

# Final Project 6306

Tim Morales

12/7/2020

## Objective

Although many different regions have taken large hits to business during the global pandemic, few cities in the world were hit as hard as New York City. Claimed to be the epicenter during the initial wave of COVID-19 in the United States, New York City has seen huge losses to local business across a variety of industries including one that is a staple to the city itself, yellow cabs. Yellow cabs have covered the streets of New York City for decades, but with the rise of COVID-19, finds itself in trying times.

In this paper, we look at how the pandemic has affected NYC yellow cabs using individual pick up data from January 1, 2018 and May 31st, 2020. We compare the effects of the pandemic on yellow cabs in comparison to its main competition, for hire vehicles including Lyft and Uber. Through collecting, processing and summarizing over 33GB of individual ride information from the NYC Taxi and Limousine Commission, we use a time series approach to show how COVID-19 has hurt NYC yellow cab predictions in comparison to its for hire competitors.

## Data Manipulation

When understanding the data collection, there are a few things to note within the code. There were multiple instances of recording error especially for yellow cab individual ride data. For those errors, whether it be that the data was placed in the incorrect data frame or the year was incorrect, we dropped the observations as we cannot justify a proper way to solve or assure the true meaning behind those recordings.

Below is the Sparklyr Code used in collaboration with an AWS EMR cluster to process the for hire vehicle information. (SEE RMD CODE FOR ACTUAL CODE)

Below you can see the code used for the NYC yellow cab data. This was done within my local computer at the same time as the AWS cluster. Notice this was done in a loop and only 1 iteration is shown.

With those summary tables I have created I will use the rbind function to unit them.

Below I review each data frame individual and make sure they load properly and functionally. I also go through each to check for any issues. I do so for each month, each year, each vehicle type. (SEE RMD CODE FOR ACTUAL CODE)

## Yellow Cab

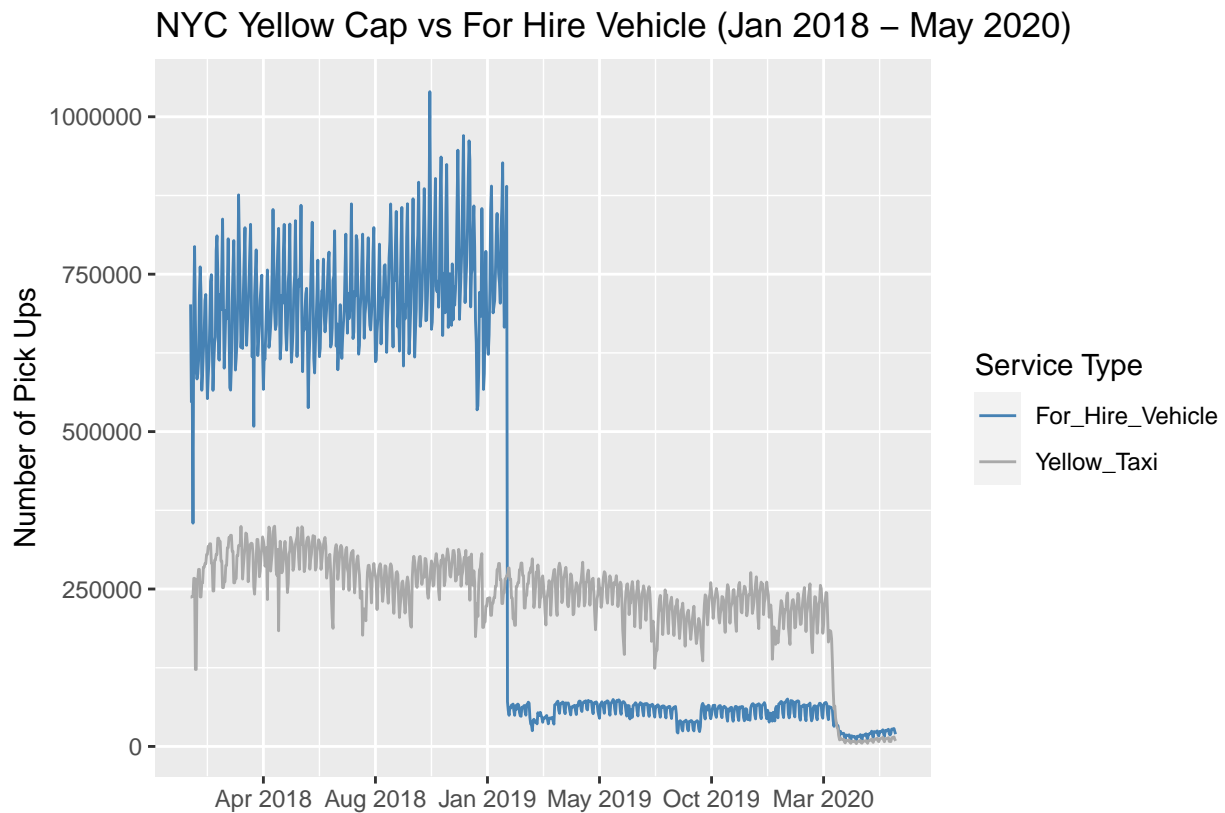
I do the same for 2019  
and 2020

## FHV

## Analysis

### EDA

I combine the data frames together below and continue with some EDA and visualization.



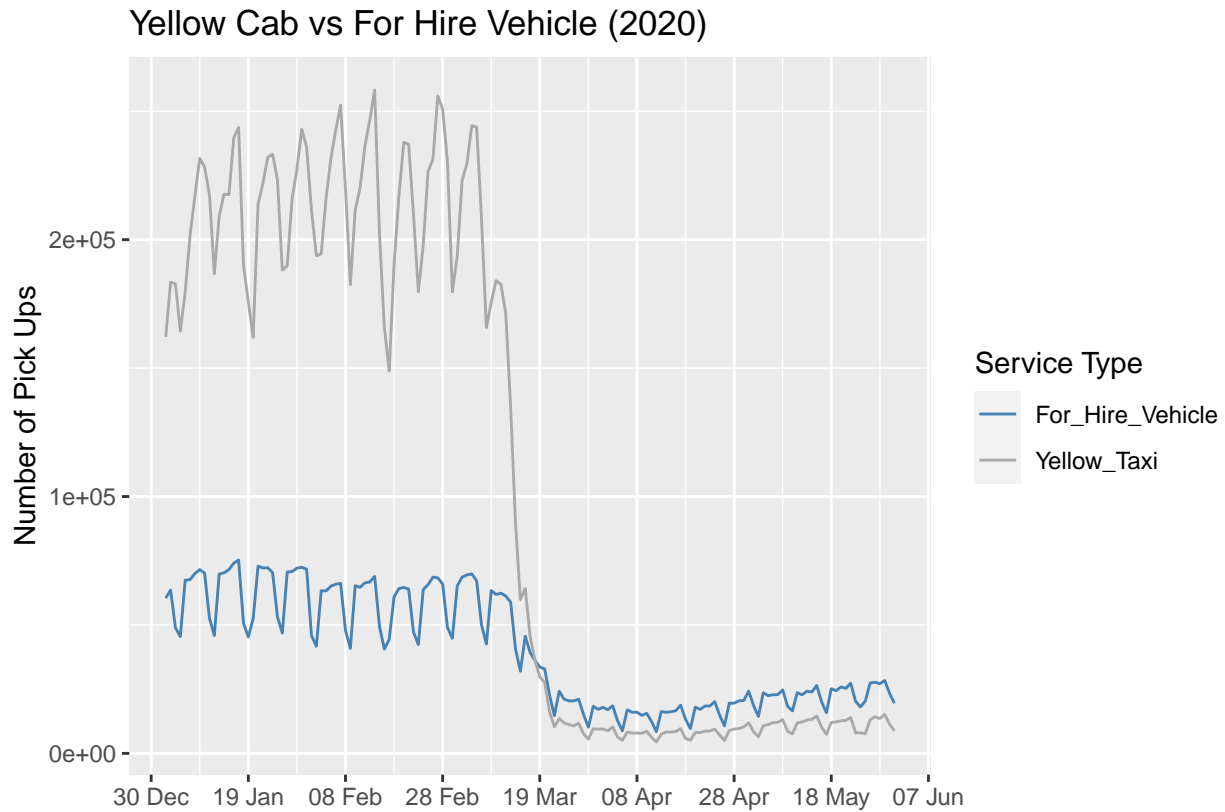
In our initial plot of the time series, we see the obvious drop off in picks due to COVID-19 for both service types. There is also a huge drop off in the for hire vehicle pick ups due to the city's cap on for hire vehicles in early 2019. This was done to protect the yellow cab system.

In terms of general trend, there seems to be a long term trend of a slight decrease in pick-ups dating to 2018 for yellow cabs. This raises a concern about stationarity in our data for the yellow cab series.

We can also see what appears to be weekly seasonality, with what we will assume to be increases in pick ups on the weekend.

Excluding 2020 and pre-2019 for FHV, we see no real signs of a multiplicative effect.

With such a large decrease in pick ups for the for hire vehicles as a result of the policy change in early 2019, we will only use the time points after the FHV cap in our modeling of the FHV estimations.

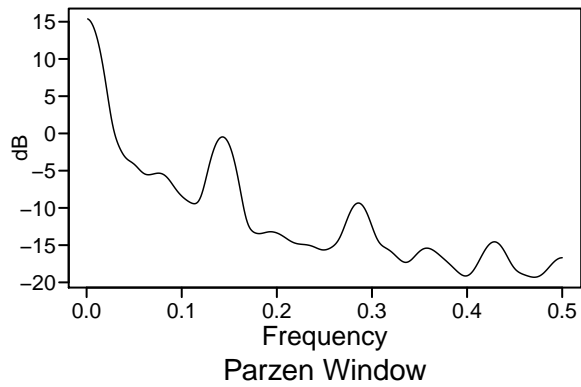
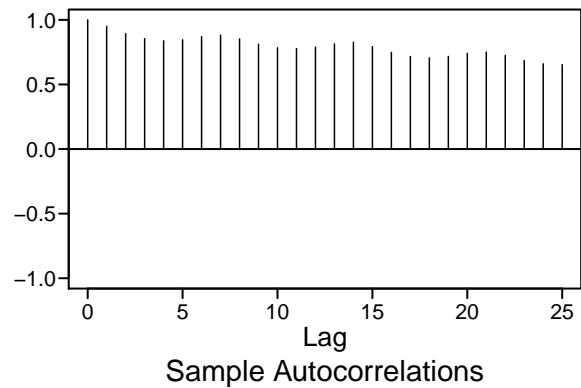


When we zoom in and look at just 2020, we can see the drop off in pick ups is much more severe for the yellow cabs. The large drop in yellow cab pick ups corresponds exactly to the pandemic striking in early - mid March. Although yellow cabs had dominated FHV in pick ups dating back to the cap introduced in 2019, once the pandemic hit, FHV took over as the main service for pickups. FHV also seem to be rebounding from the initial spike faster than yellow cabs.

## Yellow Cab

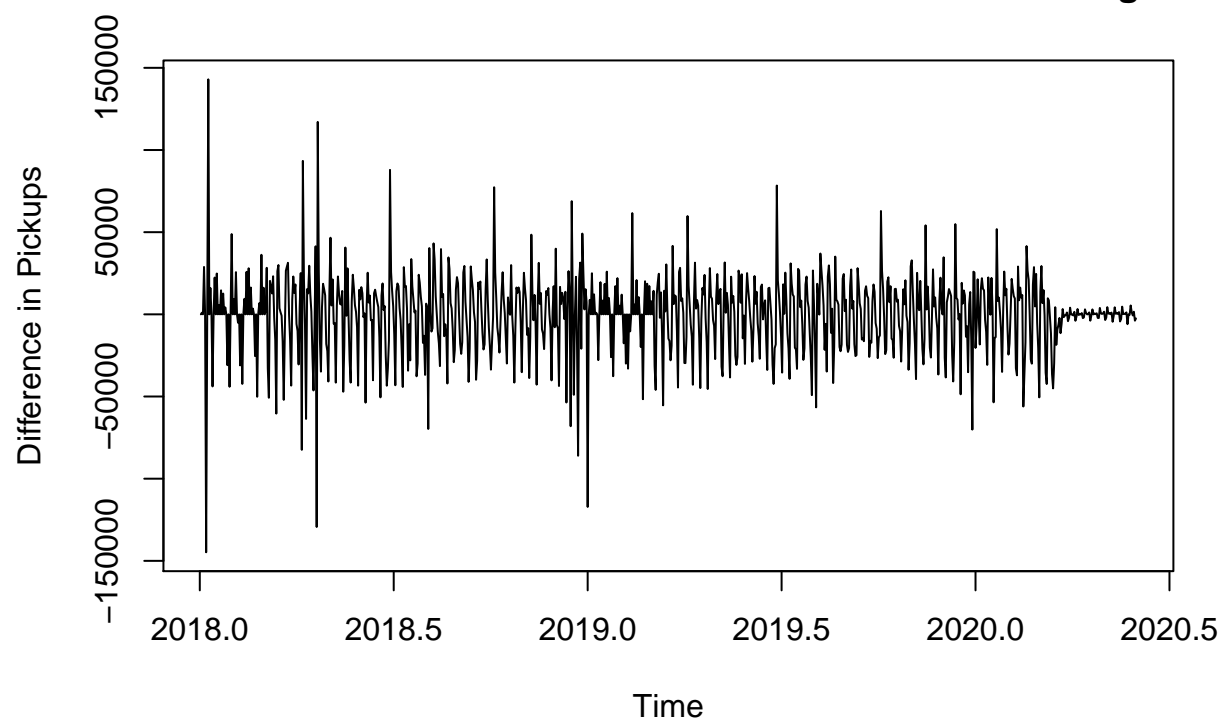
We begin the analysis conducting a SARIMA for the yellow cab time series. We want to see how the pandemic compares with what would have been the prediction without the pandemic striking.

Before the SARIMA, we look at the spectral analysis.



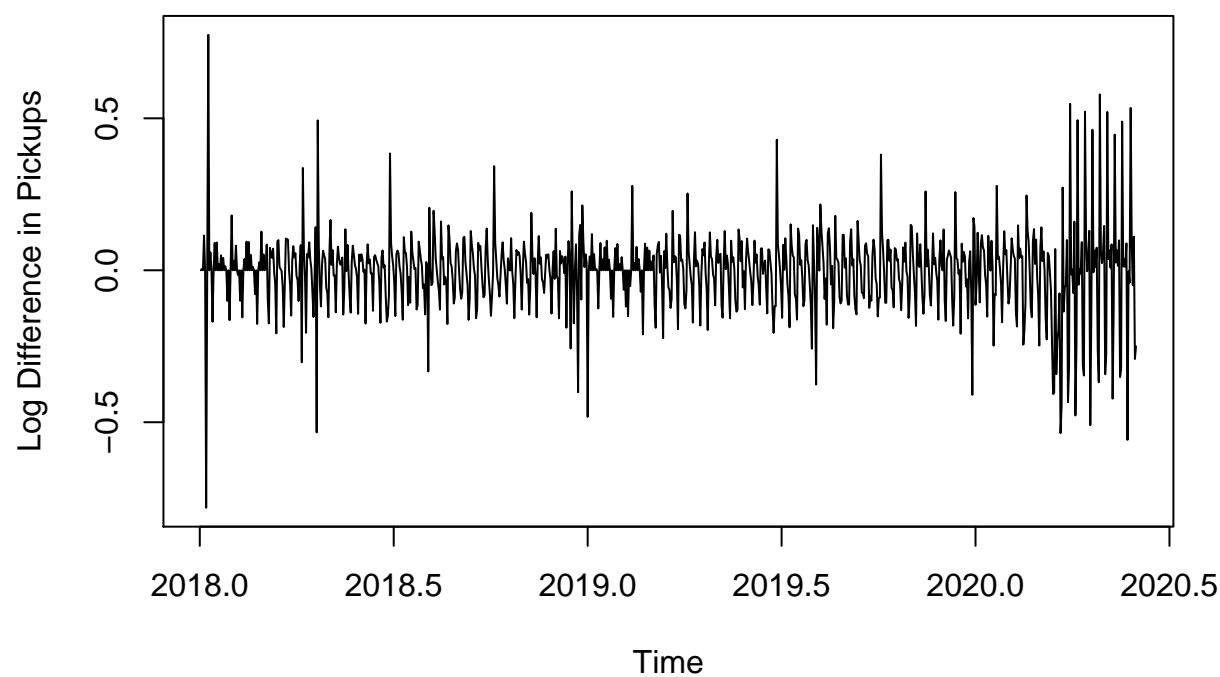
Since we saw some general decreasing trend in the data, we look to see how the graph would look if we used some first order differencing.

### Yellow Cab Time Series After First Order Differencing

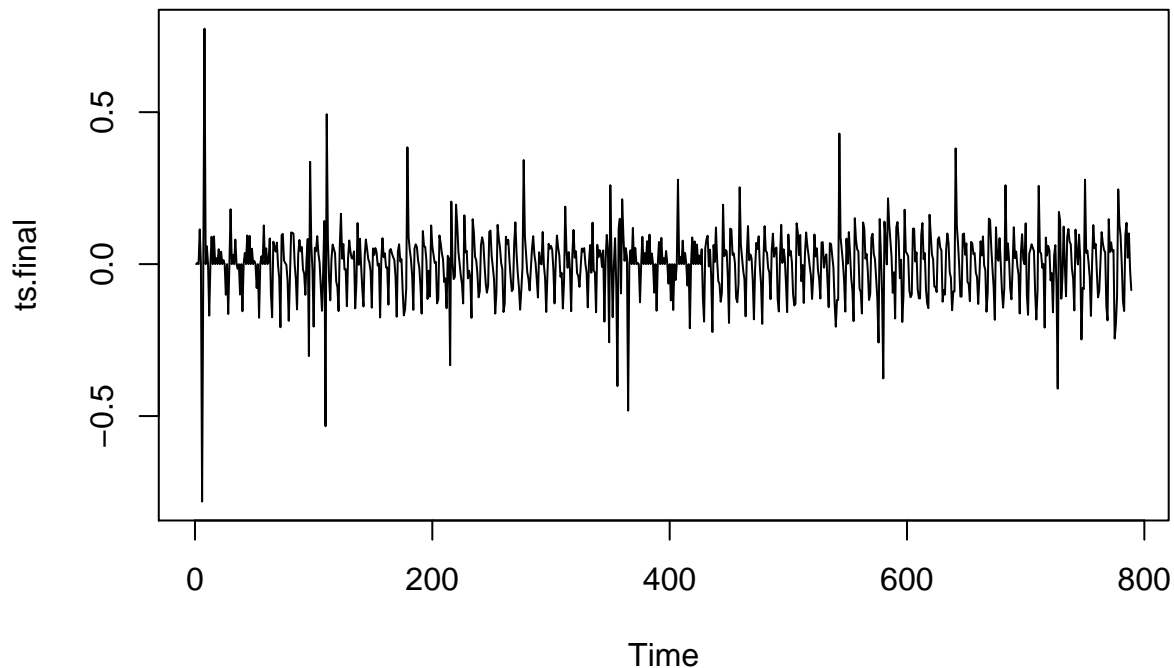


When we look at the first order difference, we actually do see some inconsistent variance, so we will perform a log transform.

### Yellow Cab Time Series After First Order Differencing and Log



Looking at the graph, the log transform and differencing creates a relatively stable time series. We just need to drop the post February information.



Finally, we fit the SARIMA using the `get.best.arma` function. The function works by finding the model with the minimum AIC when allowing parameter values to be any integer between 0 and 5. I do this in parallel. There is a maximum of 350 lags in the `arma` function, so I use the 350 lags. This also insures both yellow cab and FHV time series will be the same length.

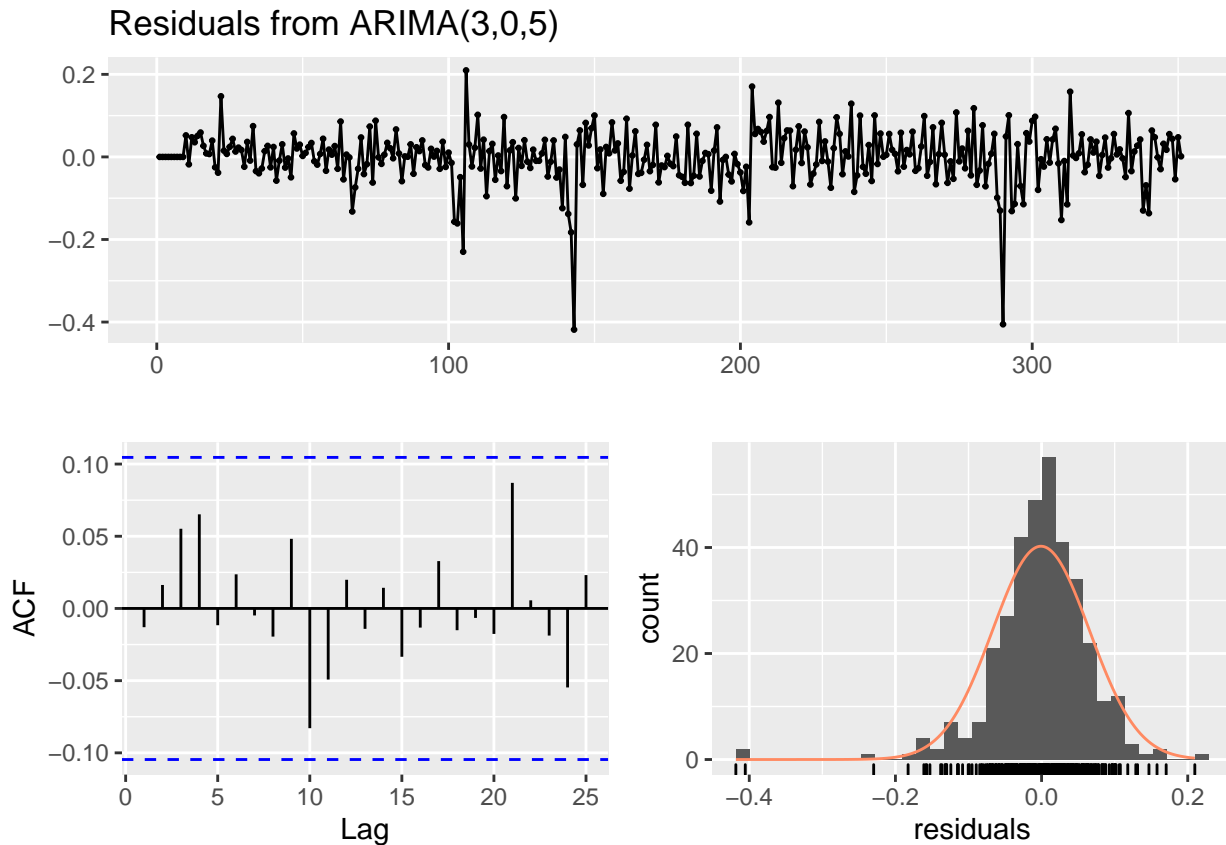
Since the above code takes so long to run, I just show a picture of the results below.

Results

We see that the best model is an SARIMA(3,0,5,5,1,3). It is interesting that it opted for the seasonal differencing.

We then fit that optimal model.

```
## Warning in arima(fitting.taxi, order = c(3, 0, 5), seasonal = list(order =
## c(5, : possible convergence problem: optim gave code = 1
```



```
##
##  Ljung-Box test
##
## data:  Residuals from ARIMA(3,0,5)
## Q* = 8.644, df = 3, p-value = 0.03442
##
## Model df: 16.    Total lags used: 19
```

The residual plots look good and very much like white noise except for a few points that we believe to be holidays.

I conduct a ljungbox to test of lag 1 autocorrelation in residuals.

```
##
##  Box-Ljung test
##
## data:  opt.taxi$residuals
## X-squared = 0.059674, df = 1, p-value = 0.807
```

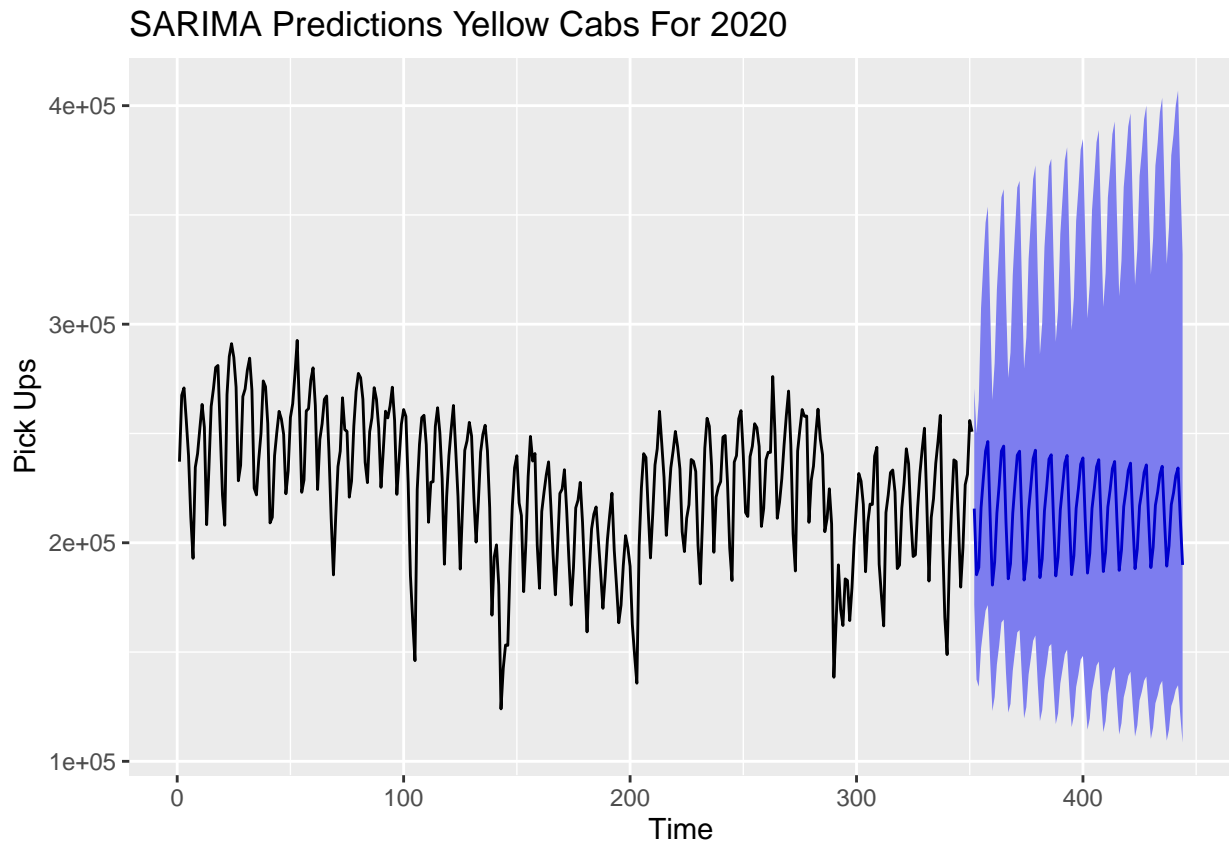
We fail to say there is significant autocorrelation. We are confident in trusting this models predictions.

Below I get some predictions for the COVID months to compare what was projected and what was actually observed for yellow cab pick ups.

Although the data was log transformed, I do not use the correction factor when converting back to original units since it was not in terms of means, but rather a count.

Below is a plot with a 99.99% confidence interval for the yellow cab pick up numbers.

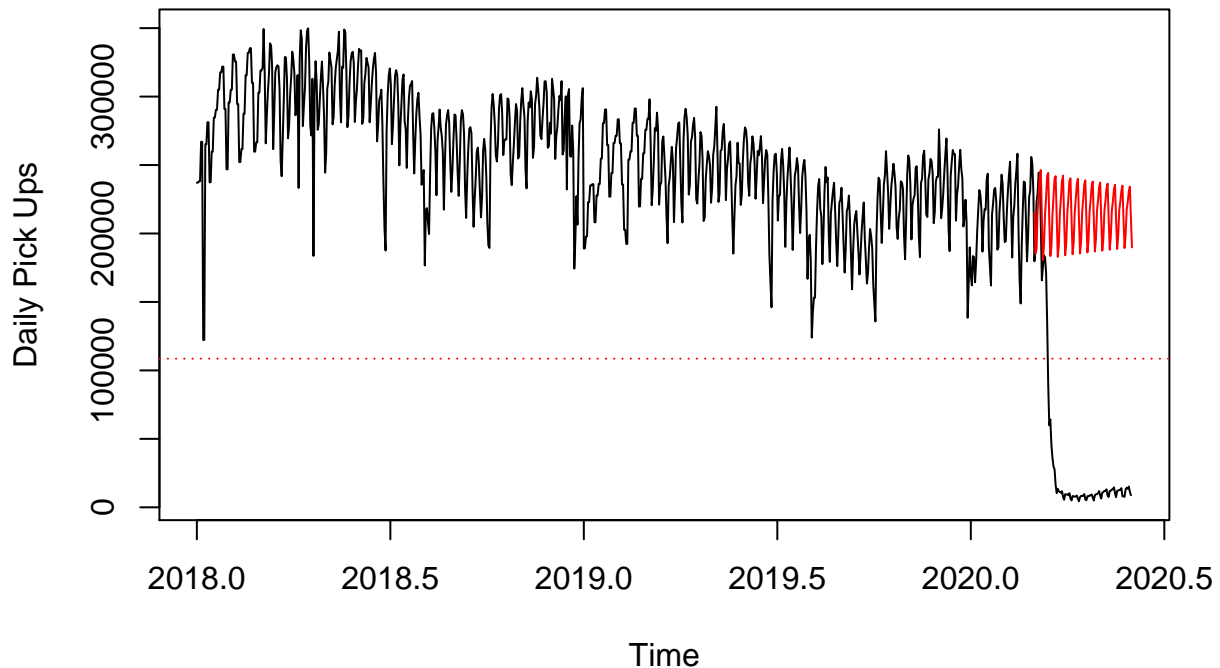
```
## Warning in predict.Arima(object, n.ahead = h): MA part of model is not
## invertible
```



We see the lowest points of the 99.9% confidence intervals through May 2020 are around 108,589 daily pick ups. We will use this minimum lower bound in the next section.



## How Drastic COVID Was For Yellow Cab Pick Ups



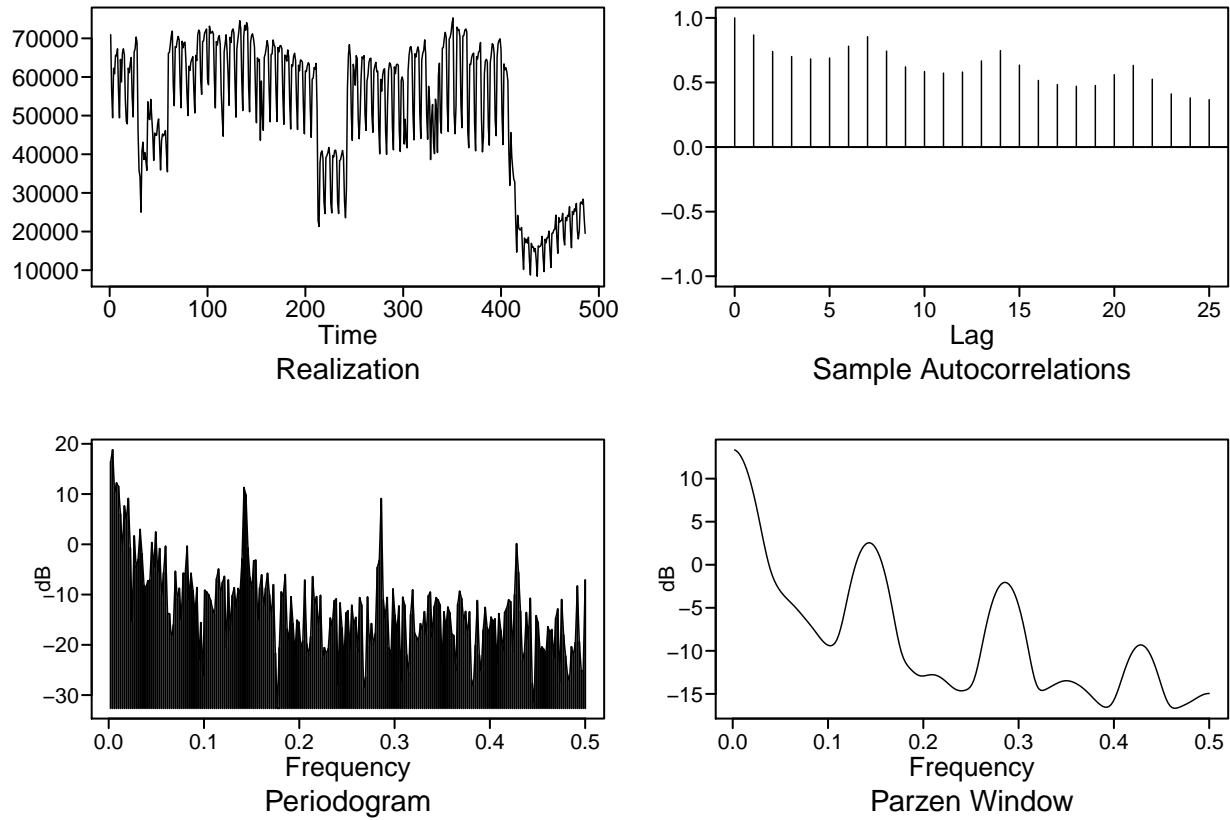
In the above plot where we superimpose predictions on the actual time series, we can see how devastating the pandemic was for the yellow cabs and how drastic the difference is between projection and actuality.

We use the red dotted line to denote the minimum bound of the 99.9% confidence interval for any day between March and May 2020. Even when having 99.9% confidence, the actual number of pick ups is well below that minimum bound of our prediction within a matter of days of the pandemic!

### For Hire Vehicles

Now we turn our focus to the FHV and see how drastic the pandemic was for this service in comparison. As mentioned before, we will only focus on post-February 2019 due to the policy changes within the city. Using this limited time frame also allows us to use 350 lags and keep a consistent length for both series.

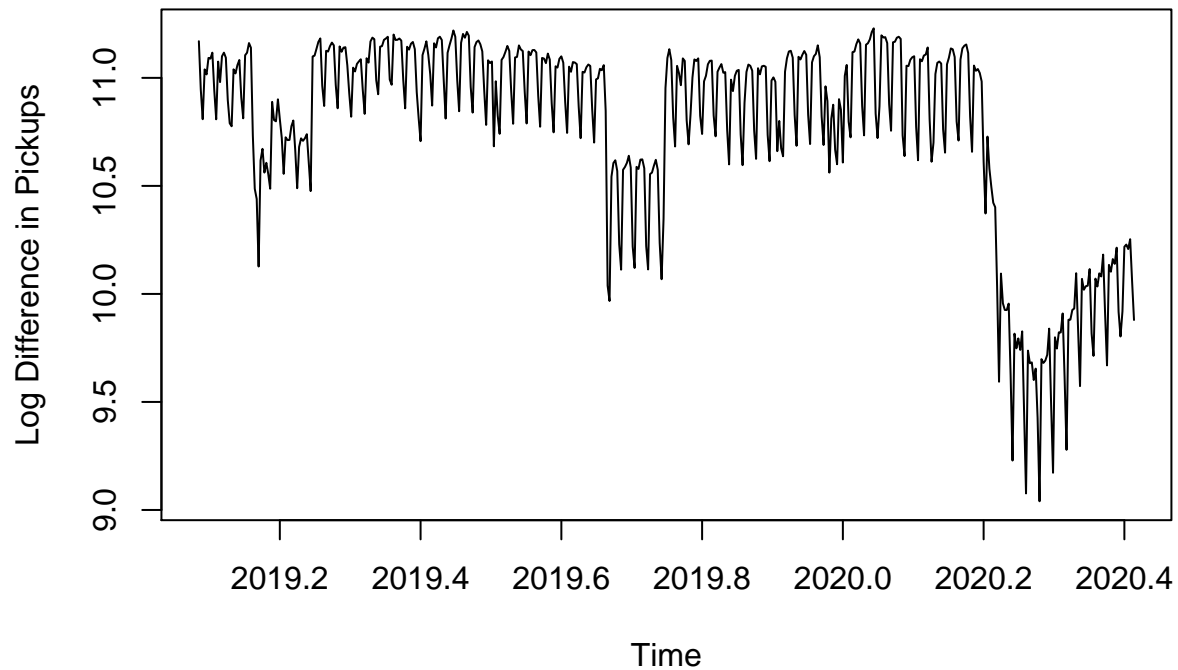
Spectral Analysis



In the spectral analysis of truncated series, we see some very similar things to the yellow cab analysis such as the large spikes in frequency corresponding to weekly cycles and the large spike around 0. The time series plot itself is interesting as there are random spikes and decreases seen at times  $\sim 30$  and  $\sim 200$ . This suggests it may be difficult to model.

We use a log transform for the inconsistent variance throughout.

## FHV Time Series After First Order Differencing and Log



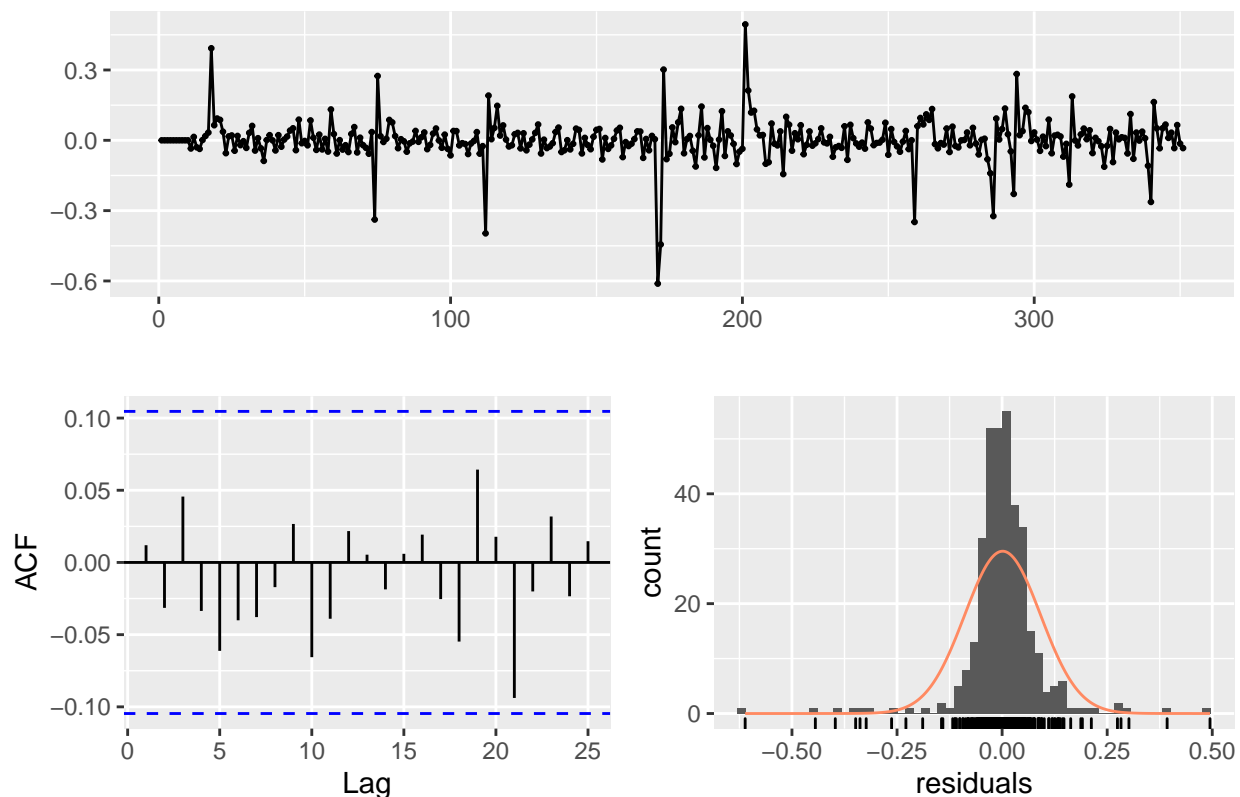
The log transformation seems to really help with the mutliplciative effect. Since there are obvious needs for some pretty severe differencing, we just fit the optimal SARIMA using the same function as above.

Results

We see the results are an SARIMA(5,0,4,4,1,5)

```
## Warning in arima(log.ts.fhv[43:393], order = c(5, 0, 4), seasonal = list(order =  
## c(4, : possible convergence problem: optim gave code = 1
```

## Residuals from ARIMA(5,0,4)



```
##
##  Ljung-Box test
##
## data:  Residuals from ARIMA(5,0,4)
## Q* = 13.246, df = 3, p-value = 0.004134
##
## Model df: 18.   Total lags used: 21
```

Although there does seem to be some random spikes in the residuals, it again may be due to something like holidays. We again follow up with a Ljung-Box test.

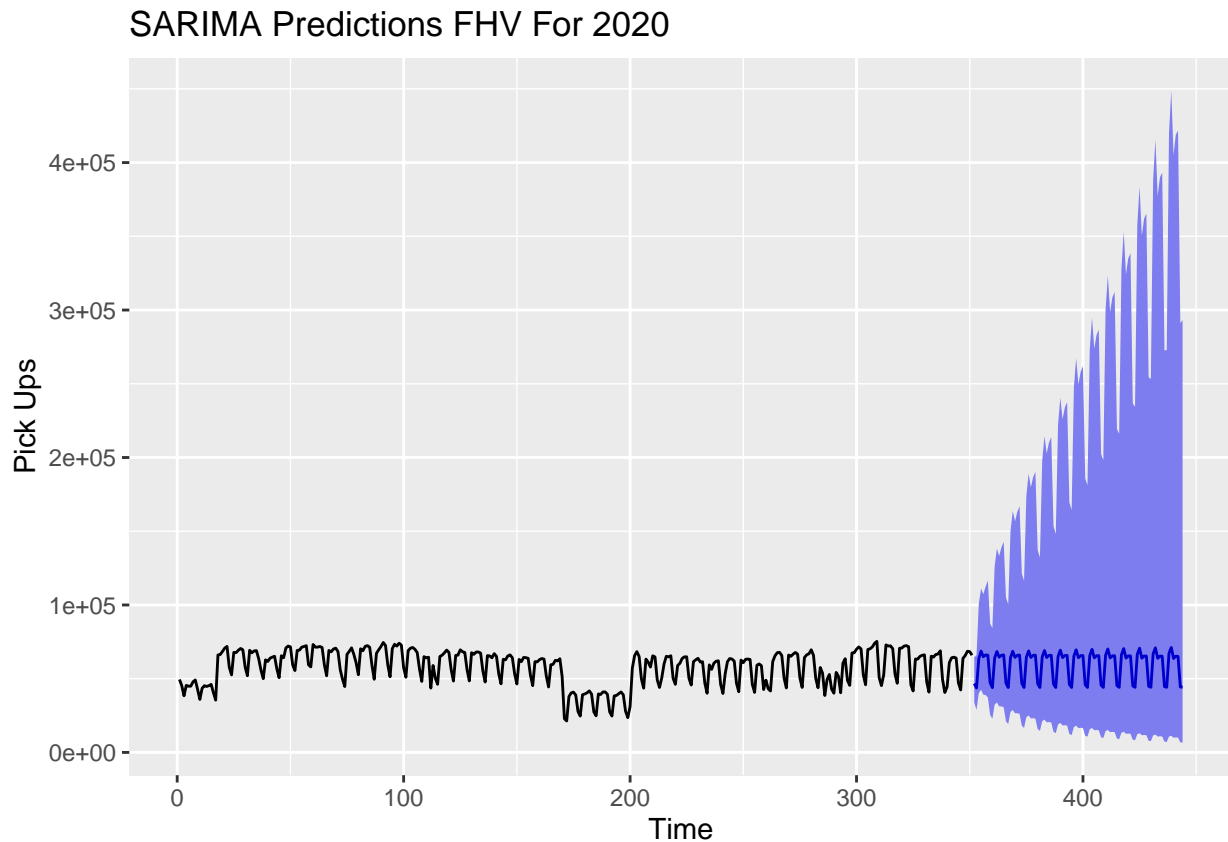
```
##
##  Box-Ljung test
##
## data:  opt.fhv$residuals
## X-squared = 0.050495, df = 1, p-value = 0.8222
```

We again fail to reject the null hypothesis that there is autocorrelation in the residuals. We therefore believe that these predictions will be reasonable and continue with the analysis.

We use the same analysis getting the forecast with a confidence interval at the 99.9% confidence level, but this time for FHV.

```
## Warning in predict.Arima(object, n.ahead = h): MA part of model is not
## invertible
```

```
## Warning in predict.Arima(object, n.ahead = h): seasonal MA part of model is not
## invertible
```



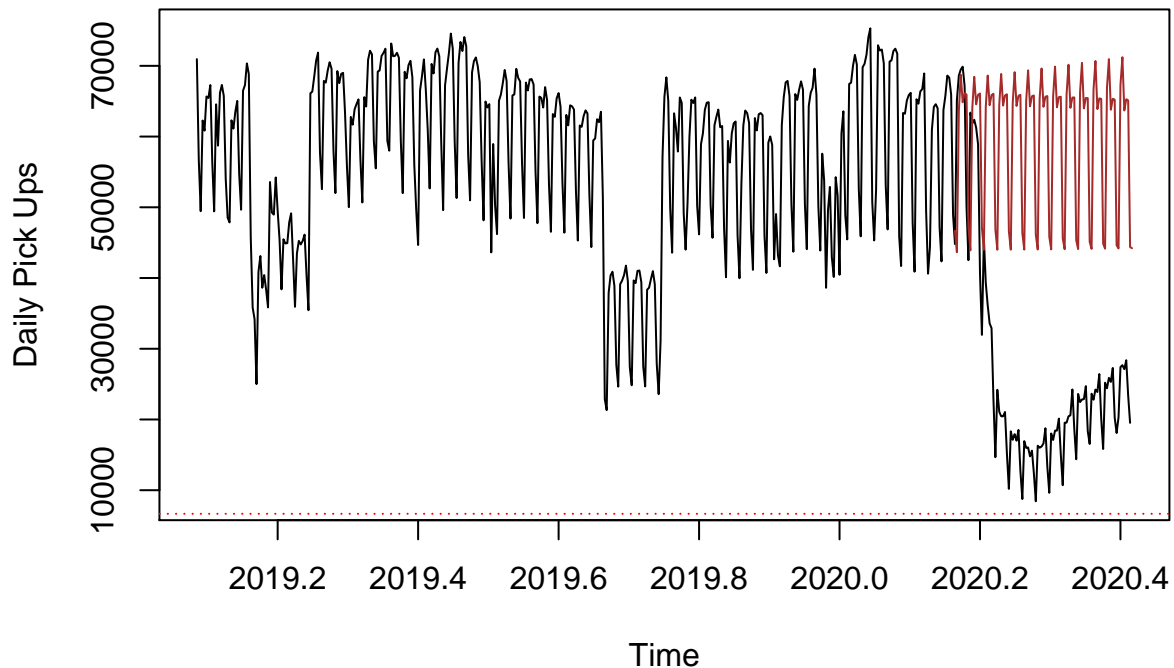
Interestingly, the model does not have very strong confidence and has a pretty large interval, especially in the increasing direction.

We see the lowest points of the 99.9% confidence intervals through May are around 6,668 daily pick ups.

```
## Warning in predict.Arima(opt.fhv, n.ahead = 93): MA part of model is not
## invertible
```

```
## Warning in predict.Arima(opt.fhv, n.ahead = 93): seasonal MA part of model is
## not invertible
```

## How Drastic COVID Was For FHV Pick Ups



When we superimposed the predictions and minimum confidence interval bound on the full time series for the FHV, we see that these COVID-19 drop offs are actually WITHIN the 99.9% confidence intervals for the FHV vehicles!

This suggests that the for hire vehicle industry may have been able to actually account for and predict losses as large as the once seen from the pandemic.

For the FHV, we also see a relatively quick increase in pick ups from March to May 2020, this may be because of how resilient the FHV industry has proven to be. After taking those huge losses in early 2019 from the NYC policy changes, maybe FHV's have been more prepared for future losses.

## Conclusion

Our analysis shows in a visual manner that yellow cabs have been hit harder by the 2020 pandemic than their for hire vehicle counterparts. This is in terms of the sheer decrease in the number of pick ups and how unpredictable the losses were for the yellow cabs. What was a true black swan event for the yellow cab industry, actually appeared to be relatively predictable for the for hire vehicle companies.

In the future, it appears that for hire vehicles will again retake control over the vehicle service industry as they did before the 2019 NYC policy change.

## Predictions

### Regression

For the regression task, I fit four different neural networks, one for each year used in the prediction. The data for each year was the corresponding month and year included in the test set data. In my training set data which was 500,000 rows from the same month as the prediction set I did an 80/20 test split. For each set, I feature engineered date and time variables to get a variable about which hour it was in the day and

created 4 different flag variables. The flag variables denoted if it was late night/past 5pm, if the drive was over a mile, under 5 minutes or over ten minutes. I centered and scaled all numeric variables.

Using tensorflow and Keras, I used keras sequential API to fit a three layer neural network for each year. The neural network used global normalization for first layer weights and had 12 neurons on the first layer. I used relu activation functions for all layers other than the final layer which I used linear. The second layer had 5 units and the final layer had 1 single unit.

I trained the model using RMS Prop with a mean squared error activation function. The model was training for 30 epochs with a minibatch size of 128. For each epoch the training set was split into a 80/20 training and validation set to monitor overfitting.

The model was then used to predict the test set and test RMSEs can be seen below. The models were used to predict the same year they were fit on for the final data predictions saved in CSVs.

List of variables included:

hour,passenger\_count,trip\_distance,fare\_amount,dropoff\_longitude,dropoff\_latitude, day,surcharge,late\_night,over\_mile,over\_five,

RMSE on test set below (additional cleaning done to just that year)

TEST SET RMSE

2020 : 1.647215 (involved outlier removal did not use long and lat)

2018 : 1.62216 (did not use long and lat)

2011 : 1.3708

2009 : 1.520322

## Classification

For classification, I again used neural networks with the same feature engineering. The main differences in this neural network was the number of neurons for each layer were 12, 4, 2 and the activation functions were relu, relu, sigmoid respectively. The optimizer for the classification was adam and the model was trained using a binary cross entropy loss function. This model only was trained for 5 epochs but with a minibatch size of 25.

List of vars included :

hour,passenger\_count,trip\_distance,fare\_amount,tip\_amount,late\_night,over\_mile,over\_ten, over\_five,

TEST SET ACCURACY

2020:

Accuracy : 0.9071

Sensitivity : 0.8297

Specificity : 0.9848

2019:

Accuracy : 0.8651

Sensitivity : 0.8033

Specificity : 0.9432

2017:

Accuracy : 0.9207

Sensitivity : 0.8411

Specificity : 0.9997

2015:

accuracy 0.9317

Sensitivity : 0.8622

Specificity : 0.9979