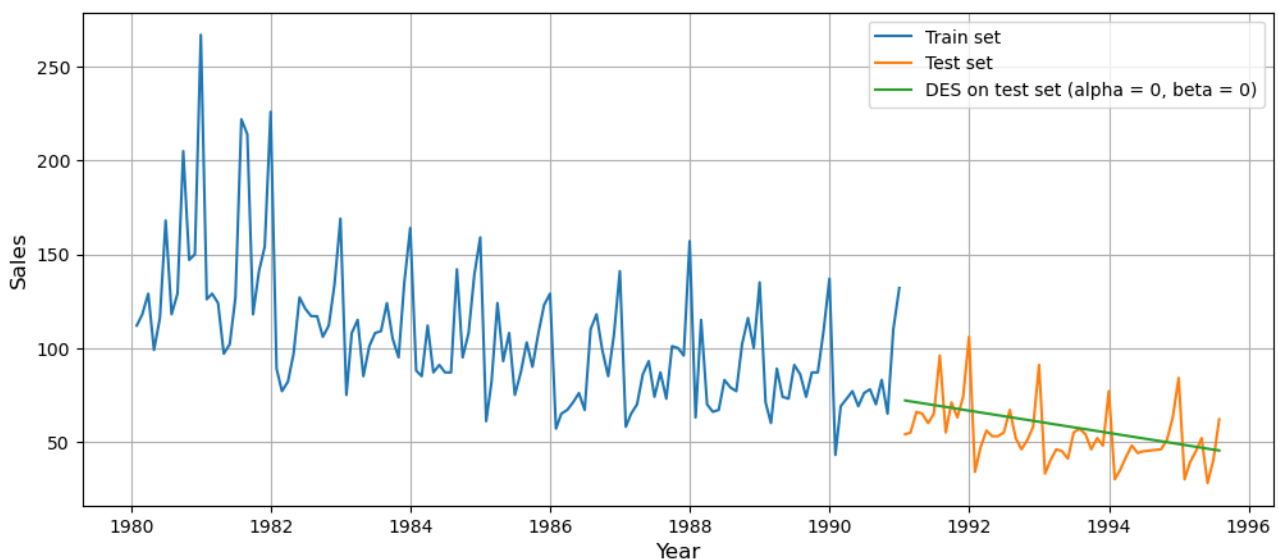


Figure 22: DES forecast ($\alpha = 0, \beta = 0$)

The DES model captures the decreasing trend. Since the time series data has seasonality as well, a Triple Exponential Model must be built.

RMSE on test set is 15.269.

This value is the same as that of the linear regression model.

We can set different values for α and β to check whether or not the performance of the DES model improves.

The details of how the different values of α and β were chosen are given in the appendix.

Alpha values	Beta values	Train RMSE	Test RMSE
0.1	0.1	34.439111	36.923416
0.1	0.2	33.450729	48.688648
0.2	0.1	33.097427	65.731702
0.1	0.3	33.145789	78.156641
0.3	0.1	33.611269	98.653317

Table 11: RMSE for different α and β values

RMSE on the test set is the lowest for $\alpha = 0.1$ and $\beta = 0.1$.

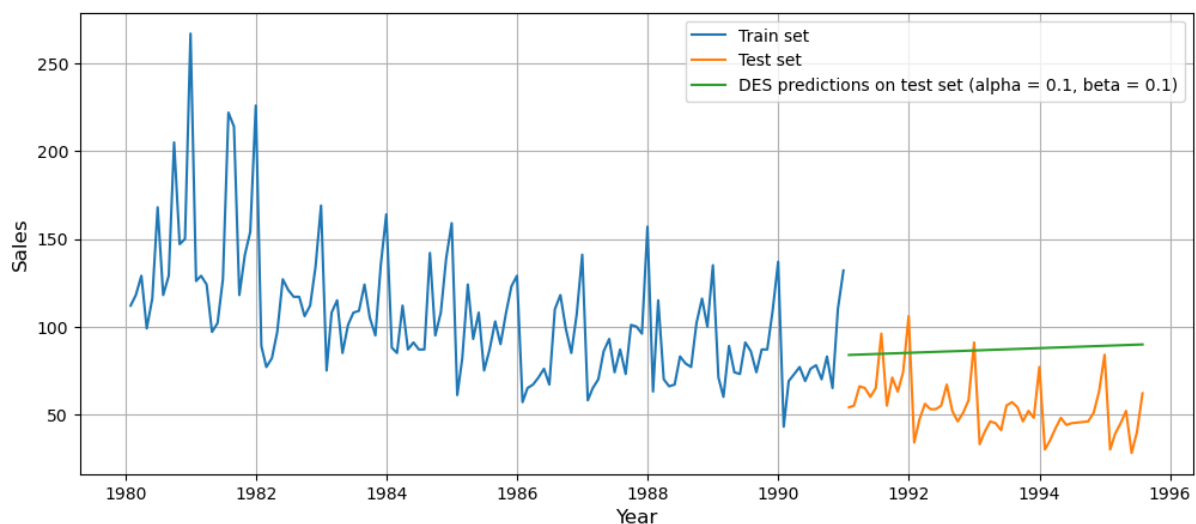


Figure 23: DES forecast ($\alpha = 0.1, \beta = 0.1$)

This model performs poorly as compared with the previous DES model. Unlike the previous model, this model is not even factoring in the trend.

RMSE on tests set is 36.923.

The value is higher than that of the previous DES model.

Triple Exponential Smoothing

Triple Exponential Smoothing (TES), also called the Holt's Winter method, is applicable when data has trend and seasonality.

It has three smoothing parameters – level (α), trend (β) and seasonality (γ). The three parameters lie between 0 and 1.

Since seasonality is multiplicative and additive, we will build one model for each.

Multiplicative seasonality

For the TES model, the parameters are as follows:

$$\alpha = 0.06, \beta = 0.051, \gamma = 0$$

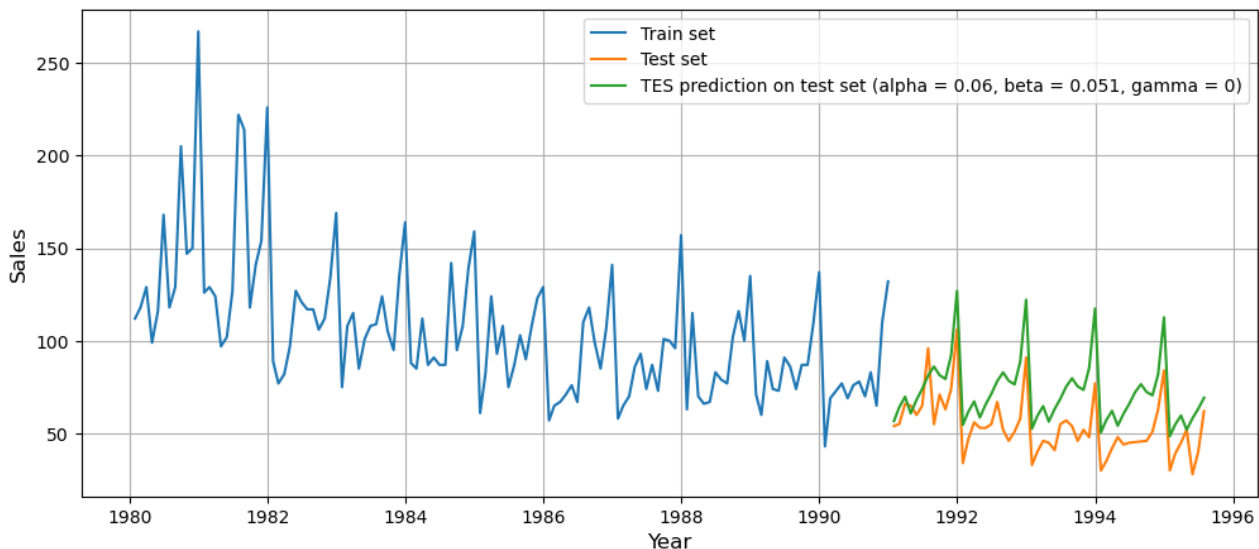


Figure 24: TES forecast ($\alpha = 0.06, \beta = 0.051, \gamma = 0$)

The TES model captures the trend as well as the seasonality. The TES model seems to be the best model so far.

RMSE on test set is 20.9497.

Additive seasonality

For the TES model, the parameters are as follows:

$$\alpha = 0.088, \beta = 0.0001, \gamma = 0.0023$$

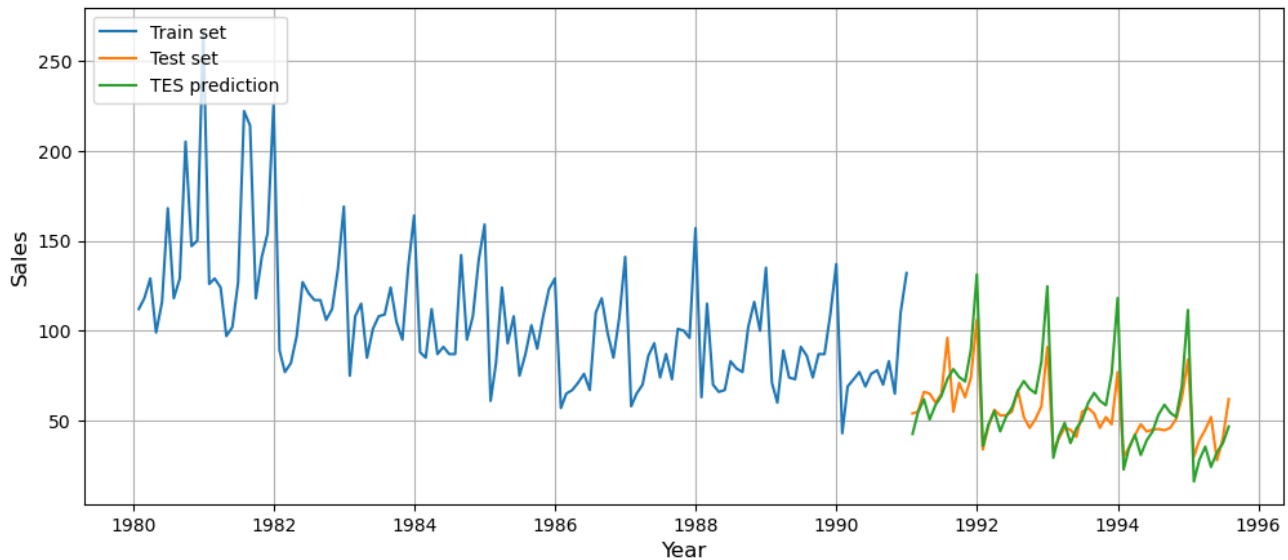


Figure 25: TES forecast ($\alpha = 0.088, \beta = 0.0001, \gamma = 0.0023$)

The TES forecast on the test set captures both trend and seasonality.

RMSE on test set is 14.298.

This is the least RMSE value so far.

We can set different values for α , β and γ to check whether or not the performance of the TES model improves.

The details of how different values of α , β and γ were chosen are given in the appendix.

Alpha values	Beta values	Gamma values	Train RMSE	Test RMSE
0.1	1.0	0.2	24.271787	27.612066
0.1	1.0	0.2	24.271787	27.612066
0.1	1.0	0.2	24.271787	27.612066
0.1	1.0	0.2	24.271787	27.612066
0.1	1.0	0.2	24.271787	27.612066

Table 12: RMSE for different α , β and γ values

RMSE on the test set is the least for $\alpha = 0.1$, $\beta = 1$ and $\gamma = 0.2$.

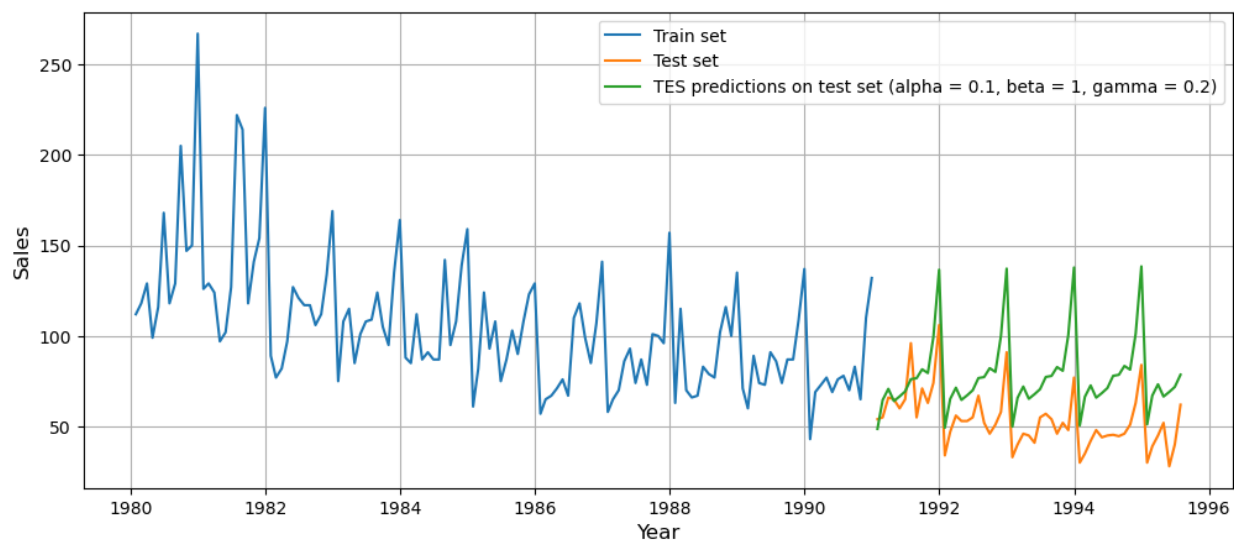


Figure 26: TES forecast ($\alpha = 0.1$, $\beta = 1$, $\gamma = 0.2$)

The TES model built on the basis of the best parameters fails to capture the decreasing trend on the unseen data. Even seasonality is not factored in entirety.

RMSE on test set is 27.612.

The performance of the model does not improve.

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 $\alpha = 0.05$.

A time series is considered to be stationary when statistical properties such as the variance and (auto) correlation are constant over time.

Dickey-Fuller Test on the time series is run to check for stationarity of data.

Null hypothesis (H_0): Time series is non-stationary

Alternative hypothesis (H_1): Time series is stationary

If $p\text{-value} < 0.05$, the null hypothesis is rejected. This means that time series is stationary.

If $p\text{-value} > 0.05$, we fail to reject the null hypothesis. This means that time series is non-stationary.

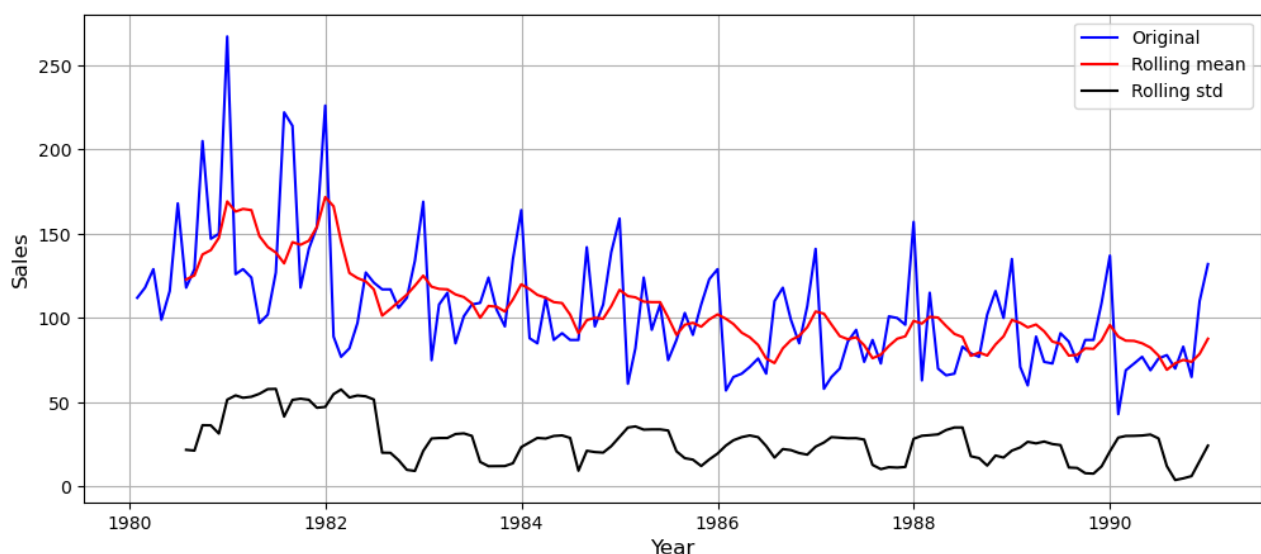


Figure 27: Non-stationary time series plot

The time series data has trend and seasonality.

At α (level of significance) = 0.05, $p\text{-value}$ is 0.219476.

Since the $p\text{-value} > 0.05$, we fail to reject the null hypothesis. This means that **time series is non-stationary**.

Differencing ('d') is done on a non-stationary time series data one or more times to convert it into a stationary time series data.

First-order differencing ($d=1$) is done where the difference between the current and previous (one lag before) series is taken and then checked for stationarity using the Dickey-Fuller test.

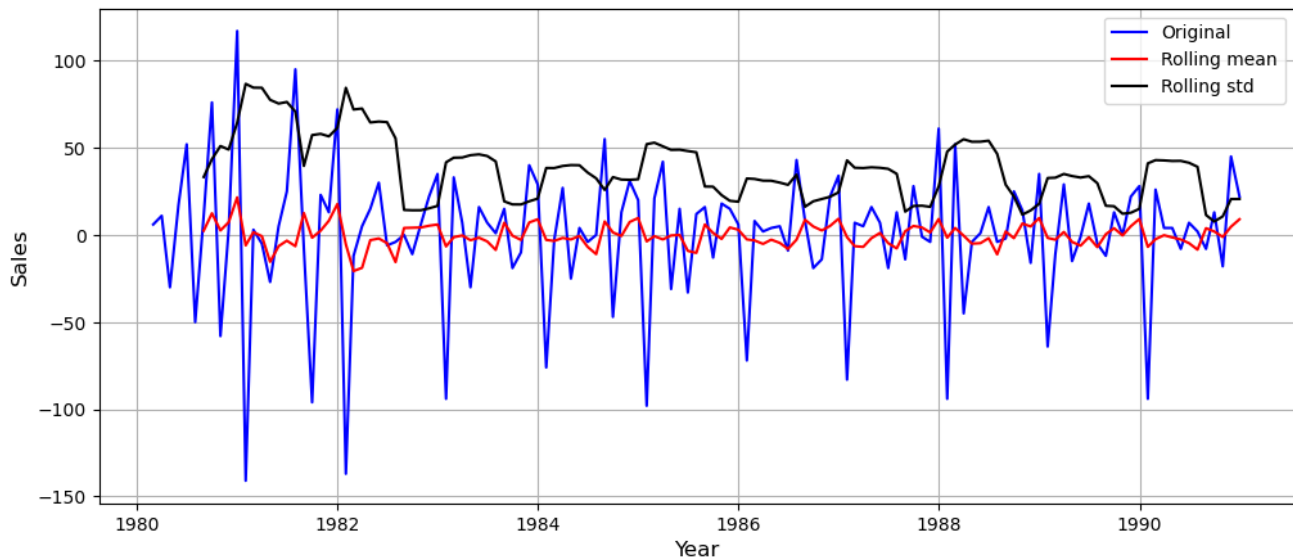


Figure 28: Stationary time series plot ($d = 1$)

At $\alpha = 0.05$, p-value is $7.061944e-09$.

Since the p-value < 0.05 , the null hypothesis is rejected. This means that **time series is stationary at $d = 1$** .

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 model

An ARIMA model has three components:

1. **Autoregressive (AR) model:** AR models use previous time period values to predict the current time period values. The number of lag observations or autoregressive terms in an AR model is denoted by '**p**'. One of the fundamental assumptions of an AR model is that the time series is assumed to be a stationary process.
2. **Integrated component:** It is the difference in the non-seasonal observations. It is denoted by '**d**'.
3. **Moving Average (MA) model:** MA models estimate the future values based on historical forecast errors. The size of the moving average window is denoted by '**q**'.

One of the approaches to find the order of '**p**' and '**q**' is the least value of **Akaike Information Criteria (AIC)**. The AIC values for different pairs of '**p**' and '**q**' are compared to find the optimum order for model-building.

For Rose wine sales time series data, we can set different values for '**p**' and '**q**' to find the combination that gives the least AIC value. The details are given in the appendix.

Here, **d = 1** because at first-order differencing, the time series becomes stationary.

Params	AIC
(0, 1, 2)	1279.671529
(1, 1, 2)	1279.870723
(1, 1, 1)	1280.574230
(2, 1, 1)	1281.507862
(2, 1, 2)	1281.870722

The AIC value is the least for **p = 0, d = 1, q = 2**.

An ARIMA model will be built on the basis of these parameters.

Table 13: ARIMA parameters

SARIMAX Results						
=====						
Dep. Variable:	Rose	No. Observations:	132			
Model:	ARIMA(0, 1, 2)	Log Likelihood	-636.836			
Date:	Sun, 14 May 2023	AIC	1279.672			
Time:	11:43:39	BIC	1288.297			
Sample:	01-31-1980	HQIC	1283.176			
	- 12-31-1990					
Covariance Type:	opg					
=====						
	coef	std err	z	P> z	[0.025	0.975]

ma.L1	-0.6970	0.072	-9.689	0.000	-0.838	-0.556
ma.L2	-0.2042	0.073	-2.794	0.005	-0.347	-0.061
sigma2	965.8407	88.305	10.938	0.000	792.766	1138.915
=====						
Ljung-Box (L1) (Q):	0.14	Jarque-Bera (JB):	39.24			
Prob(Q):	0.71	Prob(JB):	0.00			
Heteroskedasticity (H):	0.36	Skew:	0.82			
Prob(H) (two-sided):	0.00	Kurtosis:	5.13			
=====						

Table 14: ARIMA summary

The ARIMA model (0, 1, 2) has two MA terms.

It can be observed that each variable has a p value. Each variable has a null (H_0) and an alternative (H_1) hypothesis.

H_0 : Variable is not significant

H_1 : Variable is significant

At 5% level of significance, p-values exceeding 0.05 will mean that an independent variable is not significant and, hence, can be dropped.

From the table, it can be seen that p-values for all variables are less than 0.05. Therefore, we reject the null hypothesis. Hence, **the variables are significant**.

RMSE on test set is 37.327.

SARIMA model

A SARIMA model has the following components:

- Autoregressive (**p**) and moving average (**q**) components
- Seasonal autoregressive (**P**) and moving average (**Q**) components
- Ordinary and seasonal difference components of order '**d**' and '**D**', respectively
- Seasonal frequency (**F**)

To find the seasonal parameter for the SARIMA model, we need to look at the **Autocorrelation Function (ACF) plot**, which summarises the correlation of an observation with lag values.

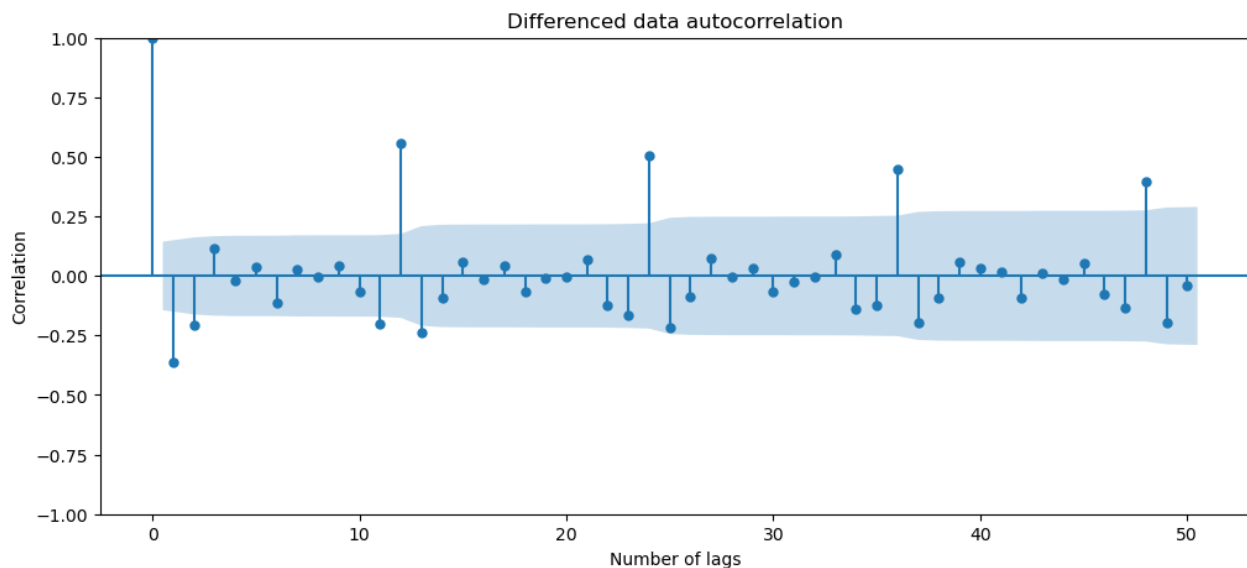


Figure 29: ACF plot

We see that there can be a seasonality of 12. Therefore, we will run an auto SARIMA model by setting the **seasonality (F) as 12** and **d = 1 (for stationarity)**.

For other parameters, we can set different values for '**p**', '**q**', '**P**' and '**Q**' to find the combination that gives the least AIC value. The details of how the parameters were chosen are given in the appendix.

Param	Seasonal	AIC
(0, 1, 2)	(2, 0, 2, 12)	887.937509
(1, 1, 2)	(2, 0, 2, 12)	889.901291
(2, 1, 2)	(2, 0, 2, 12)	890.668798
(2, 1, 1)	(2, 0, 0, 12)	896.518161
(2, 1, 2)	(2, 0, 0, 12)	897.346444

The AIC value is the least for the following combination of parameters.

p = 0, q = 2

P = 2, Q = 2

d = 1, D = 0

Table 15: SARIMA parameters

On building a SARIMA model on the basis of the above parameters, we get the following results.

SARIMAX Results						
Dep. Variable:	y			No. Observations:	132	
Model:	SARIMAX(0, 1, 2)x(2, 0, 2, 12)			Log Likelihood	-436.969	
Date:	Sun, 14 May 2023			AIC	887.938	
Time:	11:45:37			BIC	906.448	
Sample:	0			HQIC	895.437	
	- 132					
Covariance Type:	opg					
	coef	std err	z	P> z	[0.025	0.975]
ma.L1	-0.8427	189.825	-0.004	0.996	-372.894	371.208
ma.L2	-0.1573	29.823	-0.005	0.996	-58.608	58.294
ar.S.L12	0.3467	0.079	4.375	0.000	0.191	0.502
ar.S.L24	0.3023	0.076	3.996	0.000	0.154	0.451
ma.S.L12	0.0767	0.133	0.577	0.564	-0.184	0.337
ma.S.L24	-0.0726	0.146	-0.498	0.618	-0.358	0.213
sigma2	251.3137	4.77e+04	0.005	0.996	-9.33e+04	9.38e+04
Ljung-Box (L1) (Q):	0.10		Jarque-Bera (JB):	2.33		
Prob(Q):	0.75		Prob(JB):	0.31		
Heteroskedasticity (H):	0.88		Skew:	0.37		
Prob(H) (two-sided):	0.70		Kurtosis:	3.03		

Table 16: SARIMA summary

The SARIMA model has the following components:

- Two MA variables
- One seasonal AR variable with a lag of 12
- One seasonal AR variable with a lag of 24
- One seasonal MA variable with a lag of 12
- One seasonal MA variable with a lag of 24

The p-values of all components but two are less than 0.05. This means that only **two components out of six are significant**.

RMSE on test set is 26.948.

The RMSE value of the SARIMA model is less than that of the ARIMA model. In other words, **SARIMA model performs better on the unseen data.**

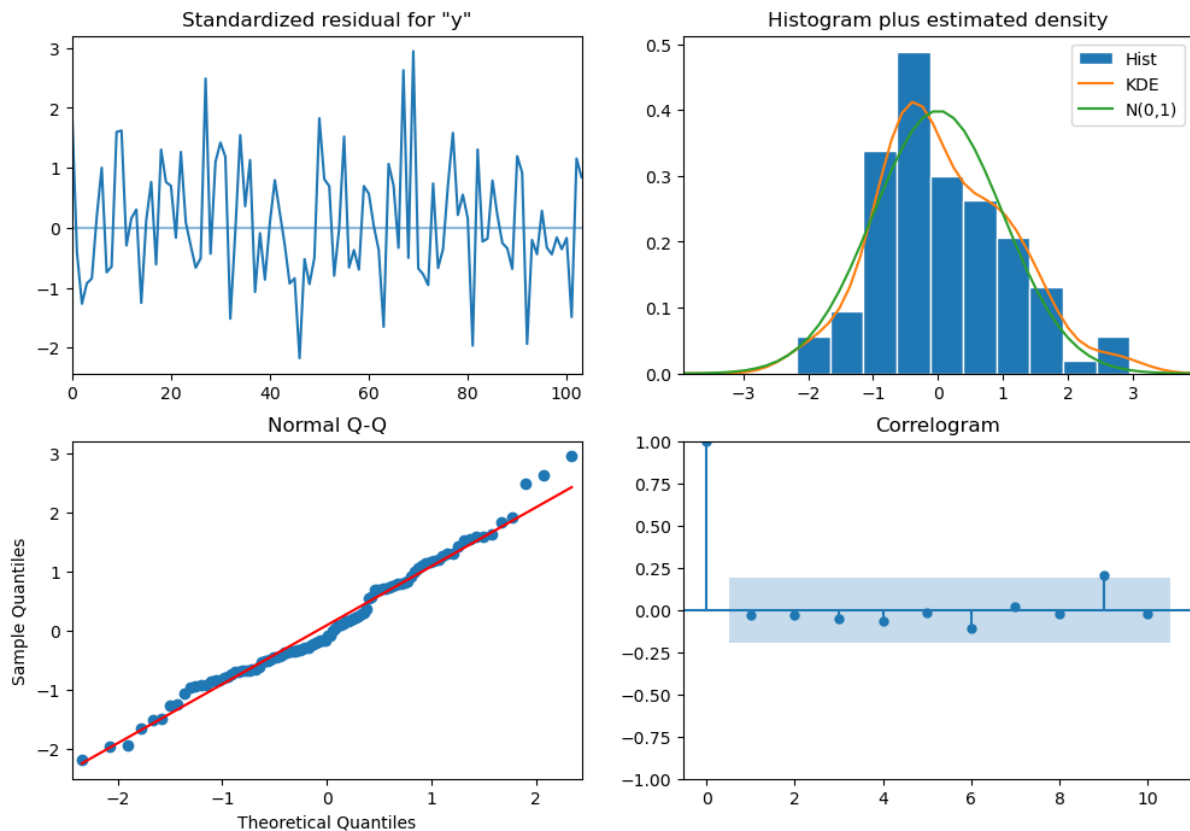


Figure 30: Diagnostic plots

The diagnostic graphs plot the residuals. All four diagnostics plots almost follow the theoretical numbers and thus we cannot observe any pattern from these plots.

The first plot shows the standardised error. Residuals do not follow a pattern.

The histogram and the normal Q-Q plot show that the residuals almost follow the normal distribution.

Correlogram tells us the correlation of residuals with time. It can be seen that all correlations are insignificant. This should be the case, as we do not want residuals to have correlation with time.

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

	Test RMSE
Triple Exp Smoothing (alpha = 0.088, beta = 0.0001, gamma = 0.0023)	14.297933
Double Exp Smoothing (alpha = 0, beta = 0)	15.275507
Regression on time	15.275520
Triple Exp Smoothing (alpha = 0.06, beta = 0.051, gamma = 0)	20.976545
SARIMA (0, 1, 2) (2, 0, 2, 12)	26.948030
Triple Exp Smoothing (alpha = 0.1, beta = 1, gamma = 0.2)	27.612066
9-point trailing Moving Average	34.431243
Simple Exp Smoothing (alpha = 0.987)	36.816502
Double Exp Smoothing (alpha = 0.1, beta = 0.1)	36.944359
ARIMA (0, 1, 2)	37.326663
6-point trailing Moving Average	39.146677
4-point trailing Moving Average	46.423686
Simple Exp Smoothing (alpha = 0.3)	47.524854
Simple Average Model	53.480456
2-point trailing Moving Average	68.989714
Naive Model	79.738146

Table 17: Comparison of all models

Of all the models built, the Triple Exponential Smoothing has the least RMSE. This model has the following parameters:

Level (α) = 0.088

Trend (β) = 0.0001

Seasonality (γ) = 0.0023

8. 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.

The most optimum model based on the least RMSE is the Triple Exponential Smoothing. A full model will be built on the complete dataset and then we will predict for future 12 months.

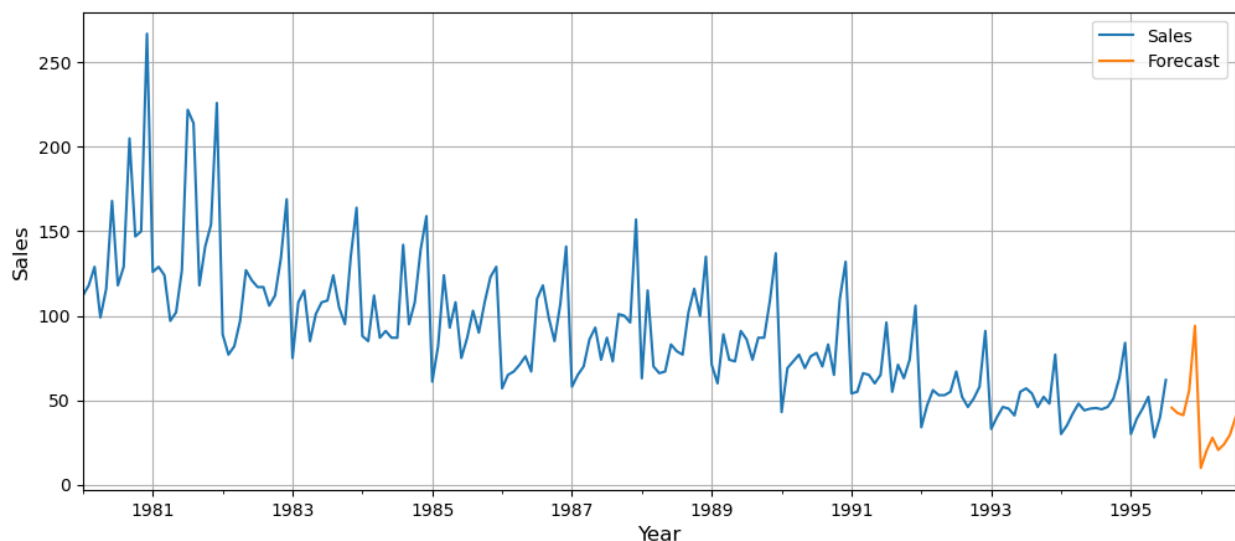


Figure 31: Sales forecast for 12 months

The graph shows the forecast for future 12 months. The prediction is in line with the time series data, as it captures both trend and seasonality.

However, this prediction (denoted by the orange graph) is very precise. In the real-world scenario, one cannot be sure about the forecast for future. Therefore, we will predict for future within a certain confidence interval.

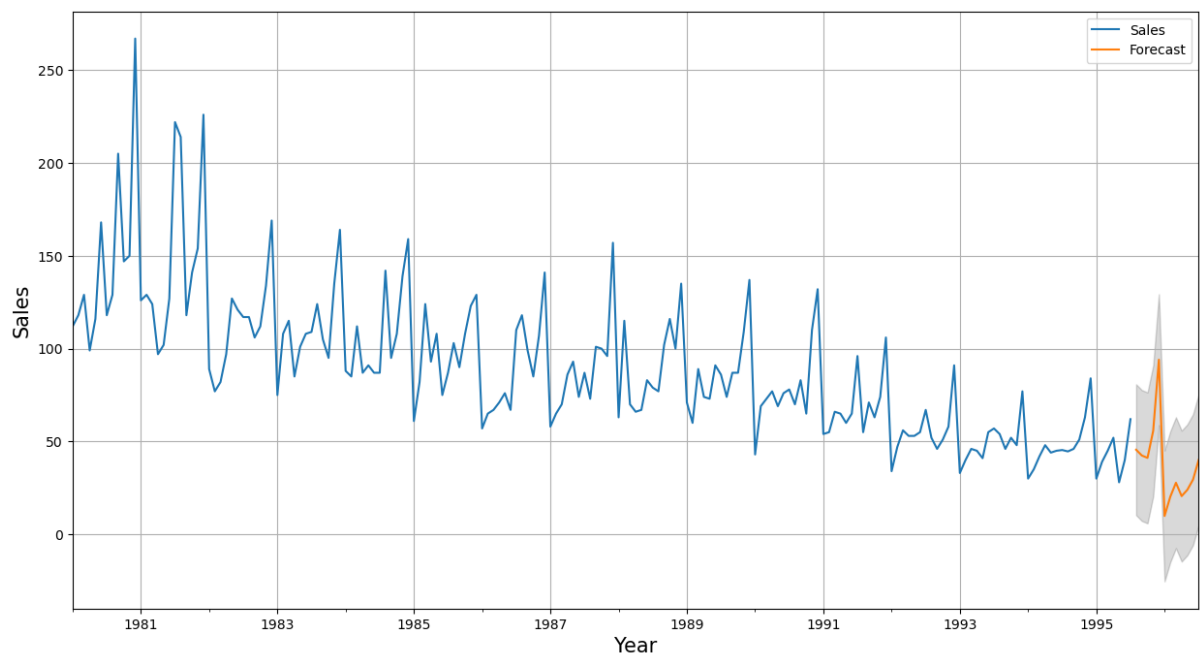


Figure 32: Sales forecast with confidence interval

The forecast values lie within the grey band, which represents 95 per cent confidence interval.

One assumption that we have made over here while calculating the confidence bands is that the standard deviation of the forecast distribution is almost equal to the residual standard deviation.

RMSE on the full model is 17.941.

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

Please explain and summarise the various steps performed in this project. There should be proper business interpretation and actionable insights present.

Exploratory Data Analysis summary

- The Rose time series data has a decreasing trend.
- As for seasonality, the sales remain low at the start of the year, but pick up pace in the subsequent months. The last two months record high sales, with the highest being in December. This pattern is repeated every year.
- The median sales were high in the first half of the decade of 1980. As the years rolled by, the sales diminished.
- The quarterly median sales are the highest in quarter four and the lowest in quarter one.
- In any given year, December has the highest sales. However, the December sales fall over time.
- Median sales on all seven days of a week is more or less the same. Surprisingly, sales on week days are more than that of Sunday's.

Model-building summary

- Most of the models, be it Naïve, moving average or Double Exponential Smoothing, fail to capture the decreasing trend and seasonality.
- Triple Exponential Smoothing models are able to factor in both trend and seasonality.
- The SARIMA model, too, captures the systematic components. However, its RMSE value is more than that of the TES model with optimum parameters.

Recommendations

- Sales in quarter four should not be the company's concern, as demand is high at the end of the year, may be because of New Year's eve and Christmas.
- However, sales in January fall drastically when compared with December of the previous year. The company can continue to offer festive season

discounts in January so that sales remain high at the start of the year as well.

- The company needs to focus on the first three quarters of the year. The advertising and marketing campaigns should be targeted at different segments of customers.
- Another strategy could be customer-specific social media campaigns.
- In summers, special discounts must be offered on crisp rose and white wines that are perfect for quenching the thirst during summers.
- Wine tasting events are a good way to attract new customers. The company can invite celebrities and sommeliers for such events.

Appendix

Simple Exponential Smoothing

```
for i in np.arange(0.3, 1, 0.1):  
    print(i)
```

Double Exponential Smoothing

```
for i in np.arange(0.1, 1.1, 0.1):  
    for j in np.arange(0.1, 1.1, 0.1)
```

Triple Exponential Smoothing

```
for i in np.arange(0.1, 1, 0.1):  
    for i in np.arange(0.1, 1, 0.1):  
        for k in np.arange(0.1, 1, 0.1):
```

ARIMA model

```
p = q = range(0, 3)  
d= range(1,2)
```

SARIMA model

```
p = q = range(0, 3)  
d = range(1,2)  
D = range(0,1)  
P =Q = range(0,3)
```

Source: Great Learning logo that has been used on the cover page has been taken from one of the monographs provided by Great Learning