

FORECAST MEMO

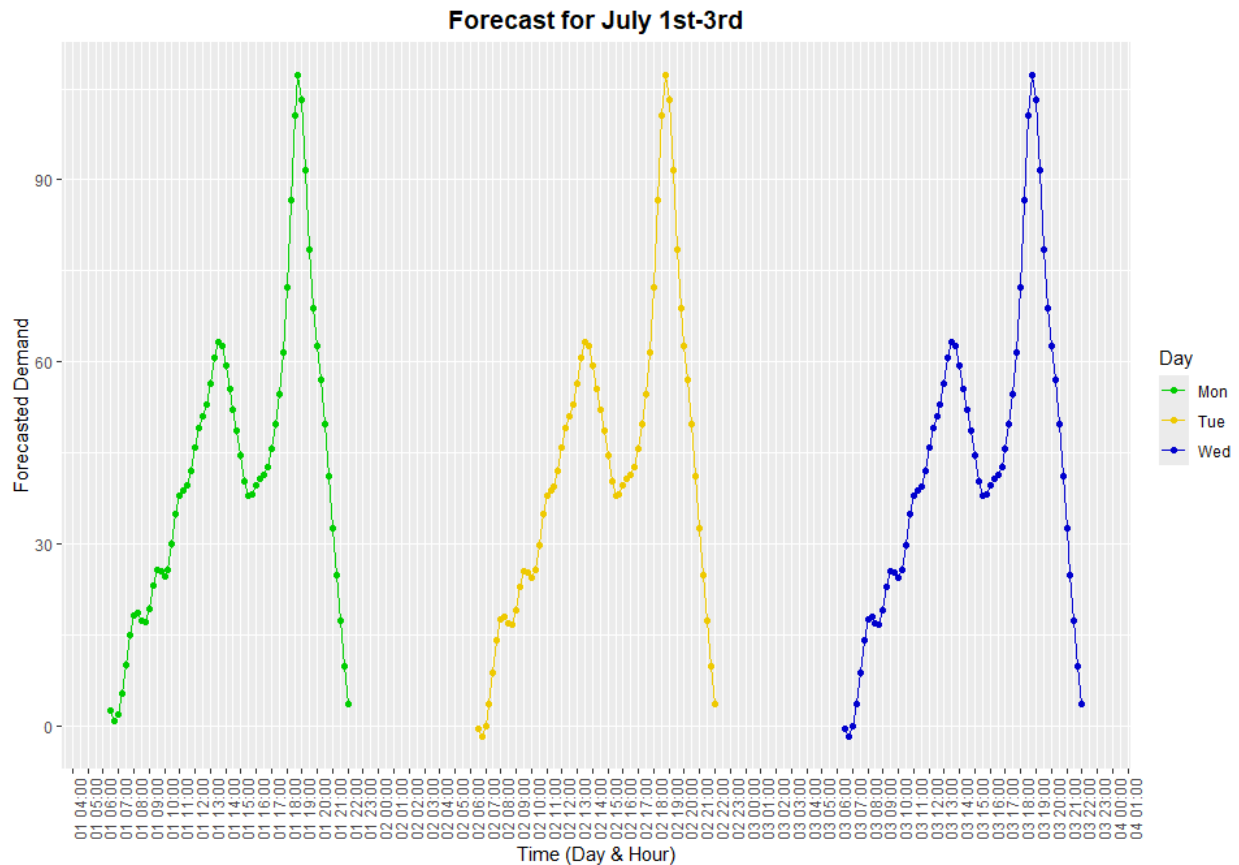
To: COO of VLT CARIOCA

From: Halan Badilla

Date: April 19, 2025

Re: Forecast of Passenger Demands from July 1st-3rd

As requested, I have analyzed the 30 days of data provided of June and created a forecast for the upcoming July 1st to July 3rd dates. Right below this I have attached a visual forecast of the expected passenger demands for July 1st–3rd. The intervals are in 15 minutes as requested, however, for the sake of visual clarity I have shown it hourly, and color coded the days. There are different patterns between the weekends and the weekdays, the three days forecasted will be weekdays and as such it is worth noting that there will be a drastic increase in passenger demands from 11:00-14:00, this is most likely during lunch breaks. And then the evening rush causes another drastic increase, doubling the demand, from 16:00-20:00. These are the most optimal times and demand the most amount of attention.



Once I received the data, the first thing I did was to try to see if there were any patterns or trends in the data provided for June. There I noticed that there was consistently an increase in demand

for both lunch breaks and during an evening rush, most likely the time when people are returning from work. These times specifically were between 11:00-14:00 and 16:00-20:00. While the lunch break displayed a mild increase, the evening rush showed demand increases creeping up to double the peak demand during lunch time. Further analysis also showed that the morning period and the end of the night were shown to be the slowest periods while there also being a moment where the demand lowers from 14:00 to 16:00 but not to the extent as the morning and nighttime.

Regarding the modeling, to make sure I could get an accurate forecast I tested two different models that would help me forecast, one is called TBATS and another which is a mix of STL and ETS. Running both forecasts on the data provided showed me that TBATS gave the better results that were more accurate and varied less from the actual data I tested on, as such, I decided to stick with it and it is what I used in the forecast that I provided.

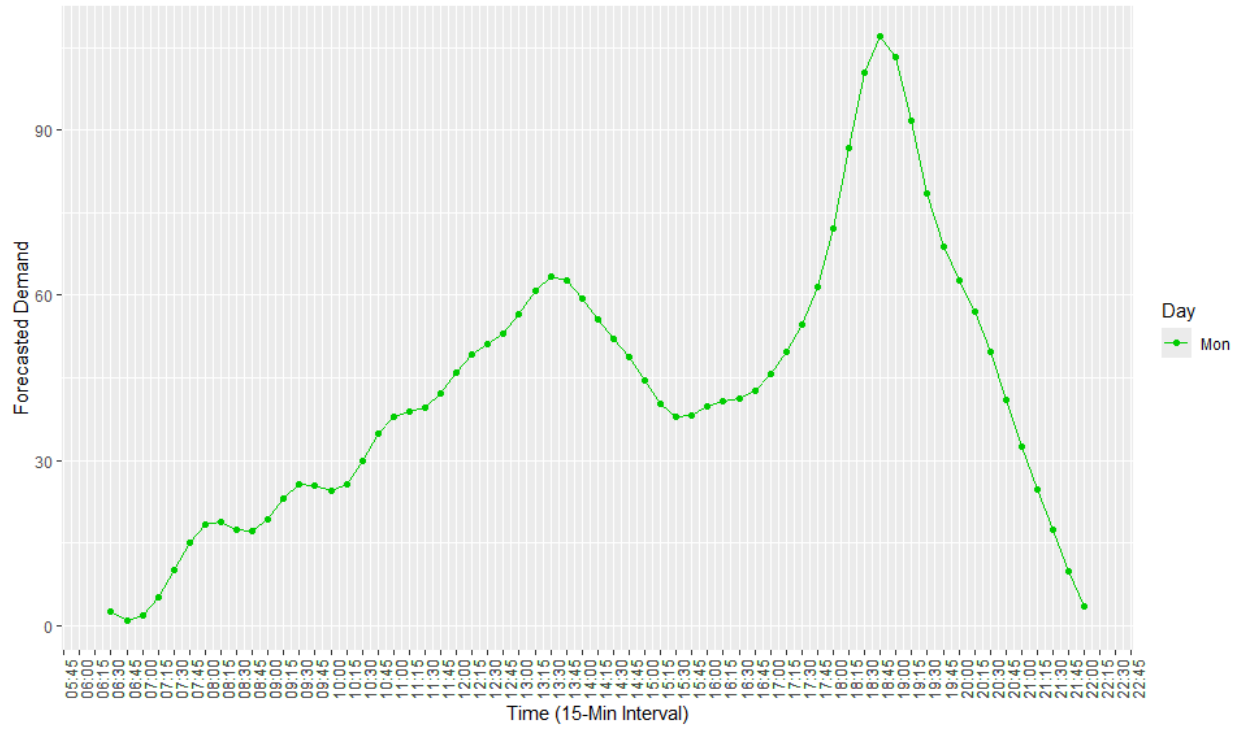
Forecast Summary:

- Consistent Peak Hours between 11:00-14:00 and 16:00-20:00
- Highest Activity during the evening rush of 16:00-20:00
- Lowest Activity during the mornings and late night with there being a sharp decline during the nighttime after the peak finishes.
- Midday Activity shows some demand but not enough to warrant full capacity.

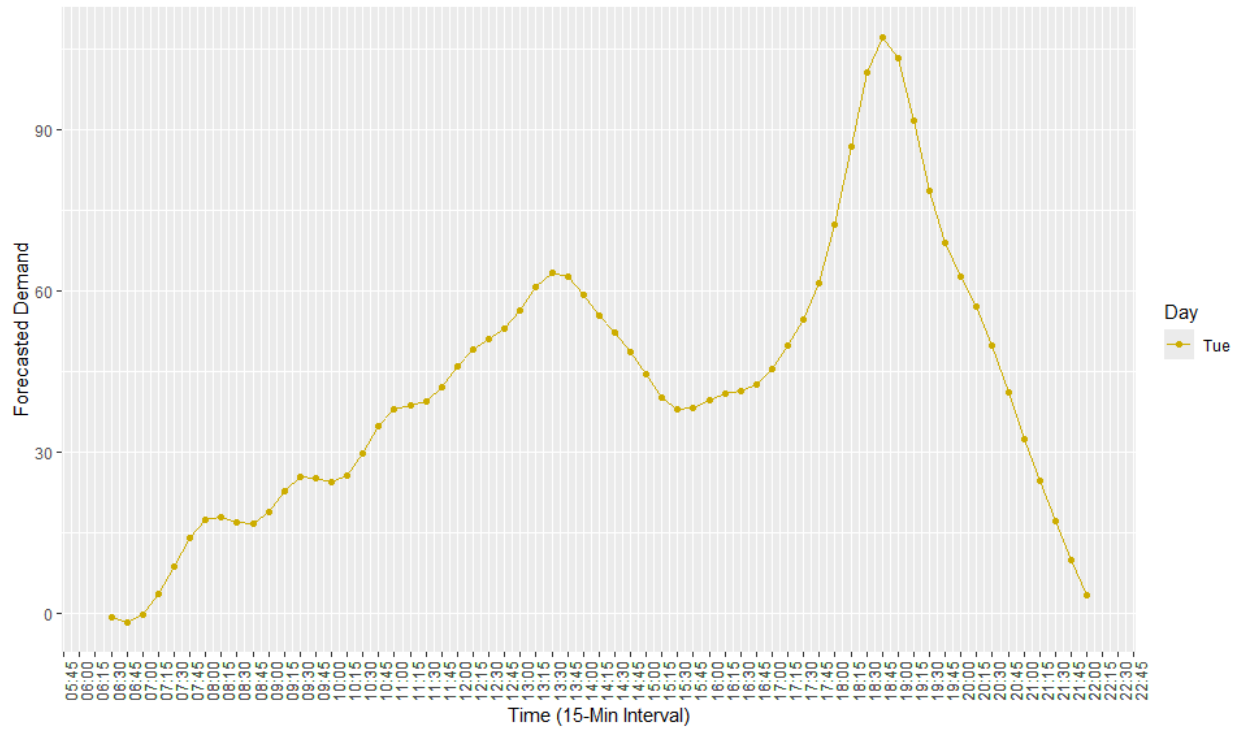
With these forecasts you should be able to keep operations running smoothly and be able to plan staffing hours to both maximize profits, maximize passenger retention and minimize costs.

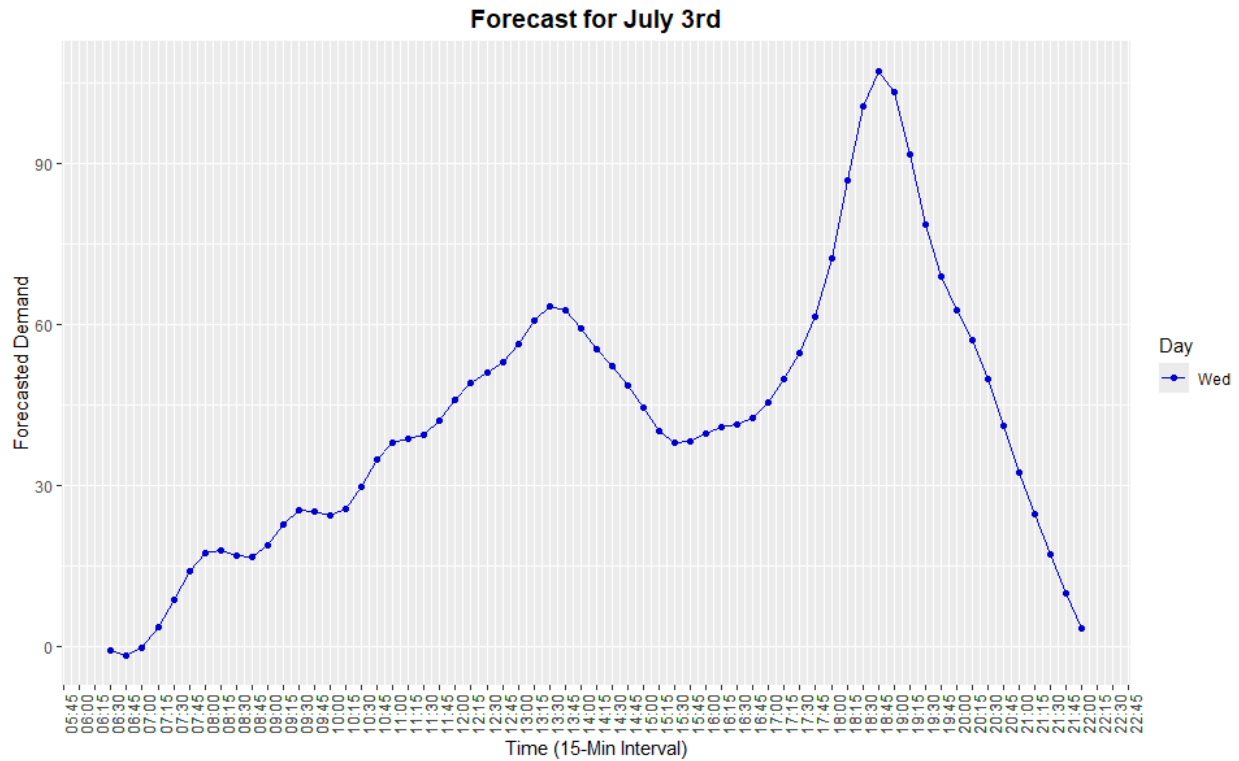
Furthermore, here are the forecasts for individual days, if any further clarity is needed then please let me know.

Forecast for July 1st



Forecast for July 2nd

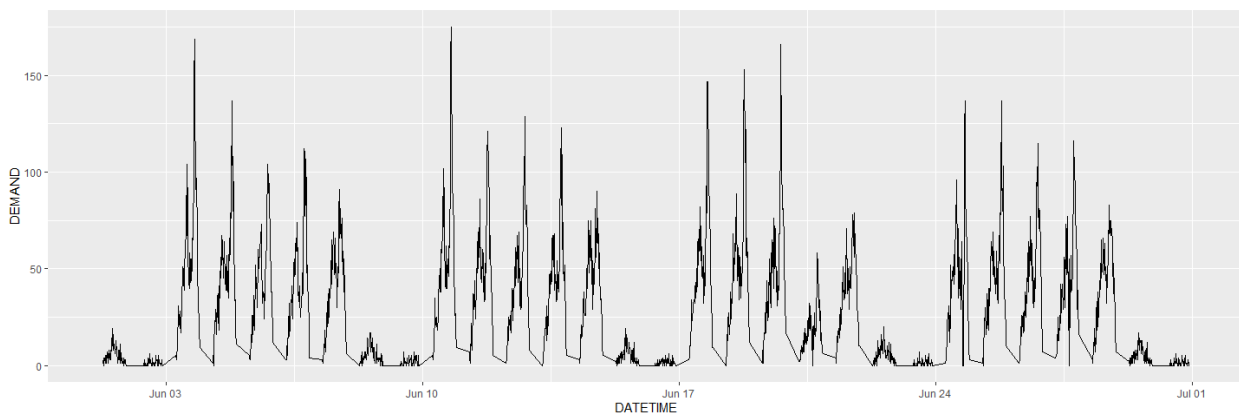




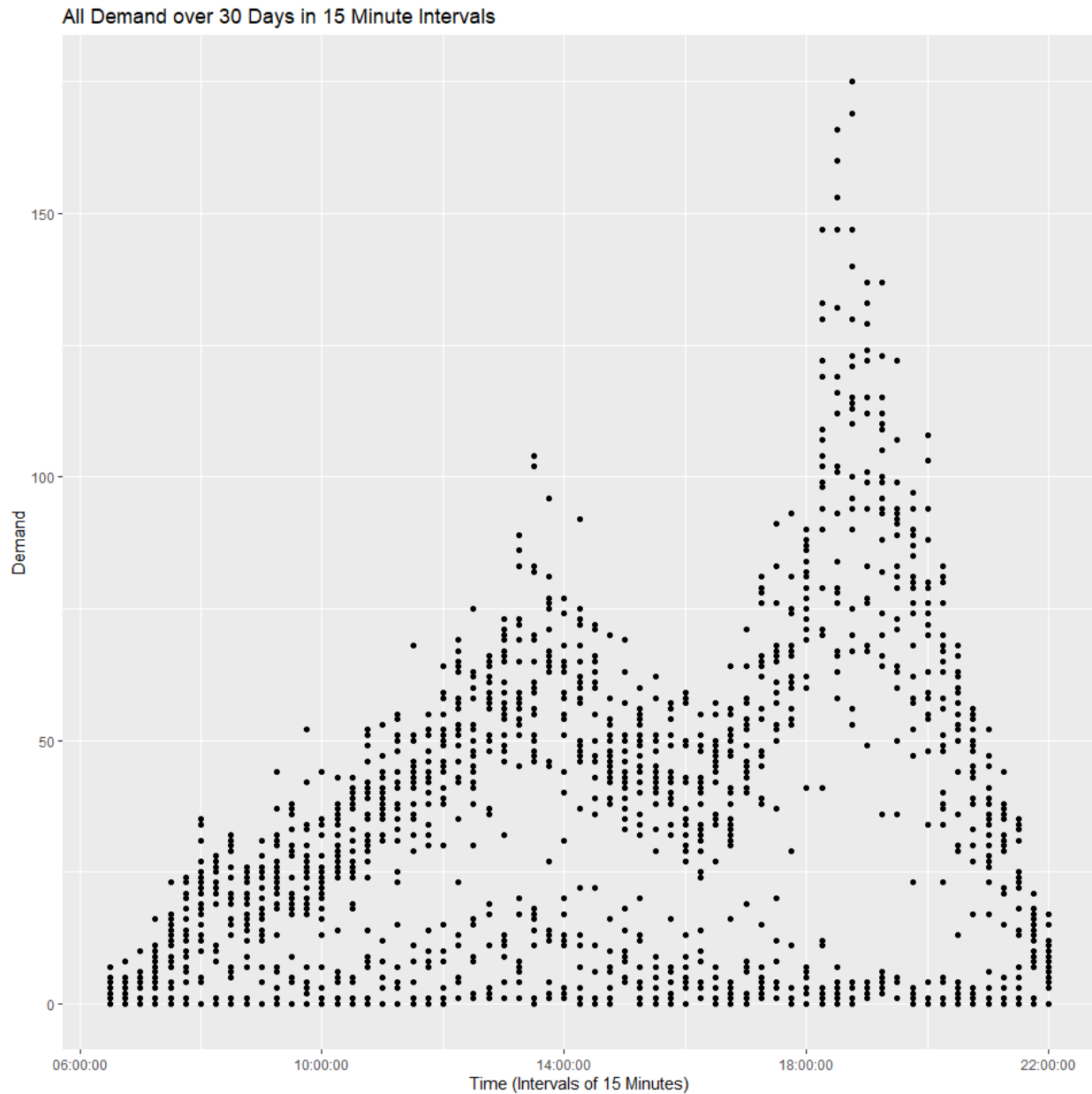
Technical Appendix:

Time Series EDA:

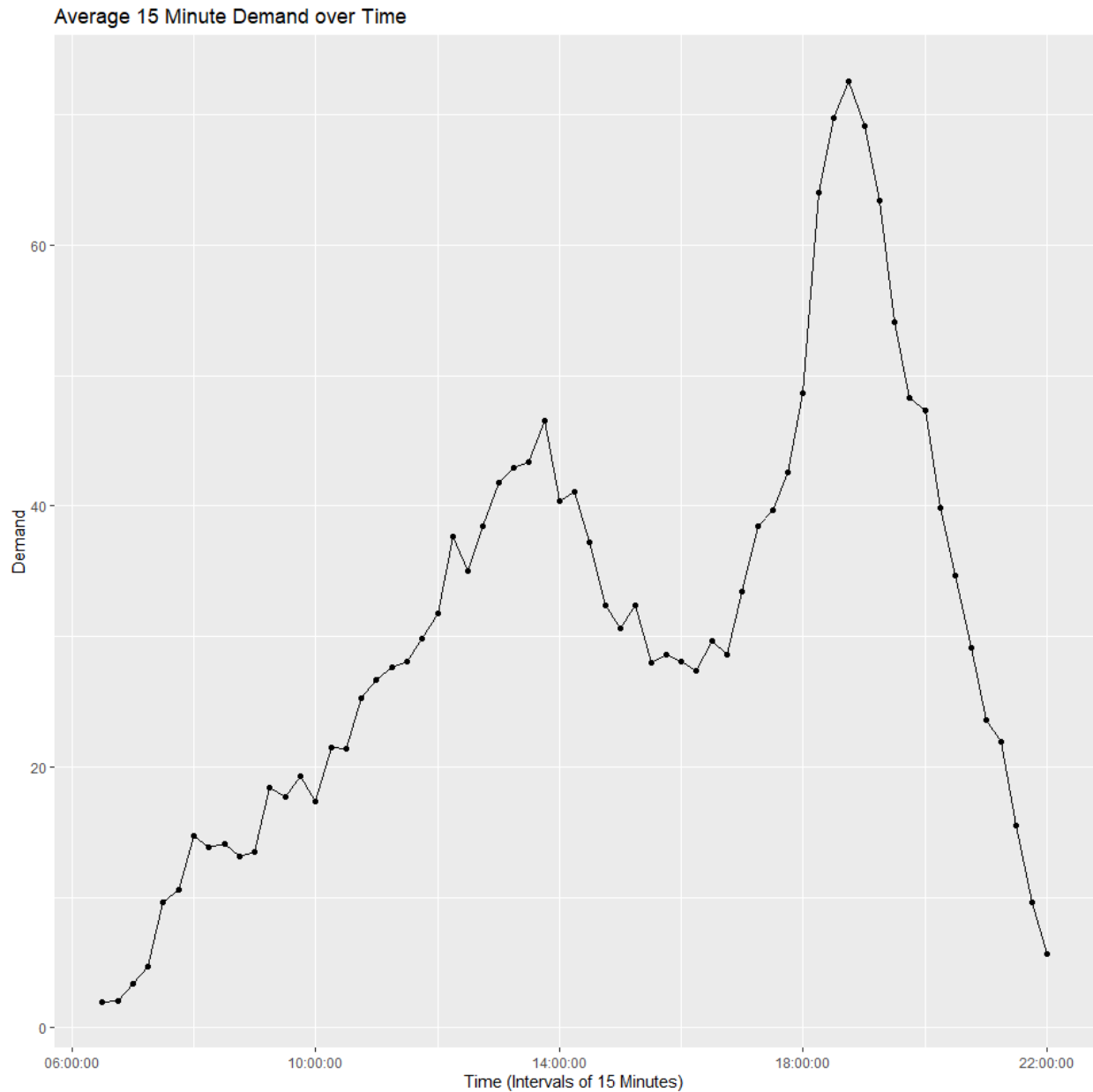
After briefly looking at the values and columns given in the 30 days of June, I ran an EDA showing the demand across all days.



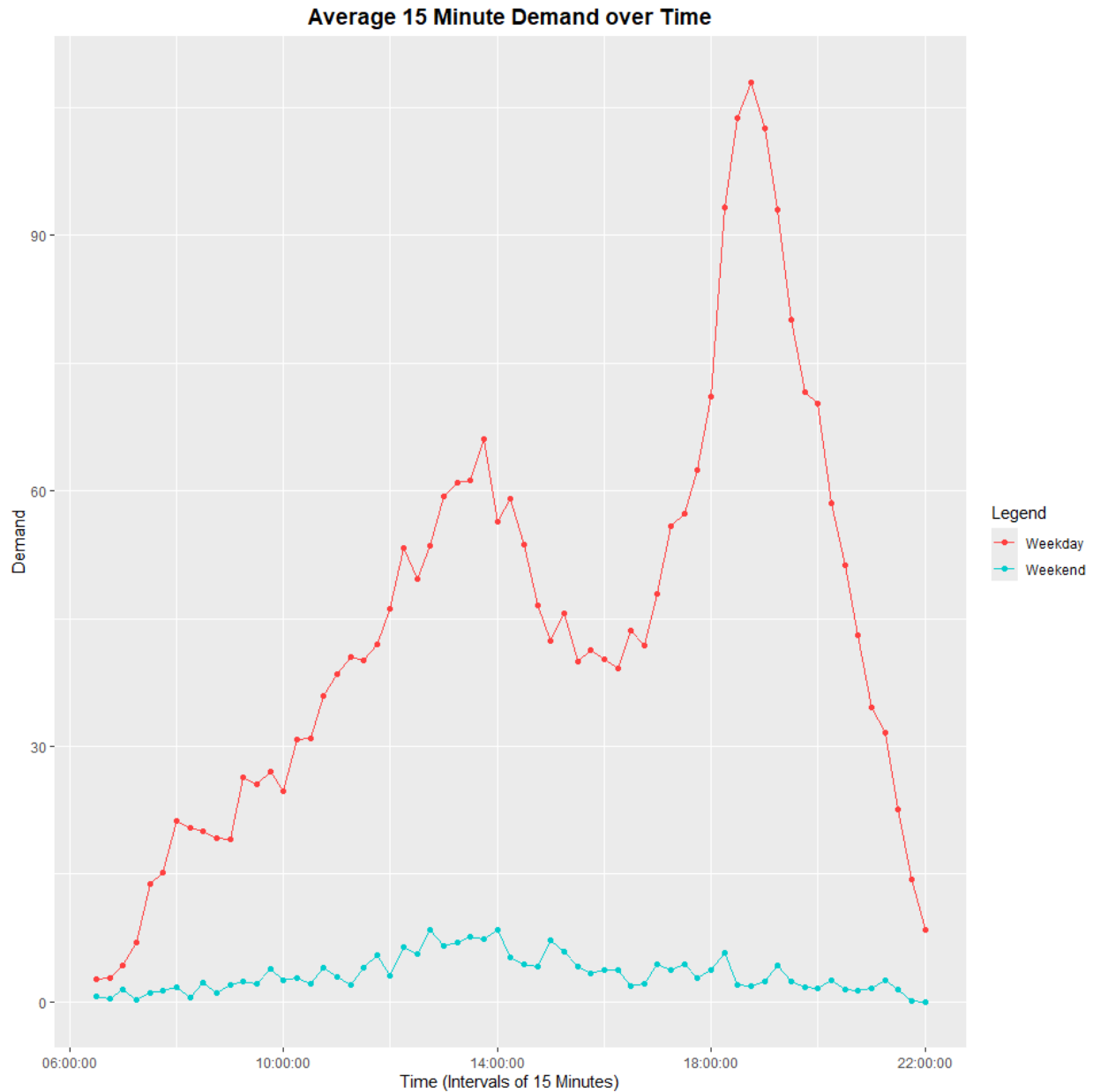
This showed me that there was a very important seasonality across what looks like the weekdays given how many spikes there were, furthermore, it looks like there was also seasonality within the days themselves. To further analyze this, I then did a scatterplot that showcased all the demands relative to the time and not the datetime.



This scatterplot showed me that there indeed existed clear seasonality during the peak hours of 11:00 to 14:00 and another from 16:00 to 20:00. To get a better idea of this I then ran an average by grouping it by time and getting an average of the demand along those times. After plotting this I could see a much clearer uptick in demand during those times.



This gave a much clearer picture but there was another problem, this took into account the weekends. As seen in both the scatter plot and the initial plot, there is a clear difference between demands on the weekends and the weekdays. I then created two columns in the data, “DAY“and “ISWEEKEND “. One displayed what day of the week it was and then another that was a Boolean/Logical showing whether it was in the weekend or not. With this I then made a separate data frame that averaged out the times. With this I could now plot them separately and see it much more clearly.



Viewing this, there is a very clear seasonality to the weekdays and essentially none with the weekends as well a slight negative trend but there is not enough data to be certain.

To delve into the forecast now, I found out what days July 1st-3rd were. These three days are all weekdays and as such I made my forecast using only the weekday data.

Data Modification:

Datetime Column: This was created to be able to plot data across multiple days properly.

Day: To see what day the data was from:

Isweekend: To verify whether the data is during a weekend or not.

Weekday filtering: This was specifically because weekends and weekdays have much different trends, this is especially noticeable in my EDA.

Fit/Forecast/Acc/ETS/Decomp/SeasonAdjust: All of these were created when testing the models. These are simply necessary for the models to function and to compare.

Training/Test data: I split the weekday data into training and test; this was specifically so that I could test out different models to get the most accurate forecast. I made a 75/25 split, given I had a month, I used the first 3 weeks for training and the last week for testing.

Forecasting Methods:

- 1) TBATS Model: This model takes into account trend and seasonality; it does this innately by using state space models with Fourier terms. This model is especially useful when you have data with multiple seasonality, I assumed this would be good as there were two trends; the lunch break and evening rush.
- 2) STL + ETS Model: This model was an alternative after I was not able to run Holt-Winters. This is generally good for anything with a trend and seasonality, in this case I could not tell if the trend existed or not with the data provided, so it was another good alternative.
- 3) Holt Winters: This was not shown or used, the reason being that when I attempted to use it the frequency of 15-minute intervals was far too much for the model to function. Holt-Winters is good for anything with 24 intervals; however, this data was using 63 intervals which is well above what it can reasonably run. As such, an alternative (STL + ETS) was used.

Model Performance:

TBATS Accuracy Statistics							
	ME	RMSE	MAE	MPE	MAPE	MASE	ACF1
Training set	-0.3911711	10.31370	7.614178	-Inf	Inf	0.7768058	0.004041196
Test set	3.2883867	14.52002	9.921251	-Inf	Inf	1.0121756	NA

STL+ETS Accuracy Statistics							
	ME	RMSE	MAE	MPE	MAPE	MASE	ACF1
Training set	-0.01244256	10.17992	7.494271	5.786792	45.92914	0.8805278	0.01696856
Test set	2.91676966	26.75699	20.741650	-Inf	Inf	2.4370080	NA

ME (Mean Error): Average of the forecasts, shows bias.

- STL+ETS showed that both models seemed to over-forecast and when ran on the test data they both under-forecasted, yet STL+ETS did better.
- Both over predicted on average.

RMSE (Root Mean Squared Error): Lower is better, penalizes larger errors more.

- STL+ETS RMSE is significantly higher than TBATS.
- TBAT was ~14.5 units (passenger demand) off on average.
- STL+ETS was ~26.8 units off on average.

MAE (Mean Squared Error): Like RMSE except not rooted. Good for seeing the magnitude of the errors.

- STL+ETS MAE was twice as much as TBATS meaning STL+ETS was more prone to errors and of bigger magnitude.
- TBAT was ~7.5 units off on average.
- STL+ETS was ~20.7 units off on average.

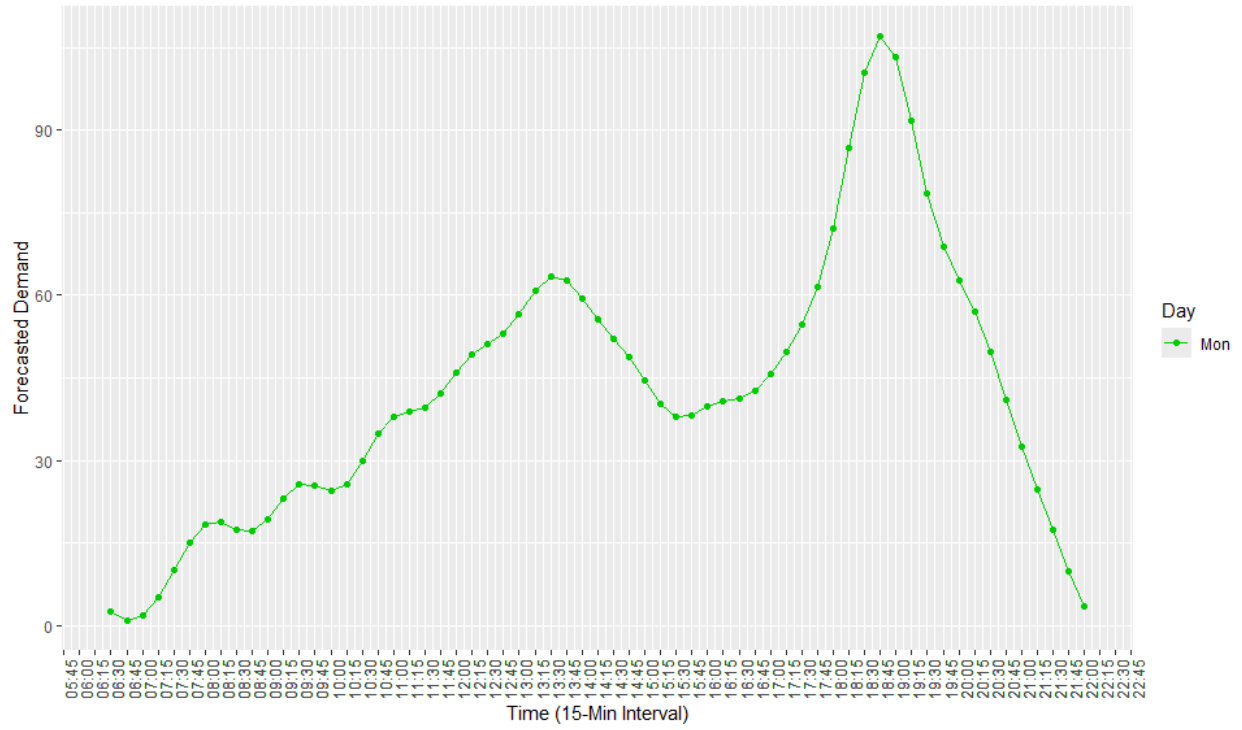
MASE (Mean Absolute Scaled Error): MAE scaled by the naive forecast. $MASE < 1$ is better, > 1 is worse than the naive forecast.

- Both showcased having great MASE in the training (Both < 1) yet when ran on the test set STL+ETS did far worse while TBATS was almost to 1.
- TBAT $1.01 > 1$, very slightly worse than naive forecast.
- STL+ETS $2.44 > 1$, incredibly worse than naive forecast.

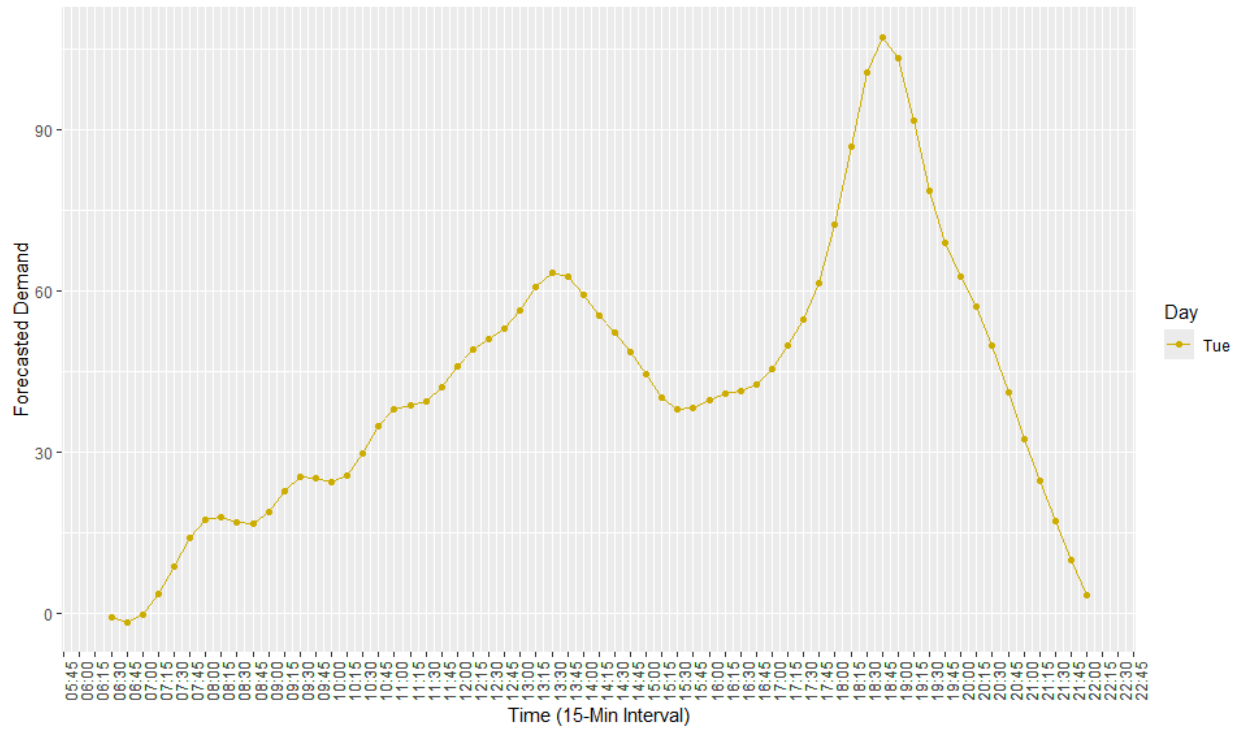
These scores were the ultimate reason for choosing TBATS over STL+ETS.

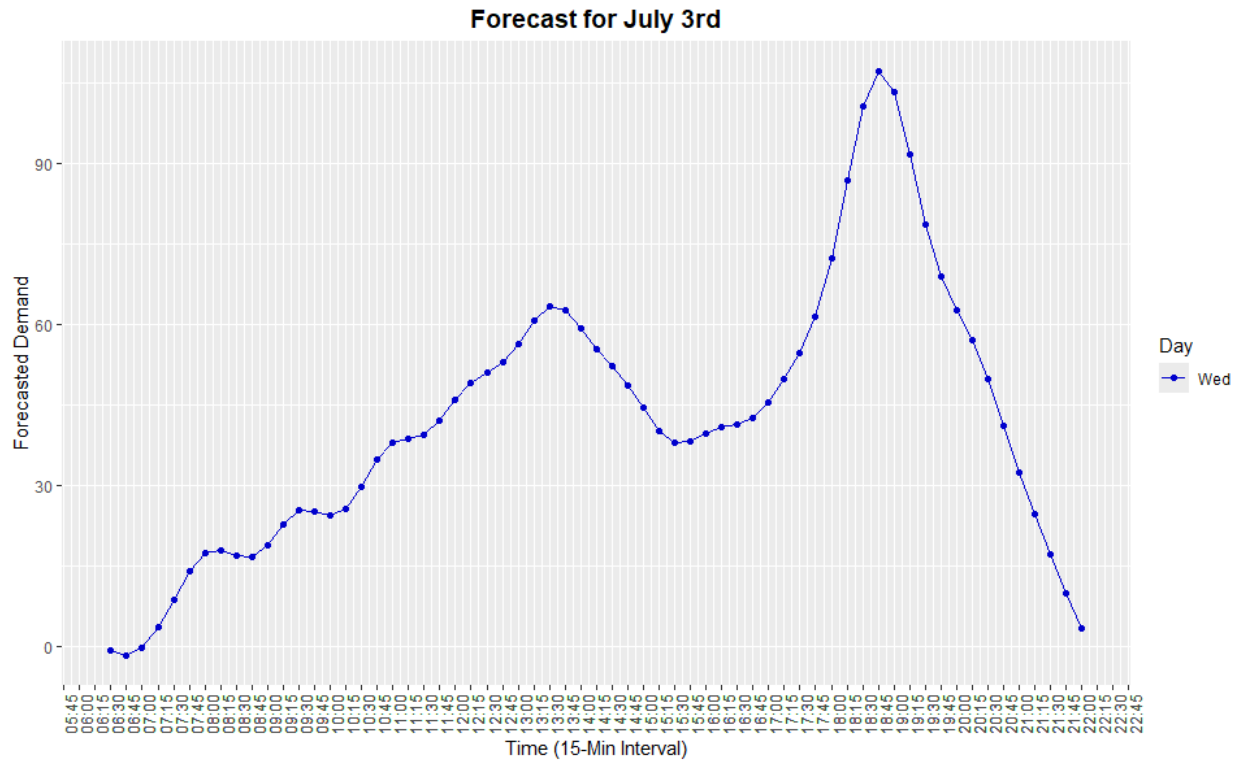
Forecasts:

Forecast for July 1st

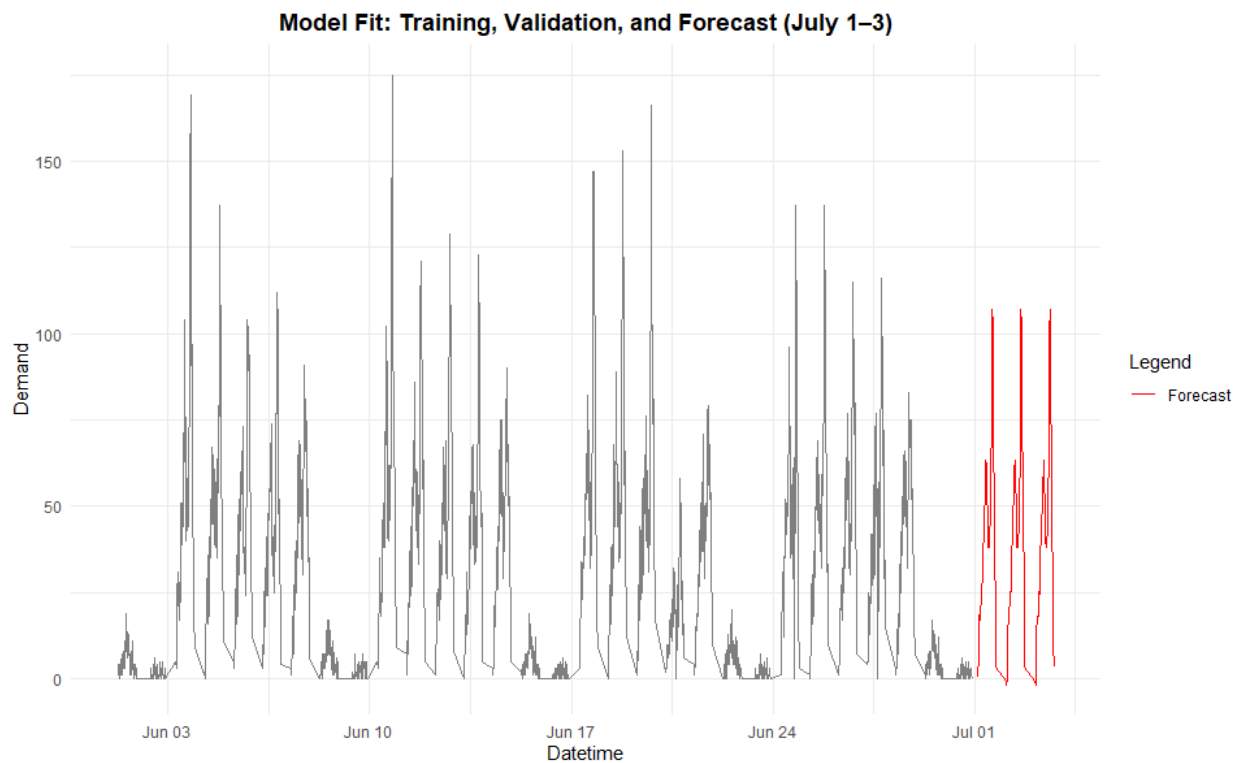


Forecast for July 2nd





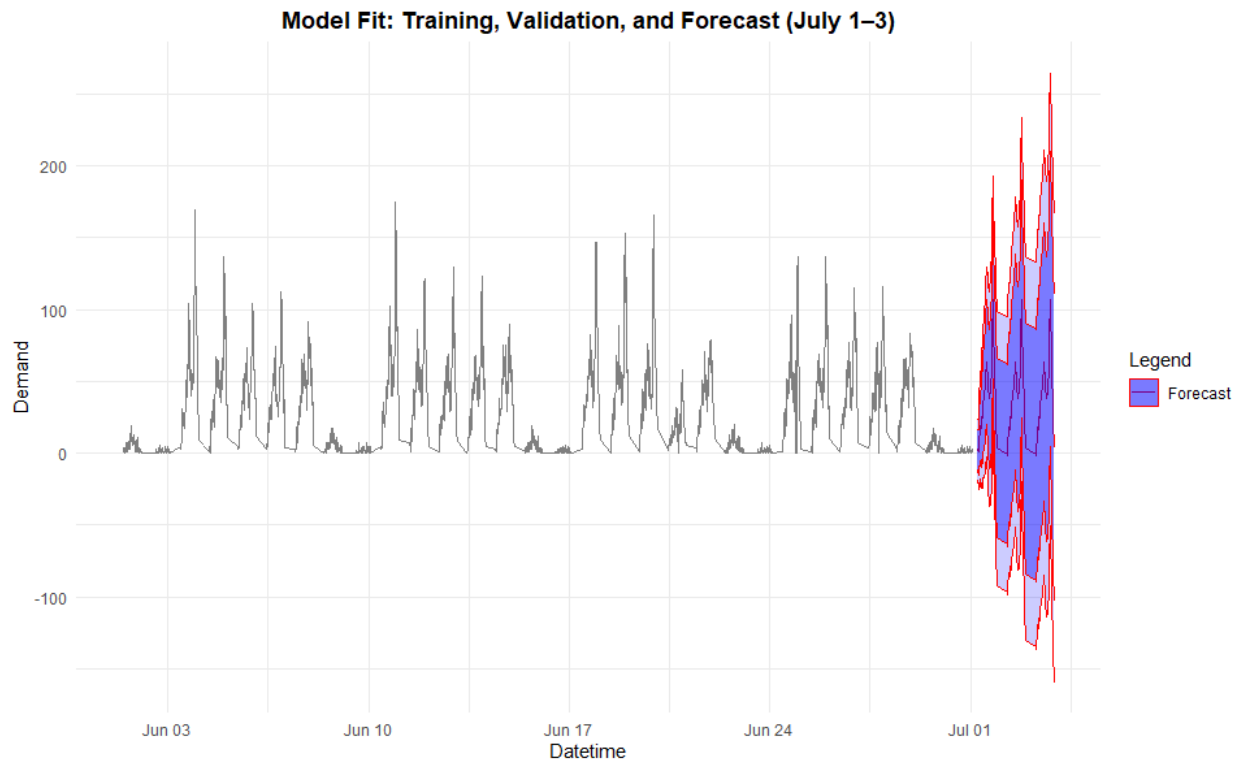
These forecasts are for the days individually while the next one shows a full fit of the existing data with the forecasted data.



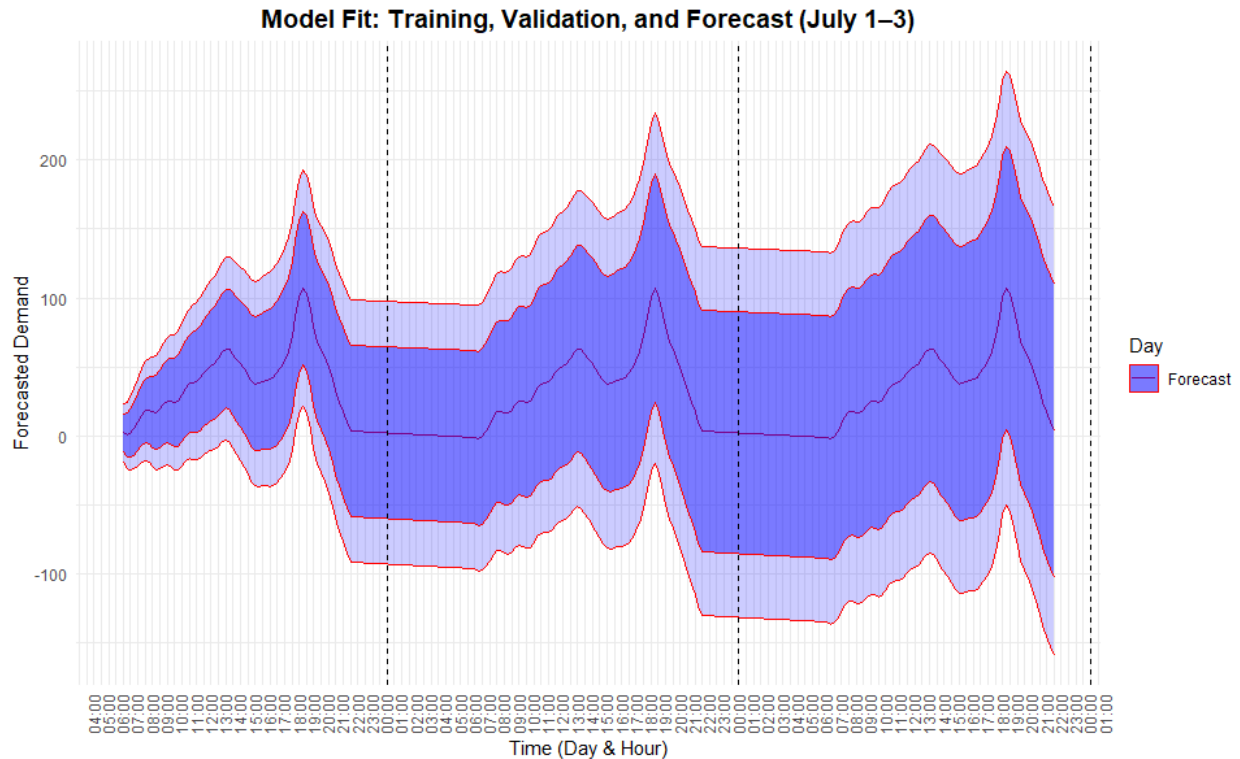
In this model, it looks like there actually was some form of negative trend. Whether this trend

was due to it being the end of the month or if in the year it actually goes down during this time, that is too hard to tell with just the 30 days provided.

Next is a model that shows the predictions including the 95% and 80% low and high predictions. The dark blue is the prediction range with 80% confidence while the light blue is the 95% confidence range.



And for a closer look that showcases just the forecasted days.



The intervals here are hourly so that it would be legible, but the data is still actually in 15-minute intervals. The sections show the days in order, the first being Monday, then Tuesday and then Wednesday. From this we can see that the confidence range in the forecast slightly increases as that is how time series forecasts generally go. The inner dark blue is the forecast with 80% confidence while everything within the light blue region is in the 95% confidence range.

References:

TBATS: <https://www.pluralsight.com/resources/blog/guides/time-series-forecasting-using-r>

<https://www.statology.org/tbats-in-r/>

STL+ETS <https://www.rdocumentation.org/packages/forecast/versions/8.24.0/topics/forecast.stl>

<https://www.statology.org/how-to-perform-time-series-decomposition-r/>

About STL forecasting <https://stats.stackexchange.com/questions/298560/is-stl-a-good-technique-for-forecasting-instead-of-arima>

Consulted ChatGPT & Perplexity for understanding and double checking their responses.