

A lot of the insights gathered so far revolve around how Divvy's usage fluctuates throughout the week and throughout the year. As the usage differs in this way, I'd like to predict station stocking requirements by the day of the week and month of the year. For example, if Divvy wanted to know what the stocking requirements should be for station x on a Monday in August, the model would be able to provide the answer.

For this model I would like to use the ratio of number of trips in to number of trips out for each station on a given day as the target variable for the station's stocking requirements. I believe this will be more valuable because ratio can relate to the number of empty and filled slots each station should have. The total number of rides associated with the station alone on will not give us the same insight.

With the goal of a supervised model in mind, I began this portion of the analysis by organizing the trip data by station, day of the week, and day of the month so I could calculate the trips in to trips out ratio. The number of rides in and number of rides out were found for each station by day of the week and day of the month, but looking closer I found that not every day of the week was accounted for each month for every station. The reason behind this was that not every station had a ride come in or out on every day of the year. This trip information we'd like to have marked as 0 was not available in the current grouping. By unstacking the day of the week, and the month of the year to the top of the dataframe I was then able to then view the ride counts for every day of the week for every month. Those days without trip numbers in them showed up as NaN, and I was able to fill them in with 0s, and restack the day of the week and month of the year variables.

The ratio could then be calculated using complete number of rides in and number of rides out information. A few things need to be addressed regarding this ratio however. When the rides out (the denominator) is 0 for this ratio, it produces a NaN value. To solve this, when rides out is 0, the ratio is set equal to the numerator (rides in). So if 12 rides came in and 0 rides went out the ratio would be 12, not 0 or undefined. This more accurately represents the reality of that particular day. The next instance of concern would be when the number of rides in for the day (the numerator) is 0. As 0 rides in and 1 ride out, and 0 rides in and 12 rides out would need to be accounted for differently we can't have them both with a ratio of 0. In these cases the ratio is set to  $1/\text{rides out}$ , this will more accurately reflect the reality as well. The last potential concern would be days with 0 rides in and 0 rides out. This should not read as 0 either, as an even number of rides in and rides out needs to be represented by a 1. In these cases, the ratio was calculated as a 1.

After the ratios were calculated, the following variables were then added to the newly organized dataframe; average trip duration, station capacity, station latitude, station longitude, and group number. The group numbers were assigned in previous analysis and are based on traffic to each station. 1 representing a high traffic station, and 4 representing a low traffic station.

Using longitude and latitude alone would not be a big help while building a model because each set of coordinates will be unique. Two stations right next to each other would not be seen as related by the model as the coordinates would differ. To better group the bike stations by location, I split stations into two different groups, north and south. These were split based on the median latitude. I also split the stations into four territories to further narrow down their location. These four territories are marked as 1 (northwest), 2 (northeast), 3 (southwest), and 4 (southeast). The north and south dividing line was again the median latitude, and the east to west dividing line was the median longitude for the stations.

The last thing I wanted to look at before feature selection was to see if there was a good way to segment the bike stations by clustering. Using a KMeans model I looked at grouping the bike stations by station capacity and the average in/out ratio. Using the elbow method I determined that 3 clusters would be optimal, and fit the model with the capacity and ratio data for each bike station. Using the labels from the models predictions, I assigned a cluster id to each of the bike stations in the dataframe. Cluster id, location, and the previous grouping by traffic level, are all ways help group the stations by similarities. I checked the correlation coefficients of the groups based on traffic, north/south sections of the city, the four territories, and the cluster ids with the in/out ratios. Although there wasn't a strong correlation with any of the four types of groupings, and the ratios, the groups based on traffic had the highest correlation. Territories had a slightly stronger correlation than the north/sections, so territories was chosen to be used going forward for insight based on location.

Once this information was gathered together in the new data frame, I used Lasso regularization to provide some insight on my feature selection. Looking at the Lasso coefficients, none of the variables had strong coefficients on their own. Because this was the case, all 4 of the predictor variables (average duration, station capacity, group, and territory) would be used in the models to keep as much predictive power as possible. With a smaller number of variables we don't have to worry about the dimensions being too high in this case.

As we're looking to predict a continuous variable (the trips in to trips out ratio) I started with the linear regression model. Before training this model though, I needed to change the territory, day of the week, and month of the year variables to categorical data. Although the data are numbers, they represent different categories, and don't relate to each other by their size. Next, dummy columns were added to turn the categorical data into a series of 1's and 0's depending on what category they belonged to. This will allow us to use the categorical information in the regression model.

After splitting the data into training and testing sets, and training the linear regression model, it became clear that model was not performing well. Afterwards trying ridge regularization as well, I found a similar result of just over 1% accuracy. Looking back at scatter plots and correlation coefficients for these variables it was not surprising to find that they don't have a strong linear relationship with the calculated ratio.

To help solve the main goal of predicting stocking requirements for the bike stations I grouped the ratios into 5 different classes. These classes were based on the value of the ratio as the ratios near one another are going to have very similar stocking recommendations. For example whether the ratio is .8 or 1.2, you still have approximately the same amount of bikes coming in vs leaving a station on the particular day, and you would ideally have that station stocked with equal number of filled and empty slots. Once these stocking classes were created, this became a classification problem instead of regression.

Before working with different models, I wanted to have a benchmark to test them against. The benchmark used was the average ratio for each station for all of 2017. I wanted to test if these models would be more effective than putting them in a particular stock class based on this average. Looking at all of the stations together, the accuracy of our benchmark in predicting the correct stocking class was 68%. Separating the stations into the four groups mentioned previously, the accuracy of the benchmark was 86% for group 1, 75% for group 2, 54% for group 3, and 58% for group 4.

With these benchmarks in mind, I looked at 4 models while searching for the most accurate predictions. KNeighbors, Random Forest, SVM, and Multinomial Naive Bayes. Although often used with text data, the bayes model was included as it can be of help in finding a dependent categorical variable that has more than 2 possible categories.

Each of these models' hyperparameters were tuned using GridSearchCV. If the value on either end of the grid search was found to be the best parameter, another grid search was run to check the range above or below that value. Classification reports were then run on each model as well.

The accuracy results for these models can be found below:

	Not Grouped	Group 1	Group 2	Group 3	Group 4
<b>Benchmark:</b>	68.46	86.03	75.34	54.48	57.89
<b>Models:</b>					
<b>SVM</b>	70.00	84.21	74.32	52.83	68.57
<b>KNN</b>	70.10	84.37	74.03	51.37	66.33
<b>Multinomial NB</b>	67.13	84.21	73.87	51.45	60.46
<b>Random Forest</b>	63.56	80.97	68.08	43.5	60.66

SVM and KNeighbors performed the best across the board, and both beat the benchmark accuracy score of 68% for the ungrouped model by 2%. When the models were segmented based on our 4 groups however, they scored just below the benchmark scores for each group with the exception of group 4. SVM and KNeighbors beat the benchmark score for that group by 11% and 8% respectively.