# V.Ger Travel Company - An Analysis

## ITNBD4 Assignment

Student: 2710017

10 January 2025

# Table of Contents

# 1 Summary

Summary of the report

# 2 Introduction

## 2.1 Background & Approach

The travel conglomerate V.Ger Travel (VGT) has a broad range of operations, including hotels, resorts, car rentals and also air travel through charter flights. Travel bookings originate primarily from VG's own travel web site which is supported by several operational information systems that cover all aspects of it's business from customer relations through to logistics and maintenance. There is a wealth of data available with over 10 years history of travel bookings, however, the use of modern Data Science methods to harness this data is in its infancy at VGT.

This report describes the recommendations of the new Chief Data Officer (CDO) as to how to implement modern Data Science techniques to utilise VGT's data to improve its efficiency and profit. Having completed an initial high level review, two opportunity areas where identified that appeared to offer the biggest opportunities for increasing efficiency and also that are achievable as a first step in implementing new techniques. The two use cases are:

- Hotel Demand Forecasting
- Customer Satisfaction & Loyalty

Each of these use cases are explored in the following sections, each of which describes:

- Business Scope & Benefits - What is proposed and why is this helpful
- Data Analysis Approach - The suggested analysis techniques and the data required
- Simulated Data Analysis - An example of the analysis using simulated data
- Conclusions & Next Steps - Findings and how to move forward

***????? - Assumptions made about data prep/data wrangling .. eg normalisation of data etc - Any assumptions about company size, number of hotels etc ????***

## 2.2 Results & Further Details

A summary of the approach taken and key results of the analysis are described for each use case in the individual sections of this report. Additionally, the all details for generating the simulated data and the analysis steps are described at the following:

- Apendices - Step by step approach, with data tables and plots
- Jupyter Notebooks - Full listings of Python code and results in a GitHub Repository
- Data Files - Generated CSV files in a GitHub Repository

## 2.3 Conclusions

*???? ConclusionshHere as well as last past of each use case ????*

# 3 Hotel Demand Forecasting

## 3.1 Business Scope & Benefits

The hotel operations of VGT are a significant part of its business and any efficiencies in this area have the potential to make large contributions to the overall profitability of the business. The profit contribution from individual hotels can be maximised by ensuring revenue is as high as possible and at the same time minimising the hotel's operating costs. One way to do this is to provide a hotel's management team with the tools to carry out reliable demand forecasting.

Hotel demand forecasting is the prediction of the demand for rooms and related hotel services to help a hotel's management team determine pricing, staffing and marketing strategies (Johansson, 2022). If the demand for hotel rooms can be reliably forecast then this enables:

- Dynamic Pricing - Adjust future prices in response to forecast demand. When high demand is expected then future room rates can be increased; when low demand is expected then discounts can be offered or packages can be advertised. And marketing strategies can be determined to respond to the demand forecasts. xx increase occupancy and revenue
- Staffing Levels - Hotel staffing can be adjusted to maintain customer service levels but not over staff when demand is expected to be lower. xx control costs
- Inventory Management - Similarly inventory can be adjusted, for example catering supplies maintained just sufficient to meet the forecast number of customers using catering facilities. xx control costs

There are many factors that will influence the demand for hotel rooms, some may remain stable whilst others are less so and will vary over time or in response to external factor, these include:

- Location & Market - Is the hotel a budget or a boutique hotel? Is is in a business district, near a beach or a ski resort etc. Are customers mainly business travelers or tourists?
- Economic - Macro level impact from the state of the economy.
- Local Competition - Competition with local hotels.
- Seasonality - When are the high and low seasons? Is it a summer or a winter resort? What are the local weather patterns and seasons. Are there weather related attractions?
- Local Events - When are any local festivities, sports events, school holidays, religious events, music conferences, business conferences?

### 3.1.1 Occupancy & Revenue Indicators

An important indicator of demand is the occupancy rate (Jeffrey and Hubbard, 1994) and (FHA, 2023) which is simply the percentage of the total rooms occupied for a given time period. The occupancy rate varies across the industry but a target of 60% to 80% is typical.

Occupancy is used alongside revenue related indicators to provide a measure of revenue health, the three main indicators (for a given time period, eg daily) are:

- Occupancy Rate (OCC%) - Percentage of available rooms that are occupied, or expected to be occupied.
- Average Daily Rate (ADR) - Average revenue per room occupied, across all room price bands.
- Revenue Per Available Room (RevPAR) - Revenue reflecting all available rooms. Calculated by: OCC% * ADR. A good overall indicator of revenue.

### 3.1.2 Forecasting Occupancy & Business Benefits

The time period used for forecasting will vary depending on the objectives desired (Lighthouse, 2024) and (Lighthouse, 2023) for example:

- Short Term - Forecast occupancy for next month so room pricing can be adjusted appropriately.
- Long-term - Forecast occupancy for next year so price bands and packages can be set, marketing strategies defined and required staffing levels determined.

To recap, if occupancy (and the associated revenue indicators) can be reliably forecast then plans can be put in place to maximise revenue by adjusting pricing and marketing strategies and controlling costs by flexing staffing and inventory levels.

## 3.2 Data Analysis Approach

At VGT, no rigorous forecasting is currently in place so to begin with a relatively simple approach will be implemented in a small number of hotels. If the benefits of this are confirmed, then it can then be extended in sophistication using more complex forecasting models and for longer time periods. It can then be implemented across all hotels in VGT.

The objective of this first step is to establish a model that can be used to forecast the daily occupancy (OCC%) at an individual hotel for the coming month, ie the forecast is for 30 to 40 days in the future. The forecast will then be used by hotel's management team to: i) adopt dynamic pricing; ii) execute supporting short-term marketing; iii) fine-tune staffing rotas and holiday leave for the coming weeks. This should improve the efficiency of the hotel's operations by increasing room bookings whilst ensuring staffing costs are controlled at an appropriate level.

### 3.2.1 Forecasting Model & Data Required

The scope of the envisaged forecasting model is to calculate a month of daily OCC% for each room category in an hotel. The output of the forecast model will be provided in spreadsheet form so that the hotel management team can manually make adjustments to try to improve revenue and staffing levels. The actual occupancy and revenue can then be tracked against the forecast throughout the month in order to assess how accurate the forecasting model is and to help refine it.

*Forecasting Model Data*

The data required for the forecasting model is the last 4 years of daily room occupancy. In this analysis a single hotel with two classes of room (standard and premium) will be used. A 4 year history of daily room bookings was selected with 3 years used to be used for training and 1 year for validation. This length of history was chosen as a starting point because a previous occupancy study (Phumchusri and Suwatanapongched, 2023) found that the choice could be quite dependent on the scenario; however, there needs to be a a balance of too short and missing seasonality vs too long and not being sufficiently responsive. Phumchusri and Suwatanapongched (2023) also found that using 4 years history of daily occupancy to forecast 2 to 8 weeks was a good approach.

The specific data required is:

- For an individual hotel and each room category (standard and premium)
- Daily
- Room capacity
- Room rate
- Rooms occupied

*Tracking Spreadsheets*

The data comprising the revenue forecasting and tracking spreadsheet is:

- Forecast OCC% (derived from the forecasting model)
- Daily room rates by room category that can be manually adjusted
- Daily Revenue, ADR, RevPAR (derived from the OCC% and room rates)
- Special events, a facility to mark local events that may impact demand. For example, sports events, concerts, conference, unusual weather forecasts
- Actual OCC%, Revenue, ADR, RevPAR

The data provided for the staffing forecasting spreadsheet is:

- Forecast OCC% (from the forecasting model)
- Averaged staff requirements per room
- Staffing levels (derived from the OCC%)
- Actual staffing levels

### 3.2.2 Techniques Considered

There are several potential tools that can be used with historical time series data to forecast future occupancy rates, ranging from established 'simpler' techniques such as linear regression through to more sophisticated machine learning and neural network models (Huang and Zheng, 2022). However, given the relative infancy of Data Science at VGT, the use of more sophisticated tools will be prioritised for future work. In this investigation, the following techniques will be examined:

1. Ordinary Least Squares (OLS) linear regression
2. ARIMA, SARIMA - Use to fully incorporate seasonality
3. SAIMAX - If exogenous factors need to be accounted for
4. *?? LightGBM .... if time allows!*

**OLS Linear Regression**

Investigating techniques for reliably identifying the factors influencing hotel occupancy has been ongoing for several years. See Andrew, Cranage and Lee (1990) and Jeffrey and Hubbard (1994) for earlier work focusing on the use of regression analysis of time series data. Although it is likely that the nature of hotel business means that occupancy will be seasonal, linear regression will still be investigated first.

**ARIMA, SARIMA**

Occupancy forecasting using ARIMA using factors such as room capacity and marketing expenditure has successfully been used (Chow, Shyu and Wang, 1998). A comparison of forecasting methods (Weatherford and Kimes, 2003) included historical time series analysis of occupancy using ARIMA. The best time period used was not clear and it is a balance of too short and missing seasonality vs too long and not being sufficiently responsive. Also the best analysis method appears to depend on the characteristics of individual hotels and hotel chains. Using a SARIMA approach using 4 years history of daily occupancy to forecast 2 to 8 weeks was found to be a strong approach (Phumchusri and Suwatanapongched, 2023).

**SARIMAX**

**?? LightGBM …. if time allows!**

## 3.3 Simulated Data Analysis

### 3.3.1 Data Summary

A dataset was created to simulate the data as defined in the earlier section. The data elements consist of:

|  | Count | Missing | Empty | Unique | Type | String | Int | Float | List |
|---|---|---|---|---|---|---|---|---|---|
| Standard_OCC | 1461 | 0 | 0 | 177 | int64 | 0 | 1461 | 0 | 0 |
| Standard_Capacity | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |
| Standard_Rate | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |
| Premium_OCC | 1461 | 0 | 0 | 101 | int64 | 0 | 1461 | 0 | 0 |
| Premium_Capacity | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |
| Premium_Rate | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |

The first few elements of the loaded data:

```
            Standard_OCC  Standard_Capacity  Standard_Rate  Premium_OCC  \
Date
2020-01-01           129                254            325           65
2020-01-02           126                254            325           53
2020-01-03           137                254            325           63
```

```
          Premium_Capacity   Premium_Rate
Date
2020-01-01               100            575
2020-01-02               100            575
2020-01-03               100            575
```

And key descriptive statistics:

```
        Standard_OCC  Standard_Capacity  Standard_Rate  Premium_OCC  \
count   1461.000000             1461.0         1461.0  1461.000000
mean     135.646133              254.0          325.0    56.141684
std       50.216654                0.0            0.0    29.349613
min       38.000000              254.0          325.0     0.000000
25%       88.000000              254.0          325.0    31.000000
50%      136.000000              254.0          325.0    57.000000
75%      183.000000              254.0          325.0    83.000000
max      254.000000              254.0          325.0   100.000000


        Premium_Capacity   Premium_Rate
count             1461.0         1461.0
mean               100.0          575.0
std                  0.0            0.0
min                100.0          575.0
25%                100.0          575.0
50%                100.0          575.0
75%                100.0          575.0
max                100.0          575.0
```

### 3.3.2 Time Series Characteristics

The occupancy time series for the two categories of room are shown in the figure below. This indicates that the occupancy has a strong seasonality with an annual peak and trough; there are also regular spikes to full capacity bookings. The premium room occupancy hits maximum and zero occupancy on several occasions.

The two occupancy time series were examined further (see details at the appendix) using lag plots, ACF plots and ADF tests which indicated autocorrelation, annual seasonality and a positive trend. Differencing was also completed to confirm non-stationarity. Finally a decomposition was completed and this confirmed the seasonality and trend, see the figure below.
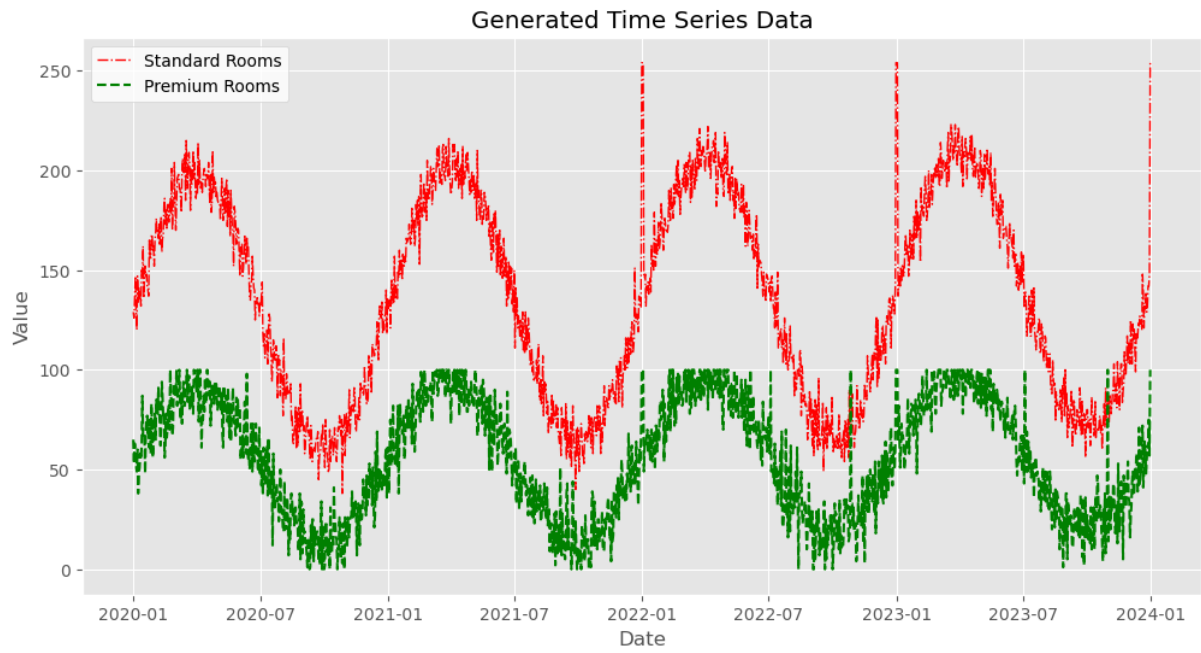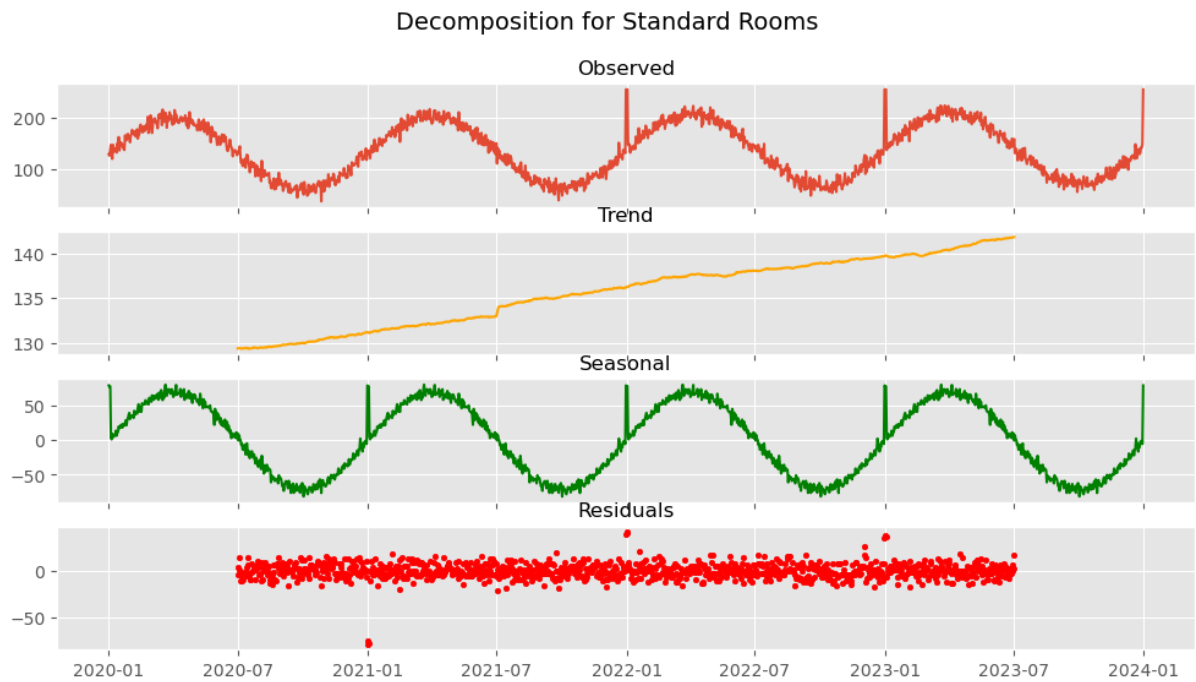
Figure 3.1: Occupancy for Each Room Category



Figure 3.2: Occupancy Time Series Decomposition

Decomposition for Premium Rooms

### 3.3.3 OLS Linear Regression

Given the strong seasonality, linear regression is unlikely to be a good model for forecasting, however a regression fit was calculated to double-check. See the plot below which shows that the fit lines could not reliably provide a forecast and even show a downward trend. And the Durban Watson statistics are less than 2 which confirms evidence of autocorrelation.

### 3.3.4 ARIMA, SARIMA

### 3.3.5 SARIMAX

Test see earlier Section 3.3.3

and earlier **?@fig-test**

### 3.3.6 Execution - Data Analysis

- Use simulated data to carry out analysis

- Show the results, forecasts ….

- *?? carry out PACF to determine the autoregressive order for ARIMA etc*

Figure 3.3: OLS Linear Regression

## 3.4 Conclusions & Next Steps

### 3.4.1 Findings

xxx

### 3.4.2 Next Steps

- ?? how are special events impactful on revenue predictions, can these be better used in the forecast going forward or in an improved forecasting model
- ?? recent flight searches, enquiries ... for the area

xxxx

- ?? track actual OCC% etc and compare to forecast to allow model refinement as well as intra-month fine tuning
- ?? Get competitors OCC% to compare and improve forecasting models
- ?? automate the manual pricing changes etc, ie no need for manual adjustment of the forecast spreadsheet
- xx LightGBM ..
- xx CNN, LSTM, RNN ...
- xx transformers, LLM ..

xxx improvements

- ?? isolate pricing strategies and refine the model
- ?? identify the main parameters / impacts on the forecasting ... what are the demand indicators?
- ?? categorisation of customers, families, demographics etc, business, tourist
- ?? identify correlations with holiday patterns, events ......
- ?? identify seasonal bands, high, low, shoulder etc
- ?? identify links between revenue and room price bands ...

# 4 Customer Satisfaction & Loyalty

## 4.1 Business Scope & Benefits

## 4.2 Data Analysis Approach

## 4.3 Simulated Data Analysis

### 4.3.1 Data Summary

A dataset was created to simulate the data as defined in the earlier section. The data elements consist of:

The first few elements of the loaded data:

And key descriptive statistics:

### 4.3.2 XGBoost ….

## 4.4 Conclusions & Next Steps

### 4.4.1 Findings

### 4.4.2 Next Steps

# 5 Conclusions

## 5.1 Conclusions

## 5.2 Next Steps

?????? - Overall conclusions - Or possibly conclusions within each section

# References

Andrew, W.P., Cranage, D.A. and Lee, C.K. (1990) 'Forecasting Hotel Occupancy Rates with Time Series Models: An Empirical Analysis', *Hospitality Research Journal*, 14(2), pp. 173–182. Available at: https://doi.org/10.1177/109634809001400219.

Chow, W.S., Shyu, J.-C. and Wang, K.-C. (1998) 'Developing a Forecast System for Hotel Occupancy Rate Using Integrated ARIMA Models', *Journal of International Hospitality, Leisure & Tourism Management*, 1(3), pp. 55–80. Available at: https://doi.org/10.1300/J268v01n03_05.

FHA, F. (2023) 'Understanding occupancy rate in hotels: A comprehensive guide'. Available at: https://fhahoreca.com/glossary/occupancy-rate/.

Huang, L. and Zheng, W. (2022) 'Hotel demand forecasting: a comprehensive literature review', *Tourism Review*, 78(1), pp. 218–244. Available at: https://doi.org/10.1108/TR-07-2022-0367.

Jeffrey, D. and Hubbard, N.J. (1994) 'A model of hotel occupancy performance for monitoring and marketing in the hotel industry', *International Journal of Hospitality Management*, 13(1), pp. 57–71. Available at: https://doi.org/10.1016/0278-4319(94)90059-0.

Johansson, A. (2022) 'Benefits and challenges in forecasting hotel demand'. Available at: https://www.demandcalendar.com/blog/benefits-and-challenges-in-forecasting-hotel-demand.

Lighthouse, A. (2024) 'Hotel ADR: Everything you need to know about Average Daily Rate'. Available at: https://www.mylighthouse.com/resources/blog/hotel-adr-everything-you-need-to-know-about-average-daily-rate.

Lighthouse, B. (2023) 'The hotelier's ultimate guide to occupancy forecasting'. Available at: https://www.mylighthouse.com/resources/blog/how-to-forecast-hotel-occupancy.

Phumchusri, N. and Suwatanapongched, P. (2023) 'Forecasting hotel daily room demand with transformed data using time series methods', *Journal of Revenue and Pricing Management*, 22(1), pp. 44–56. Available at: https://doi.org/10.1057/s41272-021-00363-6.

Weatherford, L.R. and Kimes, S.E. (2003) 'A comparison of forecasting methods for hotel revenue management', *International Journal of Forecasting*, 19(3), pp. 401–415. Available at: https://doi.org/10.1016/S0169-2070(02)00011-0.

# A Hotel Demand Forecasting - Jupyter Notebook Output

## A.1 Data Load & Characteristics

- Load time series data and look at its characteristics
- Determine autocorrelation, seasonality, stationarity
- Decomposition,
- OLS Regression

### A.1.1 Characteristics

|  | Count | Missing | Empty | Unique | Type | String | Int | Float | List |
|---|---|---|---|---|---|---|---|---|---|
| Standard_OCC | 1461 | 0 | 0 | 177 | int64 | 0 | 1461 | 0 | 0 |
| Standard_Capacity | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |
| Standard_Rate | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |
| Premium_OCC | 1461 | 0 | 0 | 101 | int64 | 0 | 1461 | 0 | 0 |
| Premium_Capacity | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |
| Premium_Rate | 1461 | 0 | 0 | 1 | int64 | 0 | 1461 | 0 | 0 |

```
            Standard_OCC  Standard_Capacity  Standard_Rate  Premium_OCC  \
Date
2020-01-01           129                254            325           65
2020-01-02           126                254            325           53
2020-01-03           137                254            325           63

            Premium_Capacity  Premium_Rate
Date
2020-01-01               100           575
2020-01-02               100           575
2020-01-03               100           575


        Standard_OCC  Standard_Capacity  Standard_Rate  Premium_OCC  \
count    1461.000000             1461.0         1461.0  1461.000000
mean      135.646133              254.0          325.0    56.141684
std        50.216654                0.0            0.0    29.349613
min        38.000000              254.0          325.0     0.000000
25%        88.000000              254.0          325.0    31.000000
```

18

```
50%           136.000000              254.0       325.0     57.000000
75%           183.000000              254.0       325.0     83.000000
max           254.000000              254.0       325.0    100.000000


              Premium_Capacity   Premium_Rate
count                   1461.0         1461.0
mean                     100.0          575.0
std                        0.0            0.0
min                      100.0          575.0
25%                      100.0          575.0
50%                      100.0          575.0
75%                      100.0          575.0
max                      100.0          575.0
```

### A.1.2 Autocorrelation, Seasonality, Stationarity

- Determine autocorrelation, seasonality, stationarity .. lag plot, ACF plot, ADF test, Differencing
- Decomposition



Figure A.1: Occupancy for Each Room Category

- Shows definite annual seasonality with peak high and low seasons
- Also some infrequent spikes in bookings
- Possibly a small upward tend over time
- Premium rooms hit max and zero bookings several times …

- Both categories of room show definite autocorrelation
- Premium rooms bunched up at max value and autocorrelation may be slight less strong
- Some outliers when rooms are fully booked
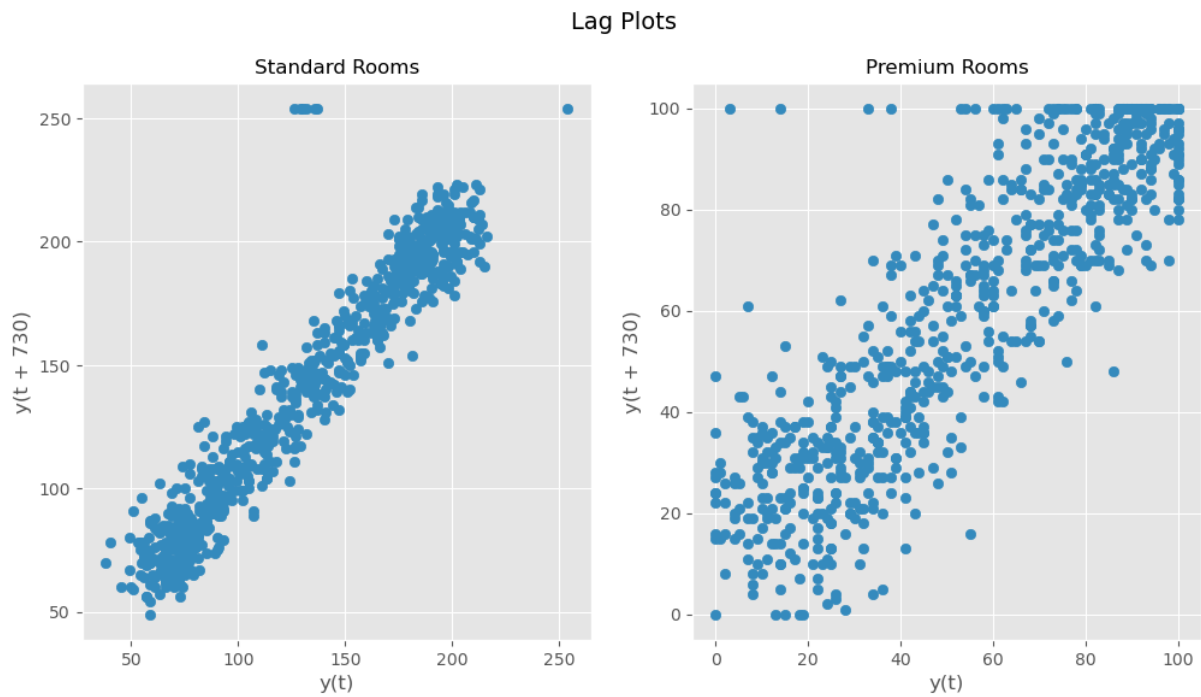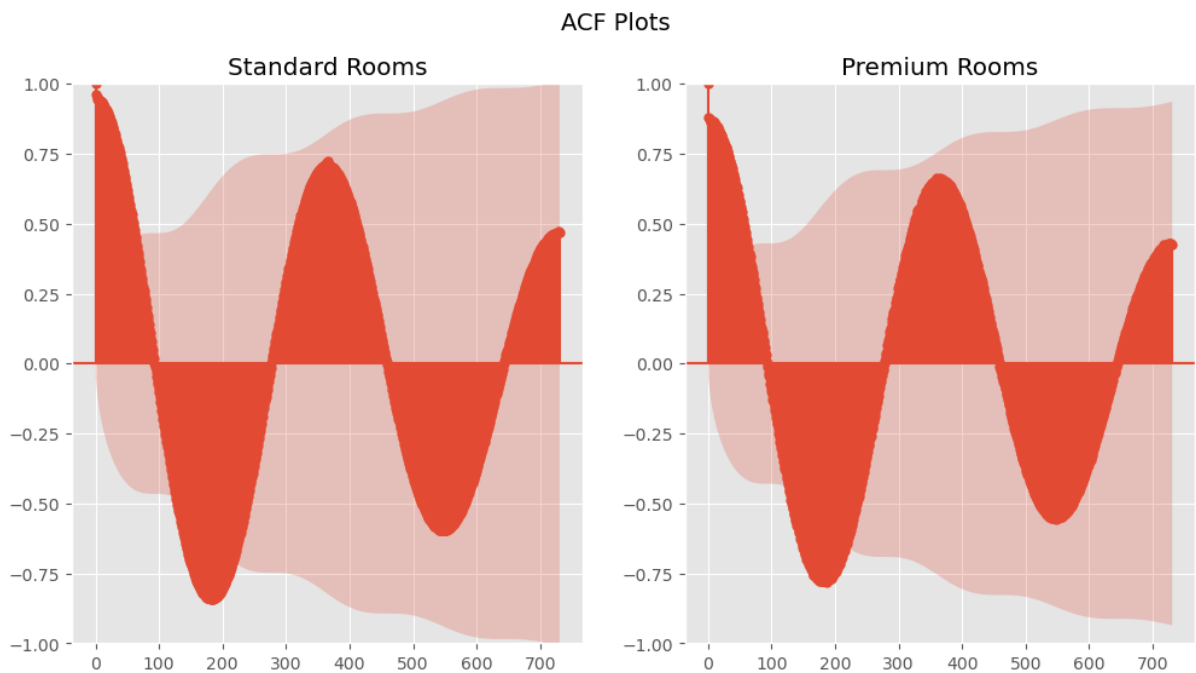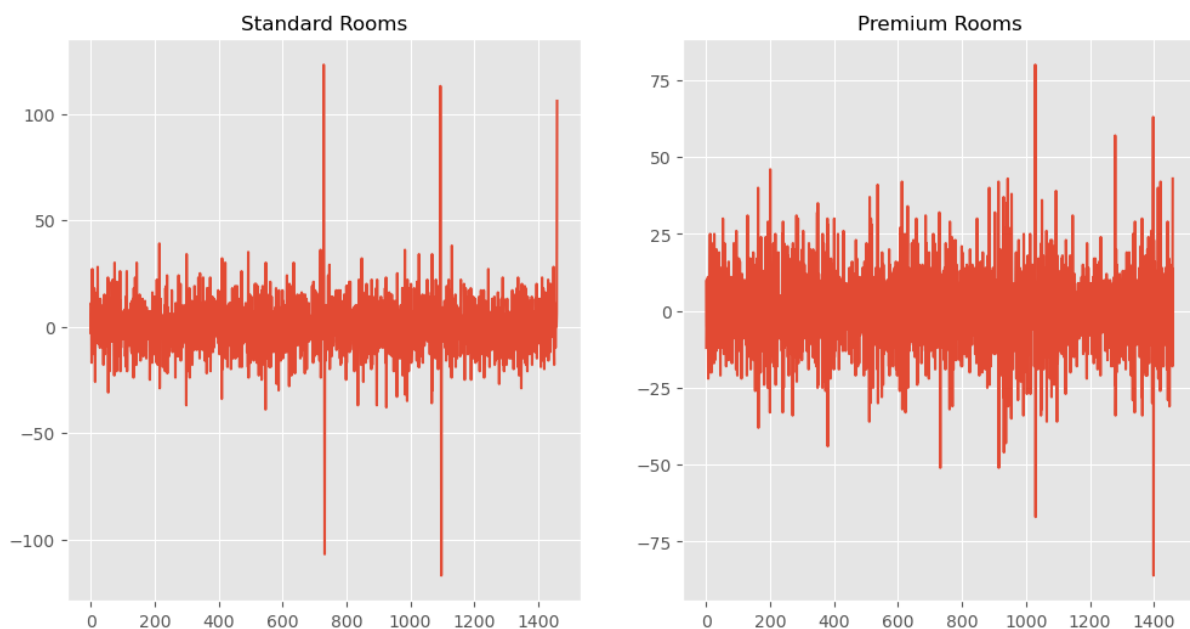
19

Figure A.2: Lag Plots - Test



- Both exhibit strong autocorrelation that diminishes slowly after approximately 250 days
- A positive trend is suggested by the slowly diminishing autocorrelation
- Multiple peaks at 350 days indicates annual seasonality
- ?? carry out PACF to determine the autoregressive order .... does indicate that it is autoregressive

ADF Test for Standard Rooms

```
ADF Statistic: -2.362457191578427
p-value: 0.15261408046089647
Critical Value 1%: -3.434908816804013
Critical Value 5%: -2.863553406963303
Critical Value 10%: -2.5678419239852994
Conclusion: Non-Stationary

ADF Test for Premium Rooms
ADF Statistic: -2.1359470749871265
p-value: 0.2303078418058474
Critical Value 1%: -3.434911997169608
Critical Value 5%: -2.863554810504947
Critical Value 10%: -2.567842671398422
Conclusion: Non-Stationary
```

### Differenced Time Series



```
ADF Test for Standard Rooms
ADF Statistic: -4.650429185341465
p-value: 0.00010417359492157454
Critical Value 1%: -3.4349151819757466
Critical Value 5%: -2.863556216004778
Critical Value 10%: -2.5678434198545568
Conclusion: Stationary

ADF Test for Premium Rooms
ADF Statistic: -6.091009300062941
p-value: 1.0365104637060615e-07
Critical Value 1%: -3.4349151819757466
Critical Value 5%: -2.863556216004778
Critical Value 10%: -2.5678434198545568
```

```
Conclusion: Stationary
```

- Confirms that both time series are non-stationary
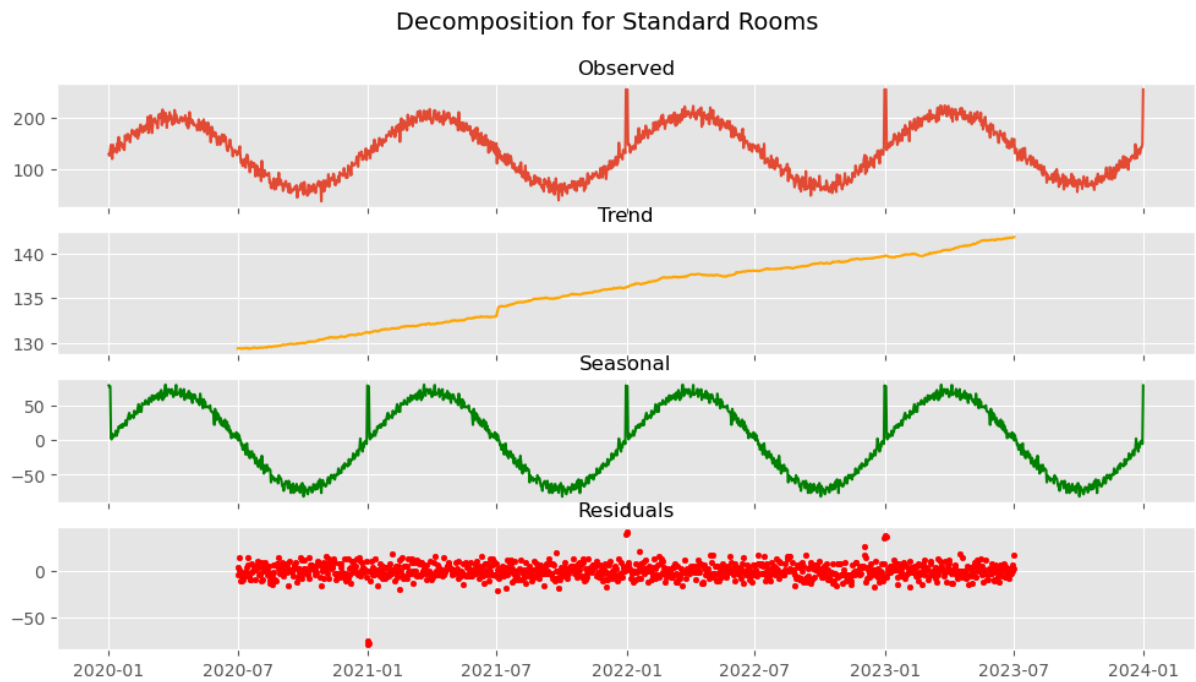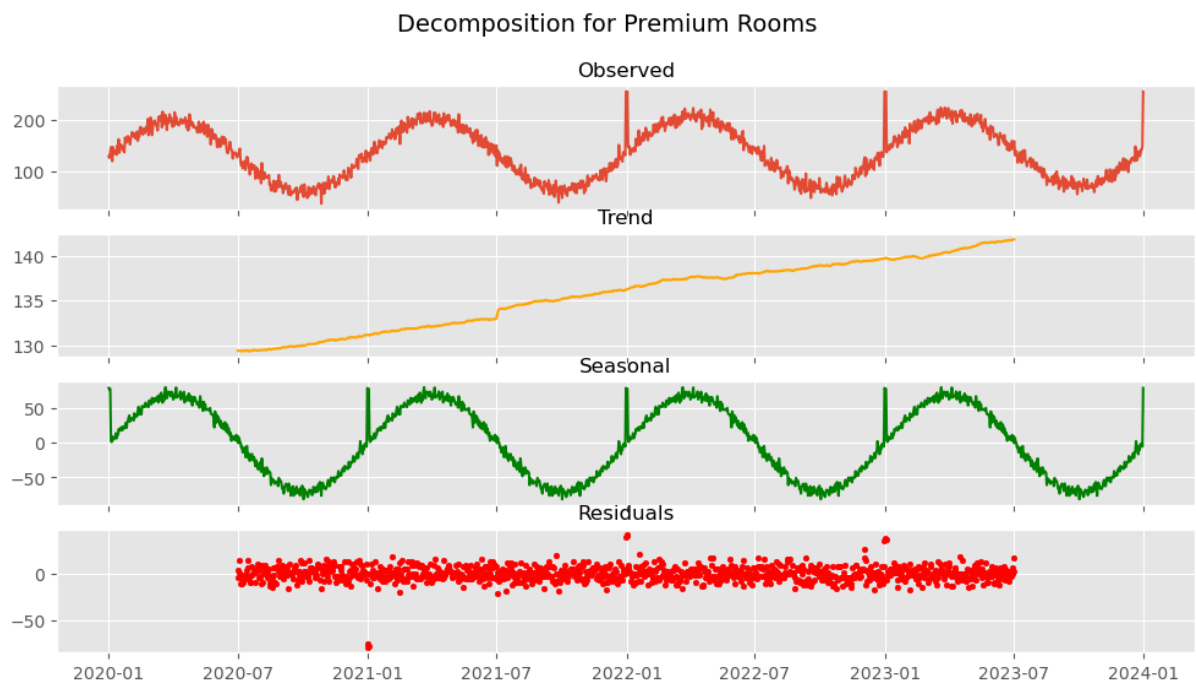


Figure A.3: Occupancy Time Series Decomposition



- Both have a small positive trend with room occupancy increasing 5 to 105% pa
- Confirms both time series are seasonal, with annual peaks and troughs
- On top of the annual seasonality, there are regular spikes leading to 100% occupancy
- Close clustering of residuals with some outliers that correspond to the seasonal spikes

## A.2 Ordinary Least Squares (OLS) Linear Regression

- Unlikely to be a good model for forecasting given the strong seasonality, but examine to confirm

```
Durbin-Watson statistic: 0.07265854622866856
Durbin-Watson statistic: 0.24391749000237792
```



Figure A.4: OLS Linear Regression

- The two fitted lines do not capture any seasonality
- Also show a downward trend line, which is not consistent with the decomposition trend line
- The Durban Watson statistics for both time series are less than 1.5 which confirms evidence of autocorrelation

## A.3 SARIMA Model Creation & Forecasting

- The strong autocorrelation for the room occupancy time series, suggest an autoregressive model such as ARMA or ARIMA or SARIMA
- The occupancy history is non-stationary and there is an upward trend, so the integrated component of ARIMA will automatically perform the differencing needed to transform the data into stationary data. So preferable to ARMA
- The strong seasonality suggests SARIMA would be most appropriate as it can process seasonal patterns. So preferable to ARIMA

23

SARIMA - nb to be clear SARIMA stands for Seasonal Autoregressive Integrated Moving Average - S: seasonal component, here handle the annual occupancy seasonality - AR: autoregressive component .... p - I: integrated, here handle the positive historical trend through differencing to make it stationary .... d - MA: moving average component ...... q - SARIMA(p,d,q)(P,D,Q),m - p,d,q non-seasonal - P,D,Q seasonal - m the seasonality period, length of the seasonal cycle - !! But this is daily, with an annual cycle so would be m=365 which is too large ??

Residuals .... ?? - Residuals plot - ACF plot of residuals - Durbin-Watson of residuals .... ?? be close to 2

Approach - Confirm/Assess the autoregressive order (AR, p) and moving average order (MA, q). Using PACF plot and ACF plot respectively [repeat autocorrelation findings from previously?] - Use the SARIMA(p,d,q)(P,D,Q,m) model - Identify the best parameters using Auto Arima and using AIC (Akaike Information Criterion) to compare

- ?? daily 365d is too much

- ?? First evaluate using residuals ....

- Split data and use 3 years history to the train the model, then 1 year to evaluate its accuracy

- ??? Use accuracy measures MAE, RMSE, MAPE ....

Forecasting With model - Make month forecast and demonstrate its use in the revenue/occupancy spreadsheet

### A.3.1 Data Load & SARIMA Model Factors

- Load time series data
- Assess the autoregressive order (AR, p) - Using PACF plot
- Assess moving average order (MA, q) - Using ACF plot

PACF plot suggests AR order, p of approximately 3

ACF plot suggests MA order q of approximately 60

### A.3.2 Create the SARIMA Model

- Save the training dataset to be used to train a SARIMA model, in a separate notebook
- Load the model in here for forecasting

### A.3.3 Forecast Using the SARIMA Model

# B Customer Satisfaction - Jupyter Notebook Output

## B.1 Data Load & Characteristics

- Load time series data and look at its characteristics

## B.2 XGBoost Model Creation & Forecasting

### B.2.1 Data Load

### B.2.2 Create the XGBoost Model

### B.2.3 Forecast Using the XGBoost Model

# C SARIMA Models Creation + Auto ARIMA - Jupyter Notebook Output

## C.1 SARIMA Model Create & Evaluate - Manual

Create the fitted model and save in model_sarima_manual

```
SARIMA Model Manual Creation using order: (3, 1, 1) and seasonal order: (1, 1, 1, 7)
RUNNING THE L-BFGS-B CODE

           * * *

Machine precision = 2.220D-16
 N =               7     M =                10

At X0          0 variables are exactly at the bounds

At iterate     0    f=  4.09053D+00    |proj g|=  3.35018D-01


At iterate     5    f=  3.94459D+00    |proj g|=  1.01366D-01

At iterate    10    f=  3.88055D+00    |proj g|=  8.90536D-02

At iterate    15    f=  3.85460D+00    |proj g|=  5.44761D-02

At iterate    20    f=  3.85114D+00    |proj g|=  5.44608D-04

At iterate    25    f=  3.85024D+00    |proj g|=  2.92939D-02

At iterate    30    f=  3.84914D+00    |proj g|=  1.15114D-03

At iterate    35    f=  3.84913D+00    |proj g|=  1.15783D-04

           * * *

Tit   = total number of iterations
Tnf   = total number of function evaluations
Tnint = total number of segments explored during Cauchy searches
Skip  = number of BFGS updates skipped
Nact  = number of active bounds at final generalized Cauchy point
Projg = norm of the final projected gradient
```

```
F       = final function value

           * * *

   N    Tit     Tnf  Tnint  Skip  Nact    Projg         F
   7     37      45      1     0     0   1.335D-06   3.849D+00
  F =    3.8491339434652918


CONVERGENCE: NORM_OF_PROJECTED_GRADIENT_<=_PGTOL
Residuals Analysis for Model:
                              SARIMAX Results
==============================================================================
Dep. Variable:                Standard_OCC   No. Observations:          1096
Model:          SARIMAX(3, 1, 1)x(1, 1, 1, 7)   Log Likelihood      -4218.651
Date:                     Wed, 08 Jan 2025   AIC                     8451.302
Time:                             17:34:39   BIC                     8486.246
Sample:                         01-01-2020   HQIC                    8464.529
                              - 12-31-2022
Covariance Type:                       opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1          0.1014      0.034      3.014      0.003       0.035       0.167
ar.L2          0.0407      0.027      1.488      0.137      -0.013       0.094
ar.L3         -0.0951      0.023     -4.168      0.000      -0.140      -0.050
ma.L1         -0.7300      0.033    -22.424      0.000      -0.794      -0.666
ar.S.L7       -0.0960      0.032     -3.027      0.002      -0.158      -0.034
ma.S.L7       -0.8760      0.022    -39.008      0.000      -0.920      -0.832
sigma2       135.0570      2.019     66.890      0.000     131.100     139.014
==============================================================================
Ljung-Box (L1) (Q):                 0.02   Jarque-Bera (JB):        23149.11
Prob(Q):                            0.88   Prob(JB):                    0.00
Heteroskedasticity (H):             1.69   Skew:                        2.08
Prob(H) (two-sided):                0.00   Kurtosis:                   25.21
==============================================================================

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

Model Evaluation on Testing Data



Testing Data Comparison

```
RMSE is: 235.92
MAPE is: 235.38
Stored 'model_sarima_manual' (SARIMAXResultsWrapper)
```

## C.2 SARIMA Model Create & Evaluate - Auto ARIMA - First

Save in model_sarima_auto

```
Performing stepwise search to minimize aic
 ARIMA(2,1,2)(1,1,1)[7]             : AIC=inf, Time=2.74 sec
 ARIMA(0,1,0)(0,1,0)[7]             : AIC=9473.205, Time=0.04 sec
 ARIMA(1,1,0)(1,1,0)[7]             : AIC=8920.643, Time=0.20 sec
 ARIMA(0,1,1)(0,1,1)[7]             : AIC=inf, Time=0.47 sec
 ARIMA(1,1,0)(0,1,0)[7]             : AIC=9250.699, Time=0.06 sec
 ARIMA(1,1,0)(2,1,0)[7]             : AIC=8802.692, Time=0.51 sec
 ARIMA(1,1,0)(2,1,1)[7]             : AIC=inf, Time=1.85 sec
 ARIMA(1,1,0)(1,1,1)[7]             : AIC=inf, Time=0.70 sec
 ARIMA(0,1,0)(2,1,0)[7]             : AIC=9022.330, Time=0.16 sec
 ARIMA(2,1,0)(2,1,0)[7]             : AIC=8765.968, Time=0.52 sec
 ARIMA(2,1,0)(1,1,0)[7]             : AIC=8880.910, Time=0.26 sec
 ARIMA(2,1,0)(2,1,1)[7]             : AIC=inf, Time=1.71 sec
 ARIMA(2,1,0)(1,1,1)[7]             : AIC=inf, Time=1.44 sec
 ARIMA(3,1,0)(2,1,0)[7]             : AIC=8730.320, Time=0.59 sec
 ARIMA(3,1,0)(1,1,0)[7]             : AIC=8844.538, Time=0.42 sec
 ARIMA(3,1,0)(2,1,1)[7]             : AIC=inf, Time=3.49 sec
 ARIMA(3,1,0)(1,1,1)[7]             : AIC=inf, Time=2.35 sec
 ARIMA(4,1,0)(2,1,0)[7]             : AIC=8693.506, Time=1.03 sec
 ARIMA(4,1,0)(1,1,0)[7]             : AIC=8807.355, Time=0.55 sec
 ARIMA(4,1,0)(2,1,1)[7]             : AIC=inf, Time=4.86 sec
 ARIMA(4,1,0)(1,1,1)[7]             : AIC=inf, Time=2.99 sec
 ARIMA(5,1,0)(2,1,0)[7]             : AIC=8676.272, Time=1.02 sec
 ARIMA(5,1,0)(1,1,0)[7]             : AIC=8794.927, Time=0.43 sec
 ARIMA(5,1,0)(2,1,1)[7]             : AIC=inf, Time=4.29 sec
 ARIMA(5,1,0)(1,1,1)[7]             : AIC=inf, Time=2.58 sec
 ARIMA(5,1,1)(2,1,0)[7]             : AIC=8616.898, Time=2.04 sec
 ARIMA(5,1,1)(1,1,0)[7]             : AIC=8718.876, Time=1.25 sec
 ARIMA(5,1,1)(2,1,1)[7]             : AIC=inf, Time=5.86 sec
 ARIMA(5,1,1)(1,1,1)[7]             : AIC=8451.943, Time=3.78 sec
 ARIMA(5,1,1)(0,1,1)[7]             : AIC=inf, Time=1.37 sec
 ARIMA(5,1,1)(1,1,2)[7]             : AIC=8453.944, Time=4.84 sec
 ARIMA(5,1,1)(0,1,0)[7]             : AIC=9024.409, Time=0.51 sec
 ARIMA(5,1,1)(0,1,2)[7]             : AIC=8452.035, Time=3.12 sec
 ARIMA(5,1,1)(2,1,2)[7]             : AIC=8455.726, Time=8.93 sec
 ARIMA(4,1,1)(1,1,1)[7]             : AIC=8449.945, Time=3.40 sec
 ARIMA(4,1,1)(0,1,1)[7]             : AIC=inf, Time=1.68 sec
 ARIMA(4,1,1)(1,1,0)[7]             : AIC=8717.998, Time=1.13 sec
 ARIMA(4,1,1)(2,1,1)[7]             : AIC=inf, Time=5.53 sec
 ARIMA(4,1,1)(1,1,2)[7]             : AIC=8451.941, Time=4.48 sec
 ARIMA(4,1,1)(0,1,0)[7]             : AIC=9024.170, Time=0.56 sec
 ARIMA(4,1,1)(0,1,2)[7]             : AIC=8450.036, Time=3.00 sec
 ARIMA(4,1,1)(2,1,0)[7]             : AIC=8615.480, Time=1.51 sec
 ARIMA(4,1,1)(2,1,2)[7]             : AIC=8453.741, Time=8.03 sec
 ARIMA(3,1,1)(1,1,1)[7]             : AIC=8451.302, Time=2.72 sec
 ARIMA(4,1,2)(1,1,1)[7]             : AIC=8446.786, Time=3.78 sec
 ARIMA(4,1,2)(0,1,1)[7]             : AIC=inf, Time=3.31 sec
 ARIMA(4,1,2)(1,1,0)[7]             : AIC=8718.912, Time=2.32 sec
 ARIMA(4,1,2)(2,1,1)[7]             : AIC=inf, Time=8.99 sec
 ARIMA(4,1,2)(1,1,2)[7]             : AIC=8444.931, Time=6.07 sec
```

```
ARIMA(4,1,2)(0,1,2)[7]              : AIC=8446.244, Time=6.68 sec
ARIMA(4,1,2)(2,1,2)[7]              : AIC=8454.879, Time=8.33 sec
ARIMA(3,1,2)(1,1,2)[7]              : AIC=8519.228, Time=6.78 sec
ARIMA(5,1,2)(1,1,2)[7]              : AIC=8455.934, Time=4.79 sec
ARIMA(4,1,3)(1,1,2)[7]              : AIC=inf, Time=7.79 sec
ARIMA(3,1,1)(1,1,2)[7]              : AIC=8453.293, Time=4.84 sec
ARIMA(3,1,3)(1,1,2)[7]              : AIC=8456.475, Time=6.33 sec
ARIMA(5,1,3)(1,1,2)[7]              : AIC=inf, Time=7.94 sec
ARIMA(4,1,2)(1,1,2)[7] intercept   : AIC=8453.667, Time=7.78 sec

Best model:  ARIMA(4,1,2)(1,1,2)[7]
Total fit time: 184.827 seconds
SARIMA Model Manual Creation using order: (4, 1, 2) and seasonal order: (1, 1, 2, 7)
RUNNING THE L-BFGS-B CODE


          * * *

Machine precision = 2.220D-16
 N =              10      M =              10

At X0          0 variables are exactly at the bounds

At iterate     0    f=  4.03105D+00    |proj g|=  3.18759D-01


At iterate     5    f=  3.91558D+00    |proj g|=  2.75819D-02

At iterate    10    f=  3.85716D+00    |proj g|=  8.66397D-02

At iterate    15    f=  3.84864D+00    |proj g|=  1.49384D-03

At iterate    20    f=  3.84831D+00    |proj g|=  1.25557D-03

At iterate    25    f=  3.84801D+00    |proj g|=  8.50538D-03

At iterate    30    f=  3.84767D+00    |proj g|=  1.84933D-03

At iterate    35    f=  3.84759D+00    |proj g|=  2.58246D-03

At iterate    40    f=  3.84749D+00    |proj g|=  3.01463D-03

At iterate    45    f=  3.84706D+00    |proj g|=  2.46234D-02


At iterate    50    f=  3.84349D+00    |proj g|=  1.32273D-02

          * * *

Tit   = total number of iterations
```

```
Tnf   = total number of function evaluations
Tnint = total number of segments explored during Cauchy searches
Skip  = number of BFGS updates skipped
Nact  = number of active bounds at final generalized Cauchy point
Projg = norm of the final projected gradient
F     = final function value

          * * *

   N    Tit    Tnf  Tnint  Skip  Nact    Projg        F
  10     50     59     1     0     0   1.323D-02   3.843D+00
   F =   3.8434904043823033

STOP: TOTAL NO. of ITERATIONS REACHED LIMIT
Residuals Analysis for Model:
```

                               SARIMAX Results
==========================================================================================
| Dep. Variable: | Standard_OCC | No. Observations: | 1096 |
|---|---|---|---|
| Model: | SARIMAX(4, 1, 2)x(1, 1, 2, 7) | Log Likelihood | -4212.465 |
| Date: | Wed, 08 Jan 2025 | AIC | 8444.931 |
| Time: | 17:49:23 | BIC | 8494.852 |
| Sample: | 01-01-2020 | HQIC | 8463.827 |
| | - 12-31-2022 | | |
| Covariance Type: | opg | | |

==========================================================================================

| | coef | std err | z | P>\|z\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| ar.L1 | 1.1291 | 0.039 | 28.783 | 0.000 | 1.052 | 1.206 |
| ar.L2 | -0.0889 | 0.026 | -3.358 | 0.001 | -0.141 | -0.037 |
| ar.L3 | -0.1674 | 0.032 | -5.157 | 0.000 | -0.231 | -0.104 |
| ar.L4 | 0.0423 | 0.025 | 1.707 | 0.088 | -0.006 | 0.091 |
| ma.L1 | -1.8247 | 0.038 | -48.012 | 0.000 | -1.899 | -1.750 |
| ma.L2 | 0.8513 | 0.031 | 27.169 | 0.000 | 0.790 | 0.913 |
| ar.S.L7 | 0.7630 | 0.217 | 3.511 | 0.000 | 0.337 | 1.189 |
| ma.S.L7 | -1.6872 | 0.209 | -8.079 | 0.000 | -2.097 | -1.278 |
| ma.S.L14 | 0.7191 | 0.183 | 3.932 | 0.000 | 0.361 | 1.078 |
| sigma2 | 132.6832 | 1.995 | 66.524 | 0.000 | 128.774 | 136.592 |

==========================================================================================

| Ljung-Box (L1) (Q): | 1.87 | Jarque-Bera (JB): | 25082.63 |
|---|---|---|---|
| Prob(Q): | 0.17 | Prob(JB): | 0.00 |
| Heteroskedasticity (H): | 1.76 | Skew: | 2.29 |
| Prob(H) (two-sided): | 0.00 | Kurtosis: | 26.07 |

==========================================================================================

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

## Model Evaluation on Testing Data



```
RMSE is: 311.34
MAPE is: 307.80
Stored 'model_sarima_auto' (SARIMAXResultsWrapper)
```

'\n# Convert the pmdarima model to a statsmodels SARIMAX model for residual analysis\nmodel_

## C.3 SARIMA Model Create & Evaluate - Auto ARIMA - Second

Save in model_sarima_auto_second

```
Performing stepwise search to minimize aic
 ARIMA(0,1,0)(0,1,0)[7]              : AIC=9473.205, Time=0.04 sec
 ARIMA(1,1,0)(1,1,0)[7]              : AIC=8920.643, Time=0.15 sec
 ARIMA(0,1,1)(0,1,1)[7]              : AIC=inf, Time=0.48 sec
 ARIMA(1,1,0)(0,1,0)[7]              : AIC=9250.699, Time=0.08 sec
 ARIMA(1,1,0)(2,1,0)[7]              : AIC=8802.692, Time=0.32 sec
 ARIMA(1,1,0)(3,1,0)[7]              : AIC=8723.919, Time=0.53 sec
 ARIMA(1,1,0)(3,1,1)[7]              : AIC=inf, Time=2.71 sec
 ARIMA(1,1,0)(2,1,1)[7]              : AIC=inf, Time=2.17 sec
 ARIMA(0,1,0)(3,1,0)[7]              : AIC=8941.142, Time=0.53 sec
 ARIMA(2,1,0)(3,1,0)[7]              : AIC=8694.720, Time=0.74 sec
 ARIMA(2,1,0)(2,1,0)[7]              : AIC=8765.968, Time=0.42 sec
 ARIMA(2,1,0)(3,1,1)[7]              : AIC=inf, Time=4.20 sec
 ARIMA(2,1,0)(2,1,1)[7]              : AIC=inf, Time=1.94 sec
 ARIMA(3,1,0)(3,1,0)[7]              : AIC=8654.898, Time=1.04 sec
 ARIMA(3,1,0)(2,1,0)[7]              : AIC=8730.320, Time=0.66 sec
 ARIMA(3,1,0)(3,1,1)[7]              : AIC=inf, Time=4.86 sec
 ARIMA(3,1,0)(2,1,1)[7]              : AIC=inf, Time=3.19 sec
 ARIMA(3,1,1)(3,1,0)[7]              : AIC=8541.718, Time=2.65 sec
 ARIMA(3,1,1)(2,1,0)[7]              : AIC=8613.852, Time=1.79 sec
 ARIMA(3,1,1)(3,1,1)[7]              : AIC=8454.819, Time=5.42 sec
 ARIMA(3,1,1)(2,1,1)[7]              : AIC=inf, Time=3.95 sec
 ARIMA(3,1,1)(3,1,2)[7]              : AIC=inf, Time=7.51 sec
 ARIMA(3,1,1)(2,1,2)[7]              : AIC=8455.068, Time=7.00 sec
 ARIMA(2,1,1)(3,1,1)[7]              : AIC=8459.822, Time=4.41 sec
 ARIMA(3,1,2)(3,1,1)[7]              : AIC=inf, Time=10.65 sec
 ARIMA(2,1,2)(3,1,1)[7]              : AIC=inf, Time=10.53 sec
 ARIMA(3,1,1)(3,1,1)[7] intercept   : AIC=8456.706, Time=12.51 sec

Best model:  ARIMA(3,1,1)(3,1,1)[7]
Total fit time: 90.579 seconds
```

```
                               SARIMAX Results
==========================================================================================
Dep. Variable:                             y   No. Observations:             1096
Model:             SARIMAX(3, 1, 1)x(3, 1, 1, 7)  Log Likelihood          -4218.409
Date:                       Wed, 08 Jan 2025   AIC                        8454.819
Time:                               17:54:27   BIC                        8499.748
Sample:                           01-01-2020   HQIC                       8471.825
                                - 12-31-2022
Covariance Type:                         opg
==========================================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------------------
ar.L1          0.1063      0.033      3.209      0.001       0.041       0.171
ar.L2          0.0429      0.027      1.582      0.114      -0.010       0.096
```

34

```
ar.L3          -0.0970      0.023     -4.245      0.000     -0.142     -0.052
ma.L1          -0.7384      0.031    -23.567      0.000     -0.800     -0.677
ar.S.L7        -0.1167      0.035     -3.297      0.001     -0.186     -0.047
ar.S.L14       -0.0196      0.045     -0.436      0.663     -0.108      0.069
ar.S.L21       -0.0310      0.042     -0.746      0.456     -0.113      0.051
ma.S.L7        -0.8523      0.030    -28.508      0.000     -0.911     -0.794
sigma2        135.0472      2.162     62.464      0.000    130.810    139.285
===================================================================================
Ljung-Box (L1) (Q):                   0.02   Jarque-Bera (JB):           22690.36
Prob(Q):                              0.88   Prob(JB):                       0.00
Heteroskedasticity (H):               1.67   Skew:                           2.07
Prob(H) (two-sided):                  0.00   Kurtosis:                      24.98
===================================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
SARIMA Model Manual Creation using order: (3, 1, 1) and seasonal order: (3, 1, 1, 7)
RUNNING THE L-BFGS-B CODE


           * * *


Machine precision = 2.220D-16
 N =            9     M =            10


At X0          0 variables are exactly at the bounds


At iterate    0    f=  4.08611D+00    |proj g|=  3.32757D-01



At iterate    5    f=  3.96176D+00    |proj g|=  5.15657D-02


At iterate   10    f=  3.89847D+00    |proj g|=  2.28767D-02


At iterate   15    f=  3.87075D+00    |proj g|=  2.68053D-02


At iterate   20    f=  3.84927D+00    |proj g|=  3.93821D-03


At iterate   25    f=  3.84921D+00    |proj g|=  2.54060D-03


At iterate   30    f=  3.84892D+00    |proj g|=  1.21032D-03


           * * *


Tit   = total number of iterations
Tnf   = total number of function evaluations
Tnint = total number of segments explored during Cauchy searches
Skip  = number of BFGS updates skipped
Nact  = number of active bounds at final generalized Cauchy point
Projg = norm of the final projected gradient
```

```
F     = final function value

          * * *

  N    Tit    Tnf  Tnint  Skip  Nact    Projg        F
  9     34     38      1     0     0   1.154D-05   3.849D+00
  F =   3.8489137053819098


CONVERGENCE: REL_REDUCTION_OF_F_<=_FACTR*EPSMCH
Residuals Analysis for Model:
```

```
==============================================================================
Dep. Variable:                 Standard_OCC   No. Observations:         1096
Model:           SARIMAX(3, 1, 1)x(3, 1, 1, 7)   Log Likelihood       -4218.409
Date:                     Wed, 08 Jan 2025   AIC                    8454.819
Time:                             17:54:34   BIC                    8499.748
Sample:                         01-01-2020   HQIC                   8471.825
                              - 12-31-2022
Covariance Type:                       opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1          0.1063      0.033      3.209      0.001       0.041       0.171
ar.L2          0.0429      0.027      1.582      0.114      -0.010       0.096
ar.L3         -0.0970      0.023     -4.245      0.000      -0.142      -0.052
ma.L1         -0.7384      0.031    -23.567      0.000      -0.800      -0.677
ar.S.L7       -0.1167      0.035     -3.297      0.001      -0.186      -0.047
ar.S.L14      -0.0196      0.045     -0.436      0.663      -0.108       0.069
ar.S.L21      -0.0310      0.042     -0.746      0.456      -0.113       0.051
ma.S.L7       -0.8523      0.030    -28.508      0.000      -0.911      -0.794
sigma2       135.0472      2.162     62.464      0.000     130.810     139.285
===================================================================================
Ljung-Box (L1) (Q):                   0.02   Jarque-Bera (JB):        22690.36
Prob(Q):                              0.88   Prob(JB):                    0.00
Heteroskedasticity (H):               1.67   Skew:                        2.07
Prob(H) (two-sided):                  0.00   Kurtosis:                   24.98
===================================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
```
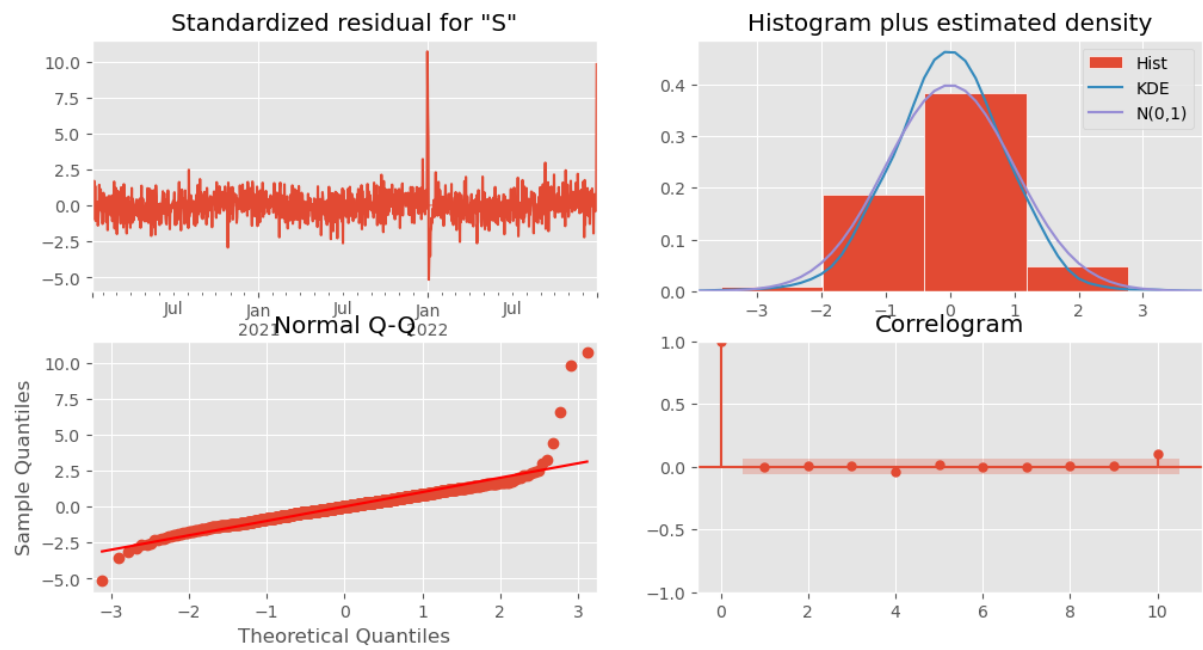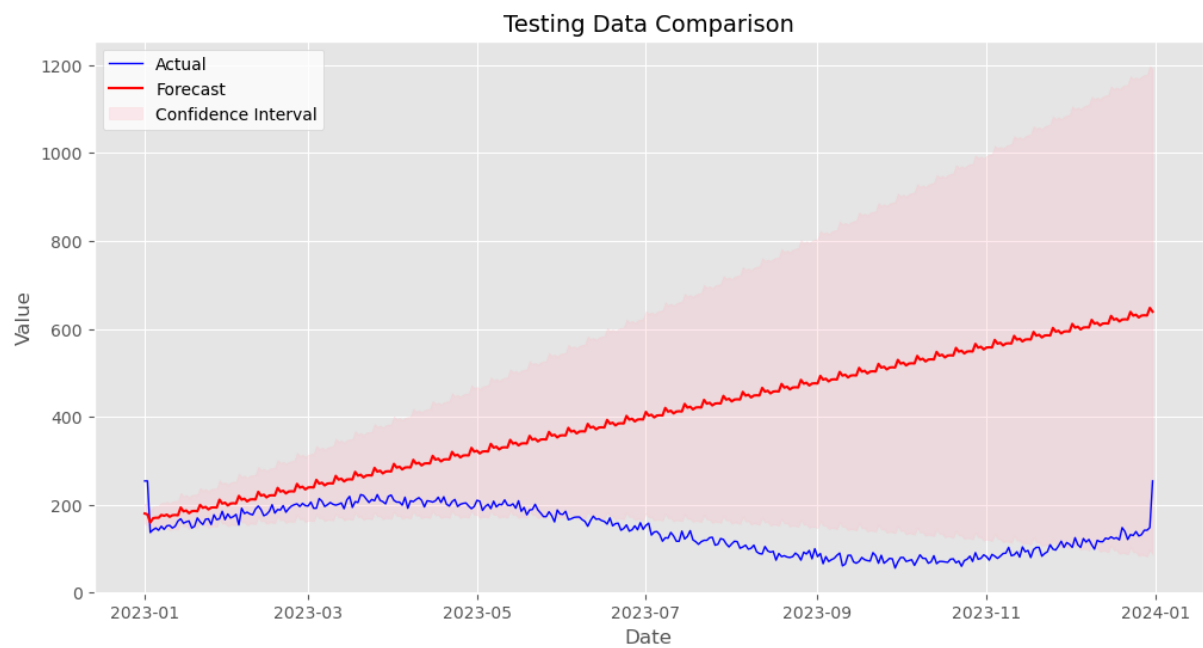
Model Evaluation on Testing Data



Testing Data Comparison

RMSE is: 260.18
MAPE is: 258.63
Stored 'model_sarima_auto_second' (SARIMAXResultsWrapper)