

# Denver 2016 B-cycle Ridership Data Exploration



## 2016 SUCCESSES



Source: [http://denver.bcycle.com/docs/librariesprovider34/default-document-library/dbs\\_annualreport\\_2016\\_05.pdf](http://denver.bcycle.com/docs/librariesprovider34/default-document-library/dbs_annualreport_2016_05.pdf)

**Denver B-cycle** is a non-profit public bike sharing organization operating an automated bike sharing system called Denver B-cycle. Its mission is to "serve as a catalyst to fundamentally transform public thinking and behavior by operating a bike sharing system in Denver to enhance mobility while promoting all aspects of sustainability: quality of life, equity, the

environment, economic development, and public health" its purpose, its organization and discuss its relevance to this exploration.

Denver B-cycle posts its trips data set on its website as soon as its annual report is released. Trips data have been available since 2010. The 2016 annual report and its associated dataset for this report were obtained from [Denver B-Cycle website](#). The original plan was to use the 2015 dataset to continue the effort by Tyler Byler who published a report, [Exploring 2014 Denver B-cycle Ridership](#). In his study Tyler indicated that “most calendar and clock variables were highly significant when predicting ridership, and weather variables such as temperature and amount of cloud cover appear to be as well”. The original plan for this report was to use 2015 data to continue Tyler’s work. However, the 2016 data became available at the end of February 2017, so gears had to be rapidly shifted to use this data instead. To this end, the reporting style will follow Tyler's study to provide seamless continuity and good reference on trends and analyses.

This study has three parts:

1. Explore the Trips datasets and visualize the data to provide useful and interesting information.
2. Deploy a variety of regression models to train and test the data.
3. Deploy a variety of classification models to train and test the data.

## Part 1: Data Exploration

### Data Acquisition

Data for this study was downloaded from several sources and combined using the following steps:

1. Downloaded B-cycle 2016 Trips and Kiosk data from [Denver B-Cycle website](#). The columns names were changed to comply with Python code best practices.
2. Created a list of the 7921 combinations of the 89 checkout/return kiosks. Used [Google Distance Matrix API](#) to provide the bicycling distance and time between each checkout and return kiosk. Adopted Tyler’s method of finding the average distance by taking the distance from each checkout-return pair’s distance separately then averaging it. As he pointed out in his study, this approach was taken “because of the large number of one-way streets in the Denver downtown area where the kiosks are highly clustered”. Google only supports a maximum of 2500 requests a day, it took four days to obtain this data.

3. Obtained daily and hourly weather data via [Dark Sky API](#) for all of 2016. Dark Sky supports up to 1000 requests per day.

## Basic Ridership Statistics

### Number of Rides

The B-cycle data, as downloaded, contained 419,611 rows of trips data. Under normal circumstances this would mean that 419,611 B-cycle trips were taken in 2016. However, the [2016 Denver B-cycle annual report](#) acknowledged 354,652 total trips for the year. The breakdown was as follows:

<u>Membership Type</u>	<u>Number of Trips</u>
Annual (And Annual Plus)	193,113
Flex Pass	3,565
30 Day	54,004
24 hour online	117
24-hour Kiosk	103,853
<b>Total Trips</b>	<b>354,652</b>

The Trips dataset reported the following breakdown:

<u>Membership Type</u>	<u>Number of Trips</u>
Annual (Denver B-cycle)	82199
Annual Plus (Denver B-cycle)	84271
Flex Pass (Denver B-cycle)	3565
Monthly (Denver B-cycle)	54004
24 hour online (Denver B-cycle)	117
24-hour Kiosk Only (Denver B-cycle)	87315
<b>Total Trips</b>	<b>311,471</b>

There were several other Membership Types that were also listed under “Denver B-cycle” in the User’s Program column:

<u>Membership Type</u>	<u>Number of Trips</u>
Denver B-cycle Founder (Denver B-cycle)	18003
Not Applicable	64959
Single Ride (Denver B-cycle)	16526

In particular, the “Not Applicable” membership type accounted for more than 15% of the 419,611 trips. Perhaps some of these trips were used in the Denver B-cycle annual report.

Also over 2.3% of the Denver B-cycle rides (9,954 rides) had the same checkout station as return station with a trip duration of only 1 minute (Figure 1). Again, Tyler’s explanation of why these trips should be removed from the dataset makes sense - “I believe these should be filtered out because I believe the majority of these “rides” are likely people checking out a bike, and then deciding after a very short time that this particular bike doesn’t work for them. I believe that most of the same-kiosk rides under 5 minutes or so likely shouldn’t count, but only culled the ones that were one minute long”.

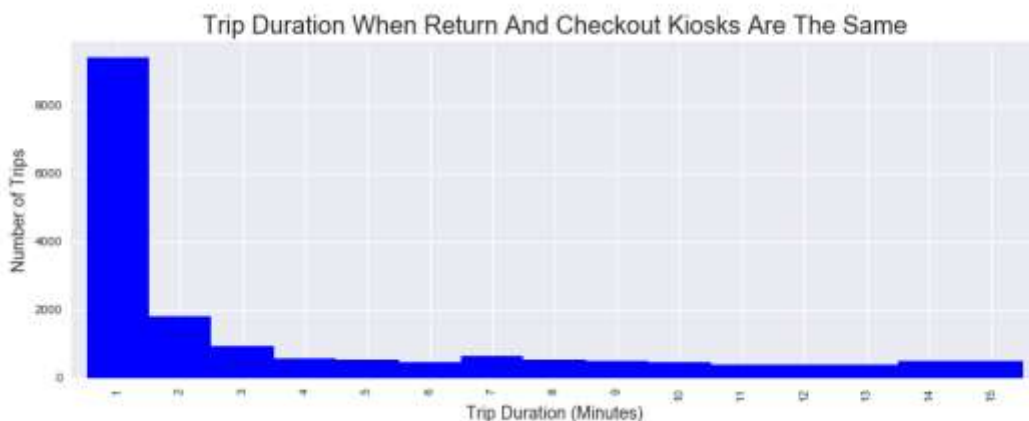


FIGURE 1: TRIP DURATION WHEN CHECKOUT AND RETURN KIOSKS ARE THE SAME

There were 6574 rows in the Trips dataset that had kiosk names not in the Kiosk Master List. These 6574 rows were removed accordingly.

Removing the 9,954 rows with a trip duration of 1 minute and 6574 rows with invalid kiosk names resulted in **394,431 Denver B-cycle rides in 2016**.

## Distance Traveled

To estimate the distance between checkout and return kiosks when they are the same, Tyler’s method of using the “average speed of all the other rides (nominal distance ridden divided by the duration), and then applying this average speed to the same-kiosk trip durations” was adopted. This resulted in **670,802 miles ridden in 2016**.

## **Most Popular and Least Popular Checkout and Return Kiosks**

### **Most Popular**

The following ten kiosks were the most popular checkout kiosks by number of total bike checkouts in 2016.

<u>Checkout Kiosk</u>	<u>Number of Checkouts</u>
16th & Wynkoop	11174
16th & Broadway	11116
1350 Larimer	10837
18th & California	9865
1550 Glenarm	9441
18th & Arapahoe	8531
20th & Chestnut	8240
13th & Speer	8228
REI	8218
16th & Little Raven	8198

The following ten kiosks were the most popular return kiosks by number of total bike checkouts in 2016.

<u>Return Kiosk</u>	<u>Number of Checkouts</u>
16th & Wynkoop	11289
1350 Larimer	10920
16th & Broadway	10870
18th & California	9863
1550 Glenarm	9501
18th & Arapahoe	8549
20th & Chestnut	8356
REI	8284
13th & Speer	8272
16th & Little Raven	8267

### **Least Popular**

The following ten kiosks were the least popular checkout kiosks by number of total bike checkouts in 2016.

<u>Checkout Kiosk</u>	<u>Number of Checkouts</u>
Pepsi Center	1795
32nd & Julian	1755
25th & Lawrence	1736
Colfax & Garfield	1725
4th & Walnut	1663
Decatur Federal Light Rail	1508
Denver Zoo	1490
Colfax & Gaylord	1421
17th & Curtis	615
39th & Fox	332

The following ten kiosks were the least popular return kiosks by number of total bike checkouts in 2016.

<u>Return Kiosk</u>	<u>Number of Checkouts</u>
21st & Market	1795
32nd & Julian	1767
25th & Lawrence	1758
Colfax & Garfield	1743
4th & Walnut	1686
Decatur Federal Light Rail	1537
Denver Zoo	1468
Colfax & Gaylord	1433
17th & Curtis	632
39th & Fox	345

## Map of Station Popularity

### Checkout Kiosks

The use of Tableau aided in the creation of the following map showing the popularity of the various Checkout Kiosks (Figure 2). The size of the circle is proportional to the number of checkouts from that kiosk in 2016.

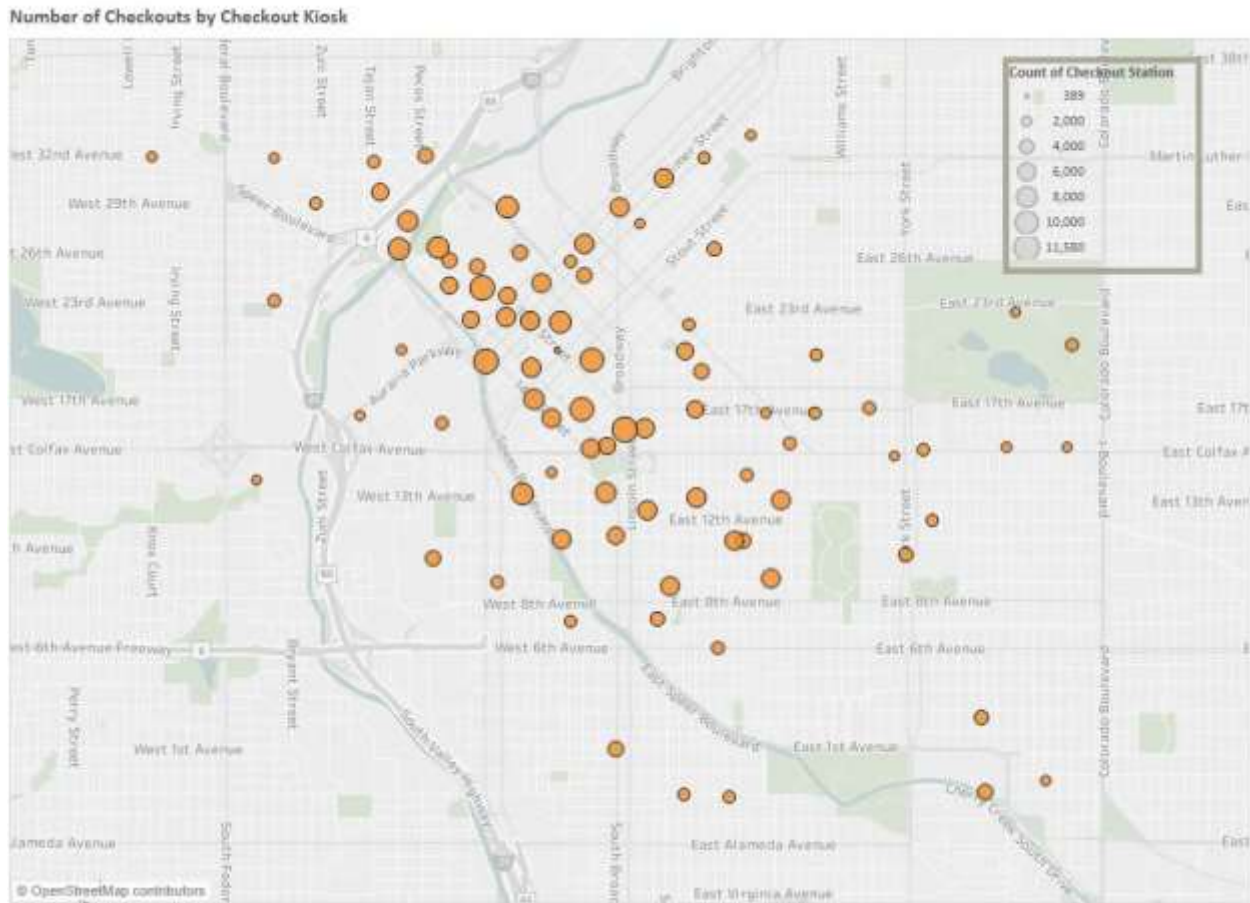


FIGURE 2: CHECKOUT KIOSK LOCATIONS AND NUMBER OF CHECKOUTS IN 2016



## Return Kiosks

Similarly, the use of Tableau aided in the creation of the following map showing the popularity of the various Return Kiosks (Figure 3). The size of the circle corresponds to the number of checkouts returned to that kiosk in 2016.

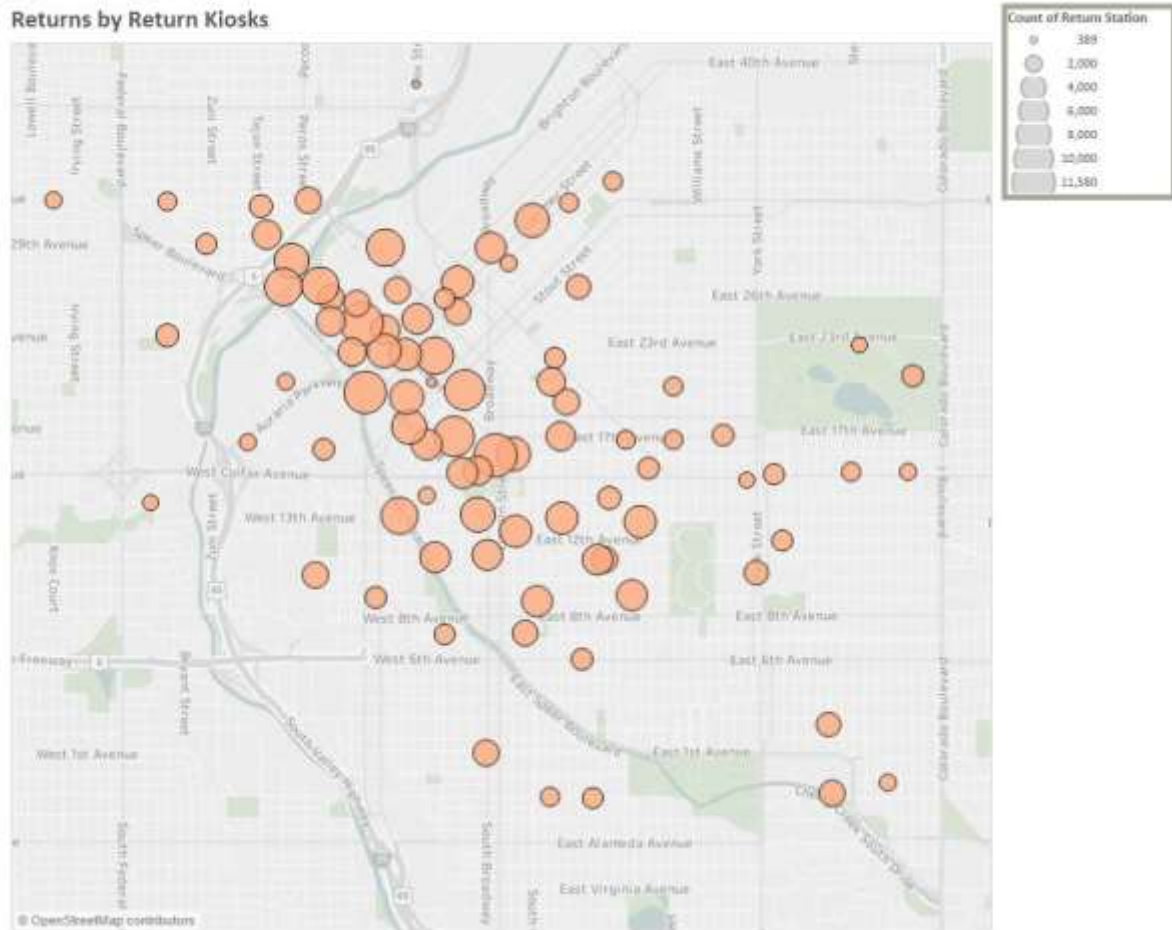


FIGURE 3: RETURN KIOSK LOCATIONS AND NUMBER OF RETURNS IN 2016



## Checkouts per Membership Type

Denver B-cycle has a number of different membership passes. The following were the top ten by number of checkouts in 2016 (Figure 4).

24-hour Kiosk Only (Denver B-cycle)	85680
Annual Plus (Denver B-cycle)	82202
Annual (Denver B-cycle)	80093
Not Applicable	56250
Monthly (Denver B-cycle)	52811
Denver B-cycle Founder (Denver B-cycle)	17675
Single Ride (Denver B-cycle)	16291
Republic Rider (Annual) (Boulder B-cycle)	5930
Flex Pass (Denver B-cycle)	3507
Republic Rider (Boulder B-cycle)	1229

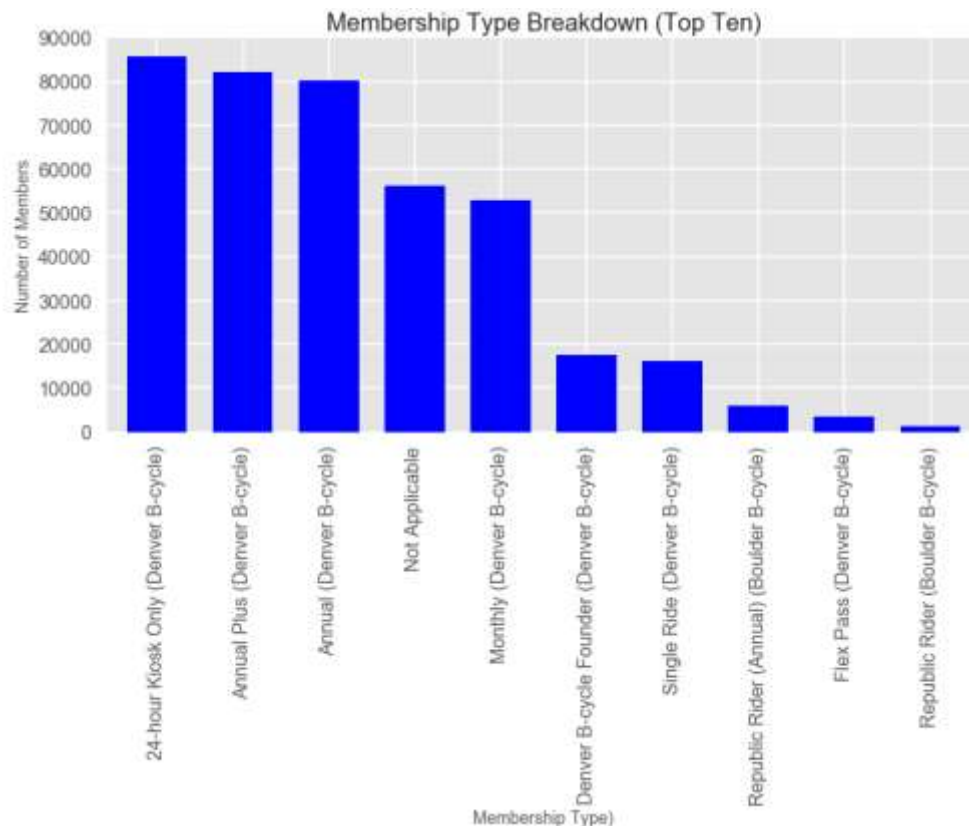


FIGURE 4: NUMBER OF CHECKOUTS BY MEMBERSHIP TYPE IN 2016

## Ridership by Calendar and Clock Variables

### Ridership by Hour

Bike checkout time is probably the most important attribute in the Trips dataset. Each checkout time was converted into its integer hour. For example, 7:02 AM or 7:59 AM would be converted to an integer of 7. In this way, total number of checkouts could be aggregated for the year and plotted against their hours of the day, as shown in Figure 5.

It appears that the highest number of checkouts occur between 4 PM and 5 PM with ridership increasing steadily from 10 AM onwards.

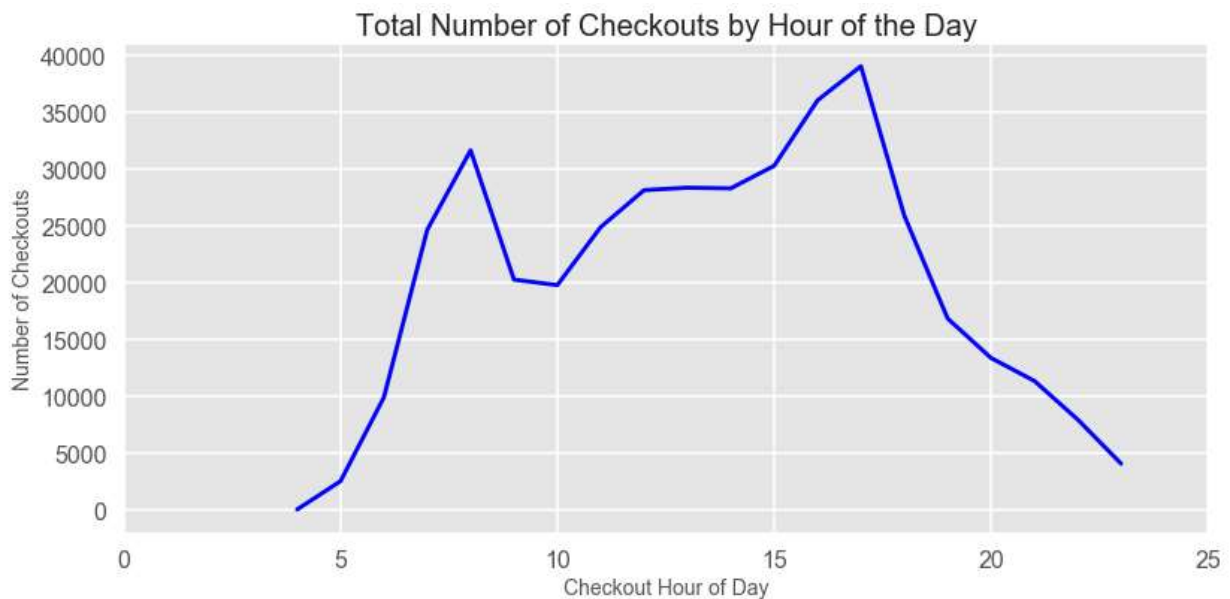


FIGURE 5: NUMBER OF CHECKOUTS BY HOUR IN 2016

Figure 6 shows the average distance ridden by the hour of the day in 2016. More distance is covered during the 10 AM period and declining steadily after 3 PM.

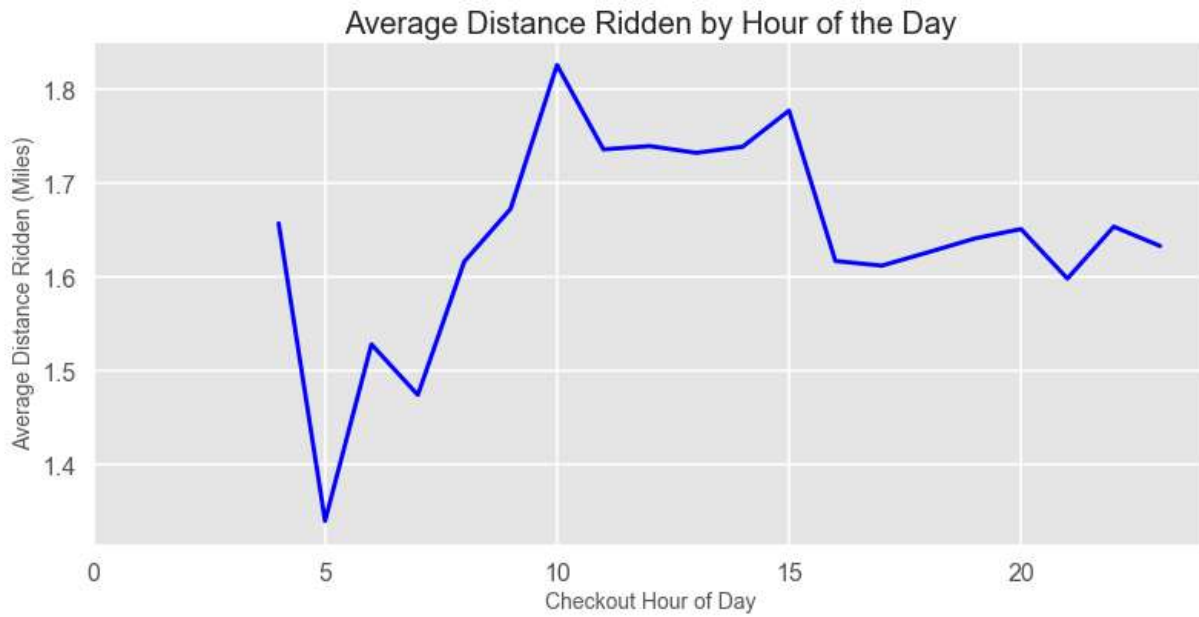


FIGURE 6: ESTIMATED AVERAGE MILES RIDDEN BY HOUR OF CHECKOUT IN 2016

## Ridership by Hour and Weekday

Figure 7 shows that weekday ridership patterns are similar. On the other hand weekend ridership demonstrate a busy afternoon (between 12 PM and 3 PM)

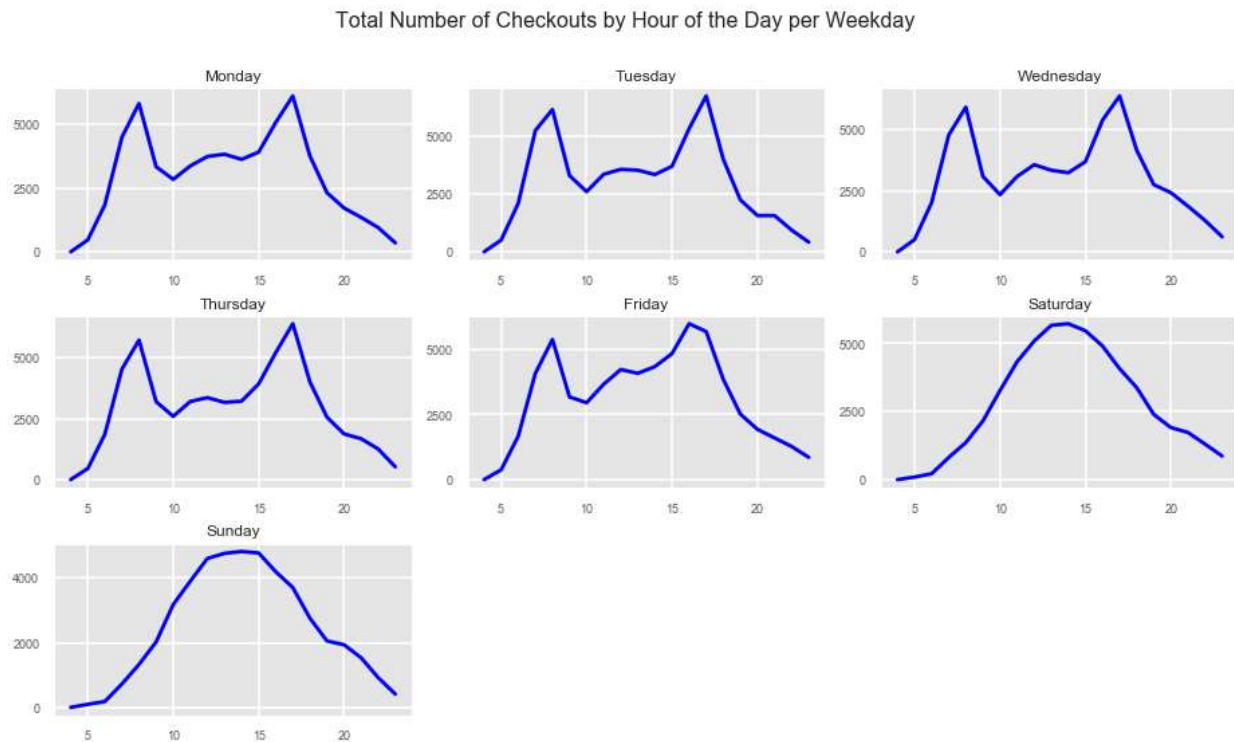


FIGURE 7: CHECKOUTS BY HOUR OF DAY PER WEEKDAY IN 2016

## Ridership by Month

Monthly checkouts, as shown in Figure 8, suggest high ridership during the summer months and low ridership during the winter months.

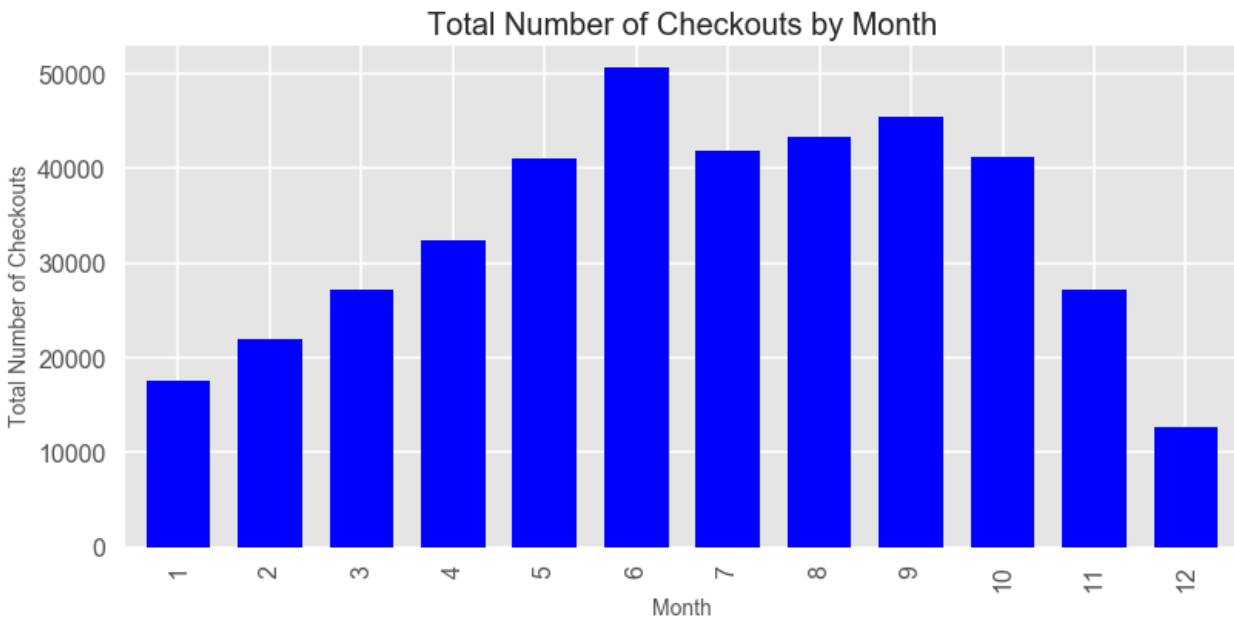


FIGURE 8: TOTAL CHECKOUTS BY MONTH IN 2016

## Merging with Weather

It is highly likely that weather plays a very important role in bike ridership and bike checkout times. This was shown in the previous plots on total checkouts per hour of the day, by weekday, and by month. To verify this, weather data obtained from [Dark Sky API](#) was merged with the Trips dataset and several graphs plotted to visualize the relationships.

## Checkouts vs. Daily and Hourly Temperature

Figure 9 shows the total number of checkouts against maximum and minimum daily temperature. It clearly suggests that ridership increases as the temperature increases and vice-versa.

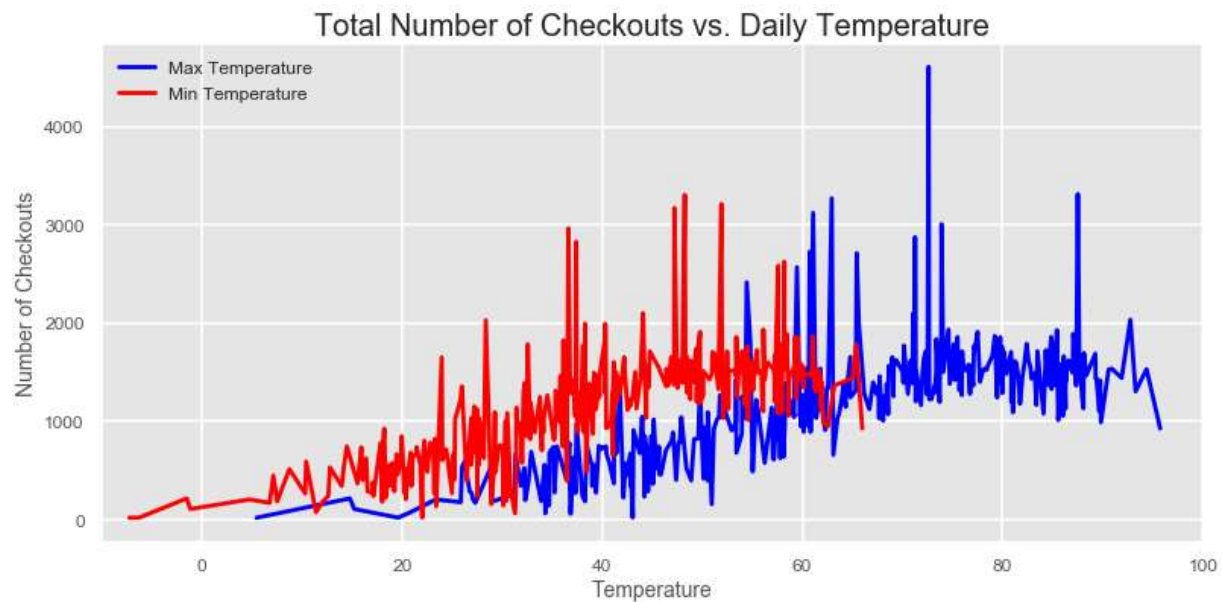


FIGURE 9: TOTAL CHECKOUTS BY DAILY TEMPERATURE IN 2016

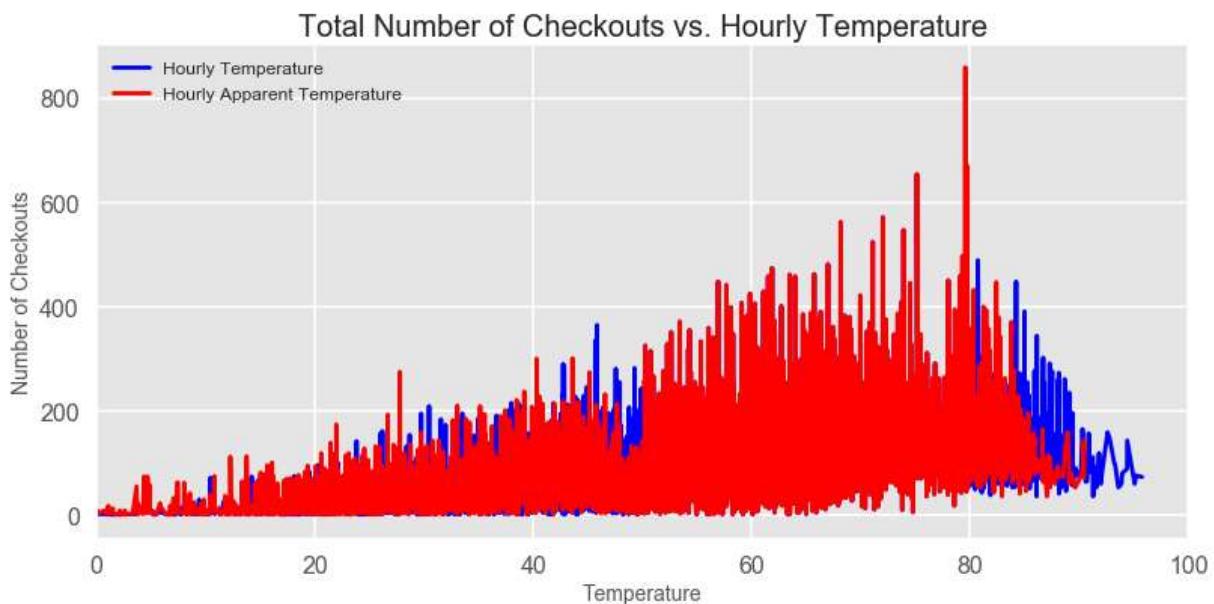


FIGURE 10: TOTAL CHECKOUTS BY HOURLY TEMPERATURE IN 2016

Apparent temperature, as defined by Dark Sky, is “apparent (or “feels like”) temperature in degrees Fahrenheit”. It appears to have a subtle effect on bike ridership as shown in Figure 11.

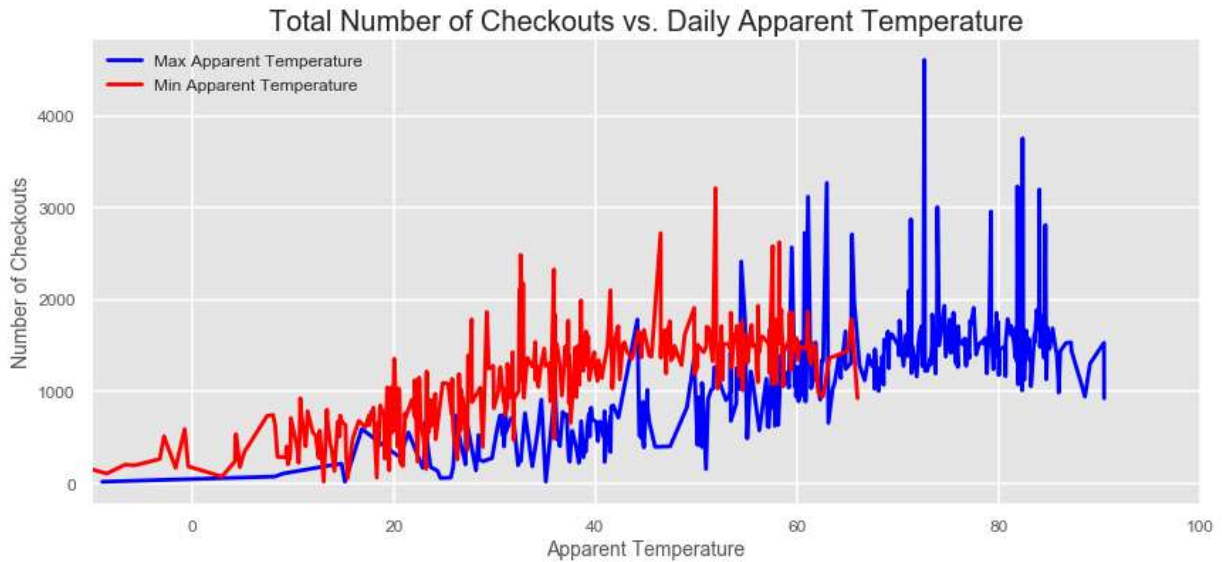


FIGURE 11: TOTAL CHECKOUTS BY DAILY APPARENT TEMPERATURE IN 2016



## Checkouts vs. Daily and Hourly Cloud Cover

Dark Sky defines Cloud Cover as “the percentage of sky occluded by clouds, between 0 and 1, inclusive”. Figures 12 and 13 show the total number of checkouts against daily and hourly cloud cover, respectively. They clearly suggest that ridership is highest as the cloud cover stays at around 0.15.

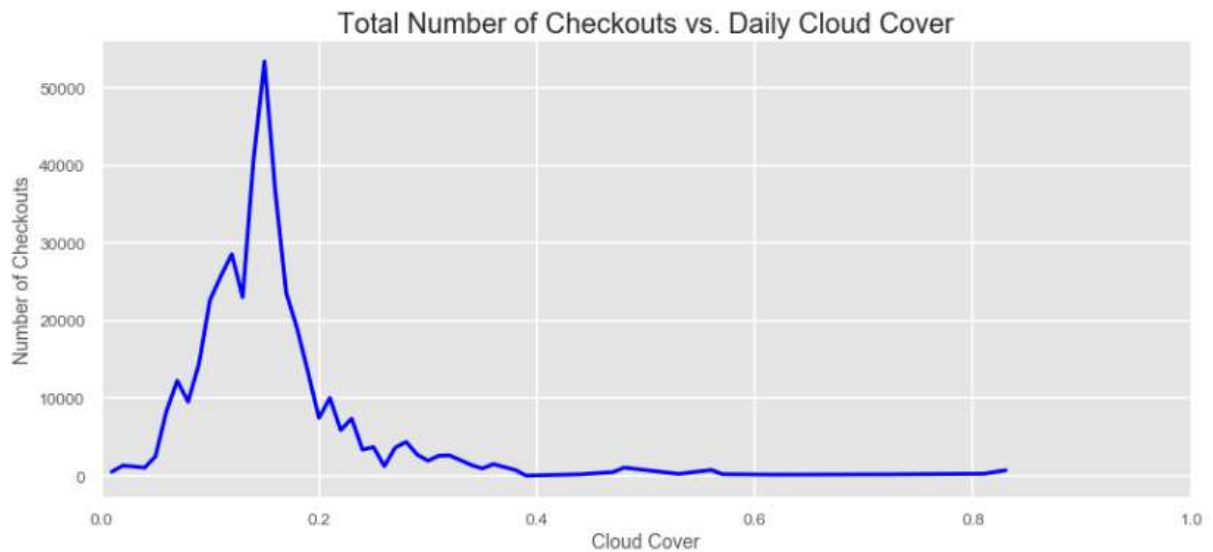


FIGURE 12: TOTAL CHECKOUTS BY DAILY CLOUD COVER IN 2016

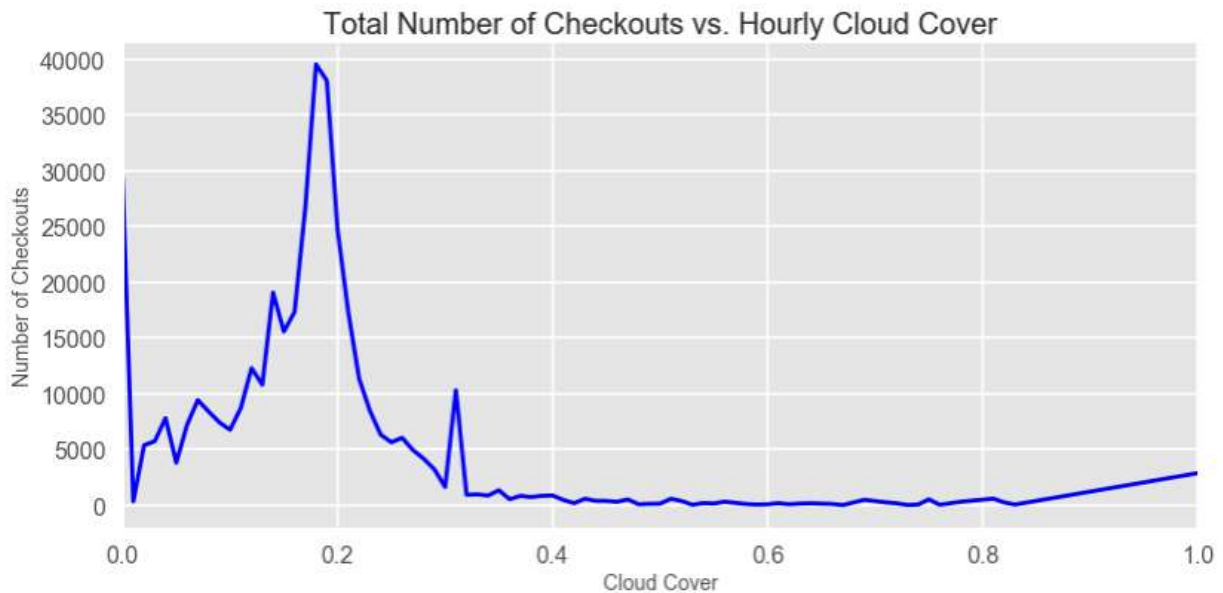


FIGURE 13: TOTAL CHECKOUTS BY DAILY CLOUD COVER IN 2016

## Checkouts vs. Daily and Hourly Wind Speed

Wind speed is reported in miles per hour. As shown in Figures 14 and 15, wind speed significantly affects ridership when it hits 8 miles per hour and higher.

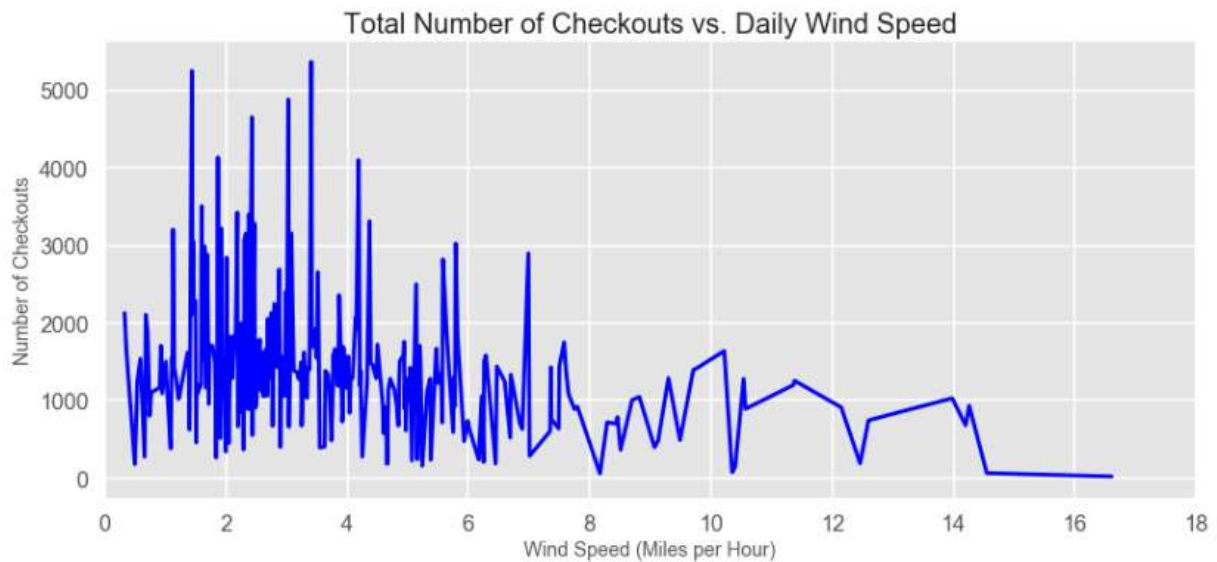


FIGURE 14: TOTAL CHECKOUTS BY DAILY WIND SPEED IN 2016

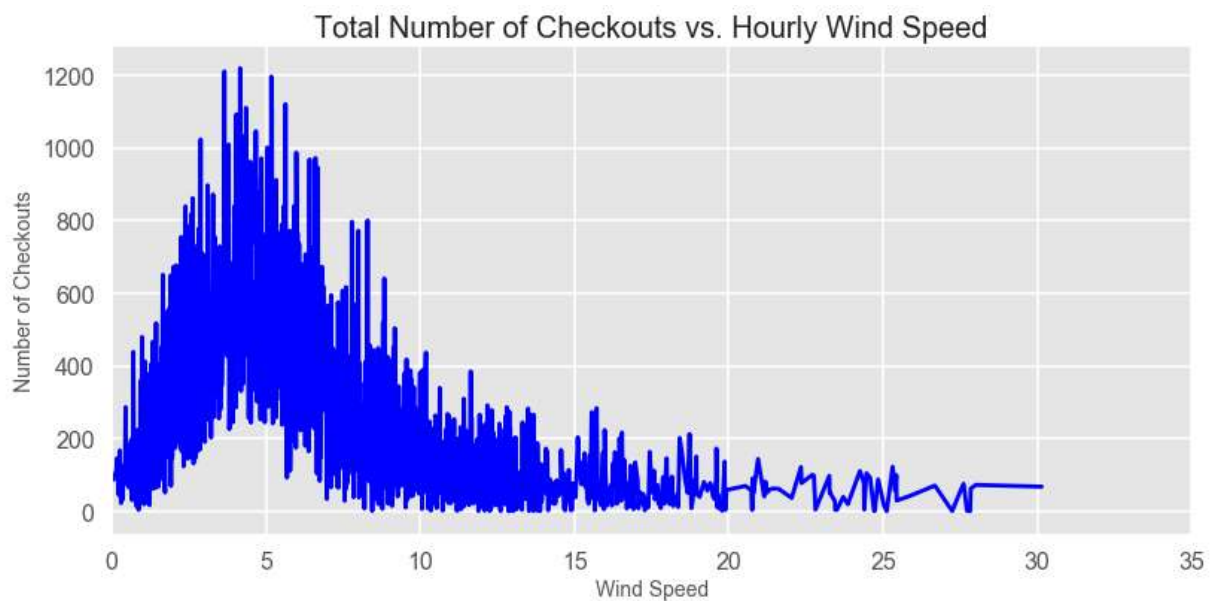


FIGURE 15: TOTAL CHECKOUTS BY HOURLY WIND SPEED IN 2016

## Checkouts vs. Daily and Hourly Humidity

Humidity is defined by Dark Sky as “relative humidity, between 0 and 1. Figures 16 and 17 show decreased ridership at higher humidity levels. The hourly humidity plot provides a better resolution than the daily humidity plot.

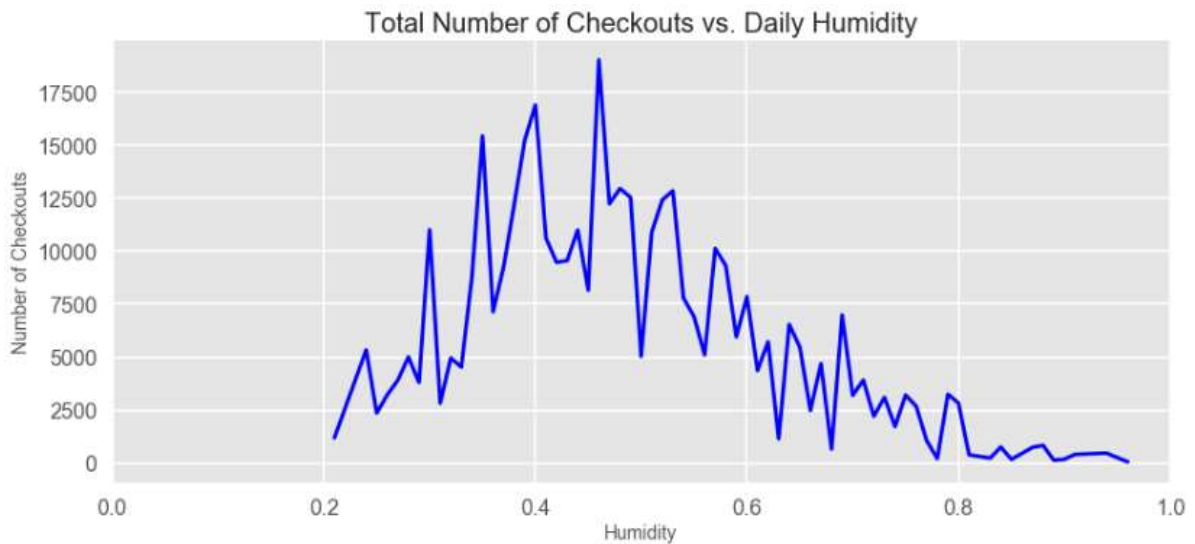


FIGURE 16: TOTAL CHECKOUTS BY DAILY HUMIDITY IN 2016

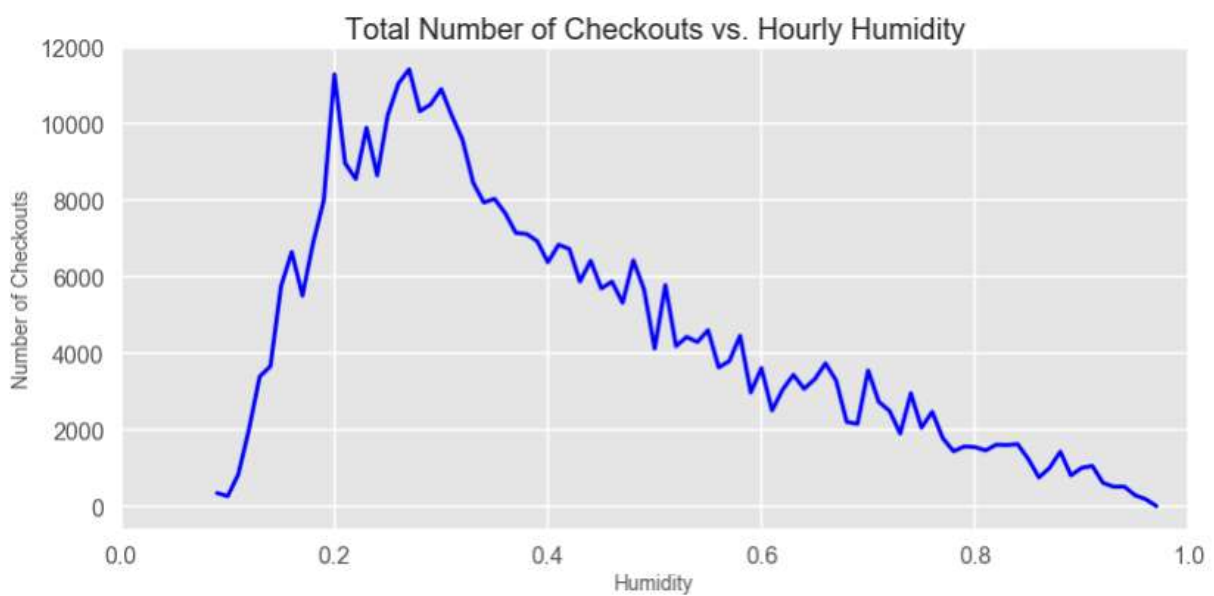


FIGURE 17: TOTAL CHECKOUTS BY HOURLY HUMIDITY IN 2016

## Checkouts vs. Daily and Hourly Visibility

Visibility is measured in miles and capped at 10 miles, according to Dark Sky. As Figures 18 and 19 show, ridership peaks when visibility is at 10 miles.

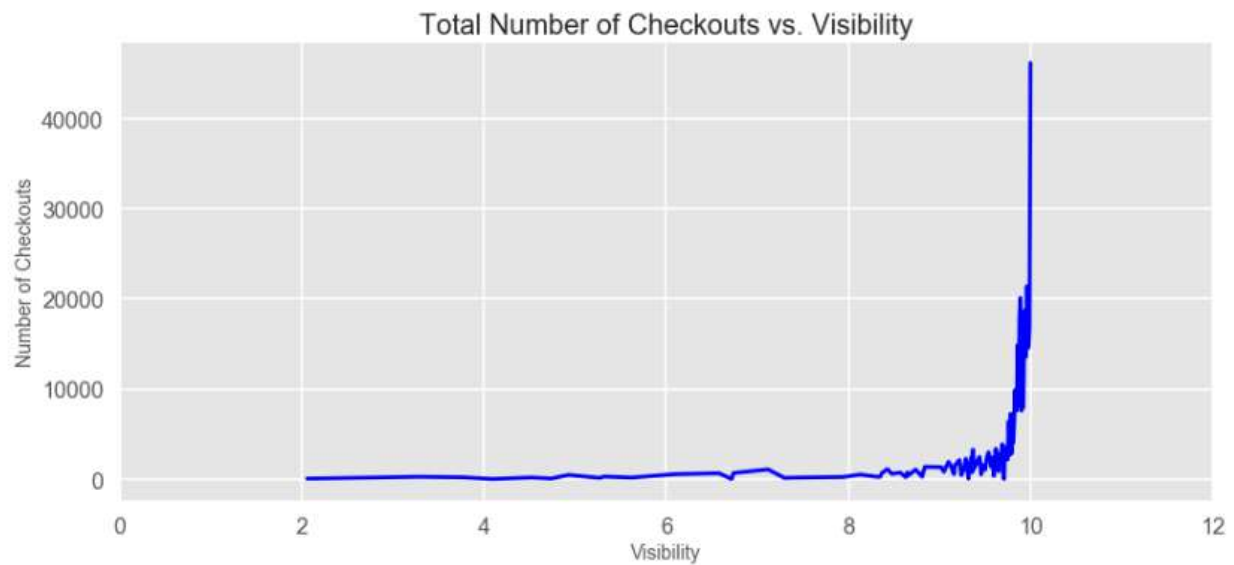


FIGURE 18: TOTAL CHECKOUTS BY DAILY VISIBILITY IN 2016

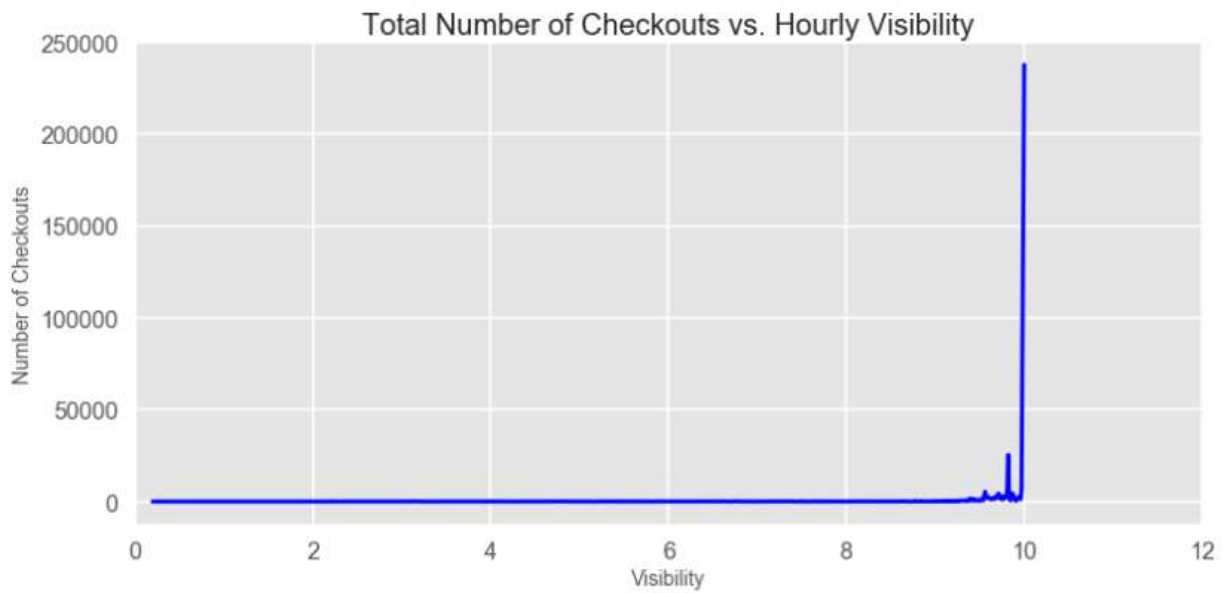


FIGURE 19: TOTAL CHECKOUTS BY HOURLY VISIBILITY IN 2016

## Days with Highest/Lowest Ridership

Another interesting data discovery was the fact that Saturdays and Sundays had the highest and lowest ridership depending upon the weather. In his study, Tyler suggests that this may be due to “weekend warriors’ who rent B-cycles for pleasure and are highly affected by the weather in their decision to ride”. This may well be the case.

<u>Checkout WeekDay</u>	<u>Checkout Date</u>	<u>temperatureMax</u>	<u>temperatureMin</u>	<u>Checkouts</u>
6	2016-05-29	71.090	44.100	2100
5	2016-05-28	65.650	40.330	1990
4	2016-06-03	74.600	56.120	1933
2	2016-06-15	85.430	51.980	1927
5	2016-05-21	77.510	49.790	1909
0	2016-06-27	87.060	58.440	1889
5	2016-06-25	79.230	61.040	1868
	2016-06-04	75.500	53.410	1857
3	2016-06-23	84.860	59.280	1857
4	2016-09-02	79.770	59.500	1855

<u>Checkout WeekDay</u>	<u>Checkout Date</u>	<u>temperatureMax</u>	<u>temperatureMin</u>	<u>Checkouts</u>
5	2016-12-24	50.960	28.940	154
6	2016-04-17	34.710	30.140	140
	2016-01-31	31.260	23.430	133
2	2016-12-07	15.250	-1.110	105
1	2016-02-02	20.870	11.430	72
5	2016-04-16	34.340	31.310	61
6	2016-12-25	36.860	25.290	56
2	2016-03-23	43.070	22.040	18
6	2016-12-18	19.640	-6.220	17
5	2016-12-17	5.490	-7.220	16

## Part 2: Regression Modeling

In this section various regression models were used to test and train the Trips data that was merged with the weather data to try to predict the checkout time based on weather conditions.

The following regression models were used in this study:

- Decision Tree Regression
- Linear Regression
- Lasso Regression
- Ridge Regression
- Random Forest Regression
- Extra Trees Regression
- Nearest Neighbors Regression
- Bayesian Ridge Regression

Prior to applying the models a feature correlation was performed on all the features to see if any of the features were highly correlated to one another. As shown in Figure 20, Temperature and Apparent\_Temperature were highly correlated suggesting that one of them could be removed from the features in the model application.

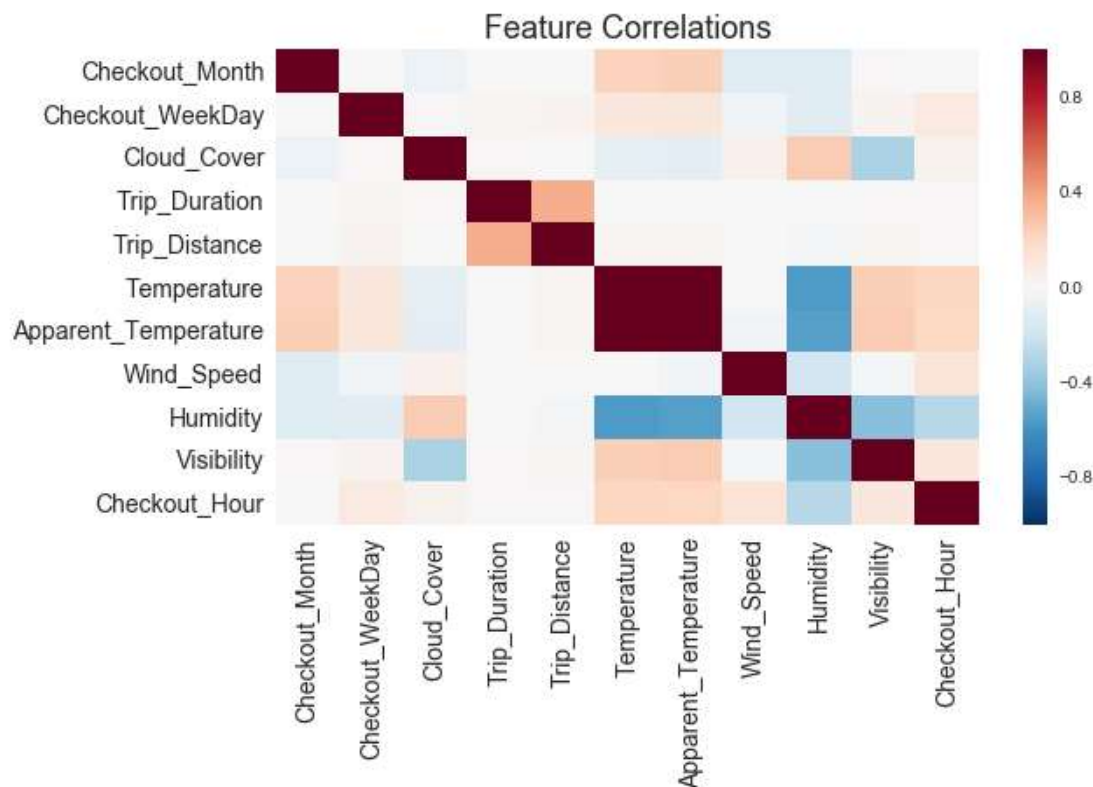


FIGURE 20: FEATURE CORRELATION

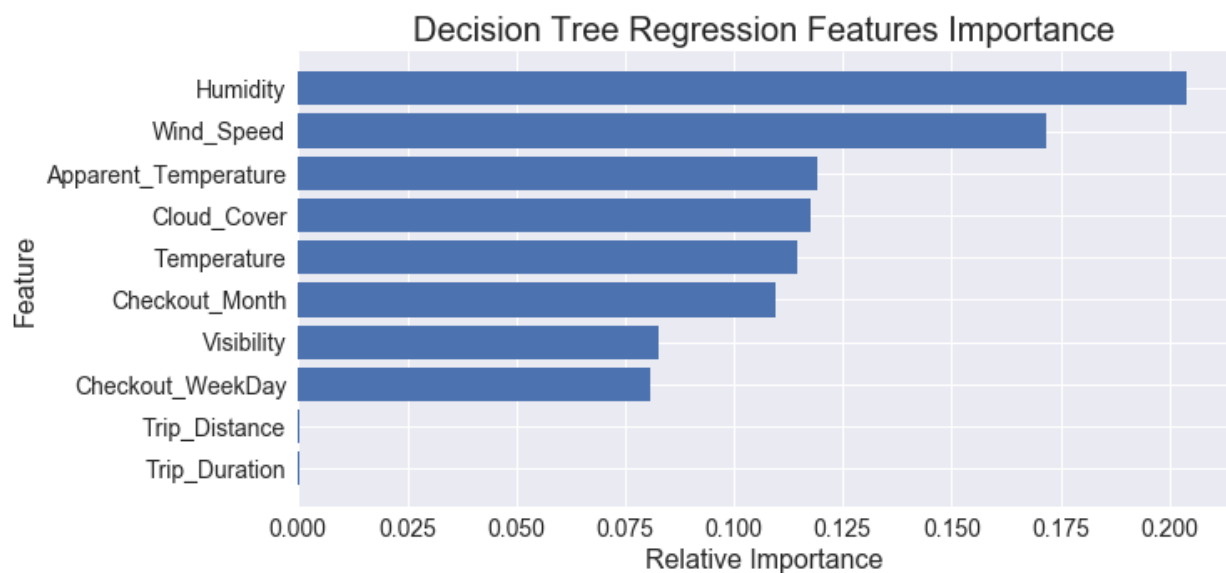
## Regression Modeling – All Features

### Decision Tree Regression

Training Set Score: 1.000

Testing Set Score: 0.997

	Decision Tree Regression
R Squared	0.99703



### Linear Regression

Training Set Score: 0.106

Testing Set Score: 0.102

	Linear Regression
R Squared	0.102423

### Lasso Regression

Training Set Score: 0.106

Testing Set Score: 0.102



	Lasso Regression
R Squared	0.102426

## Ridge Regression

Training Set Score: 0.106

Testing Set Score: 0.102

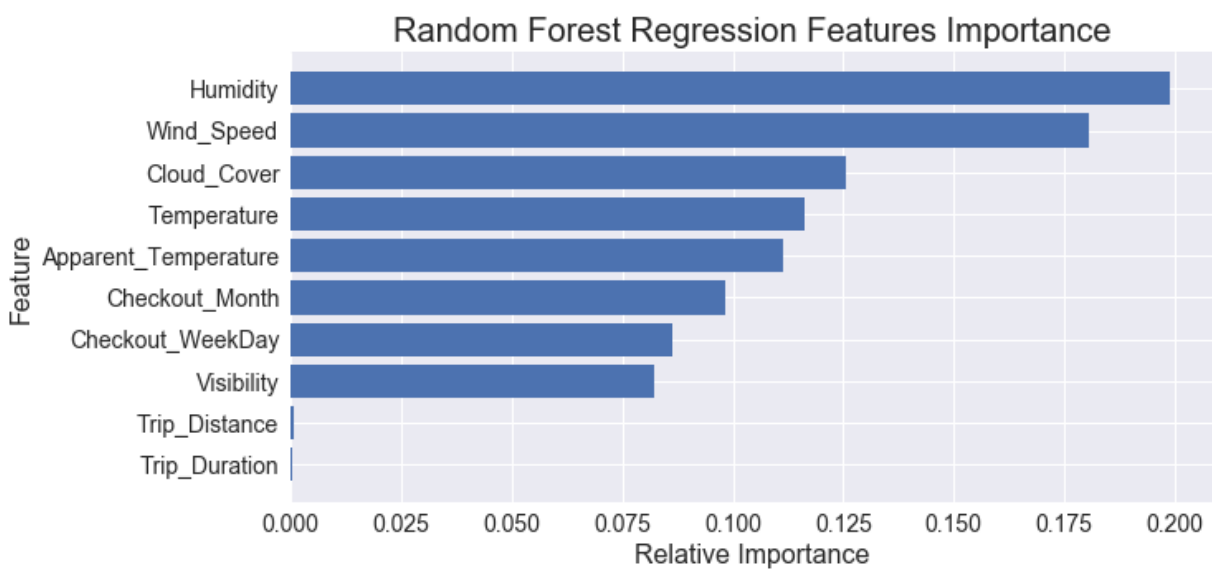
	Ridge Regression
R Squared	0.102448

## Random Forest Regression

Training Set Score: 1.000

Testing Set Score: 0.998

	Random Forest Regression
R Squared	0.99767

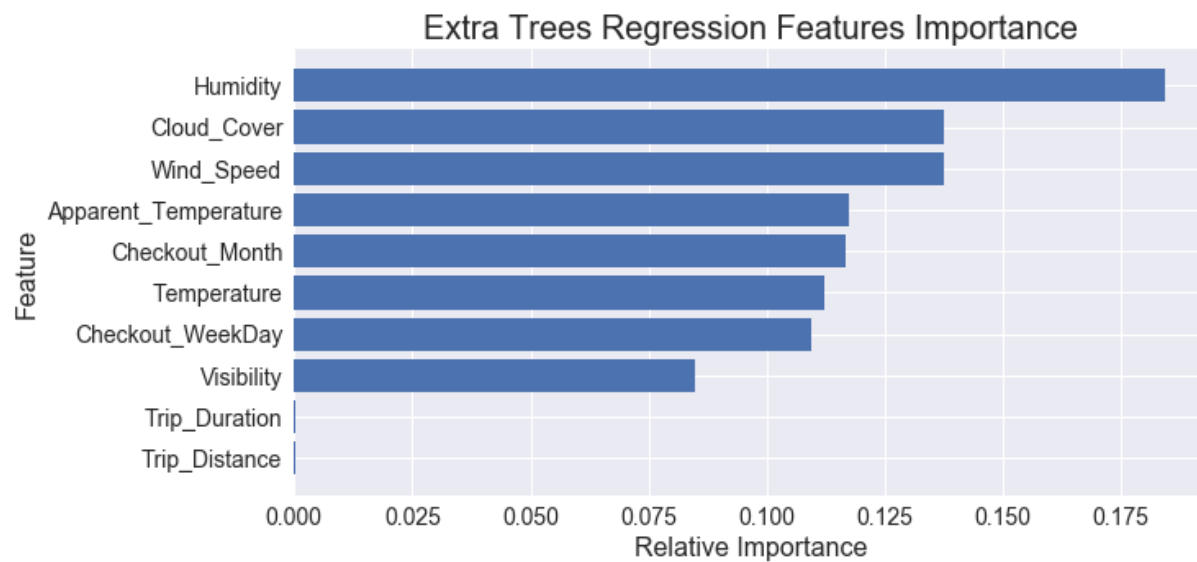


## Extra Trees Regression

Training Set Score: 1.000

Testing Set Score: 0.999

	Extra Trees Regression
R Squared	0.999004



## Nearest Neighbors Regression

Training Set Score: 0.913

Testing Set Score: 0.653

	Nearest Neighbors Regression
R Squared	0.652742

## Bayesian Ridge Regression

Training Set Score: 0.106

Testing Set Score: 0.102

	Bayesian Ridge Regression
R Squared	0.102425

## Regression Modeling Summary – All Features

	Decision Tree Regression	Linear Regression	Lasso Regression	Ridge Regression	Random Forest Regression	Extra Trees Regression	Nearest Neighbors Regression	Bayesian Ridge Regression
R Squared	0.99703	0.102423	0.102426	0.102448	0.99767	0.999004	0.652742	0.102425

## Regression Modeling - Selected Features

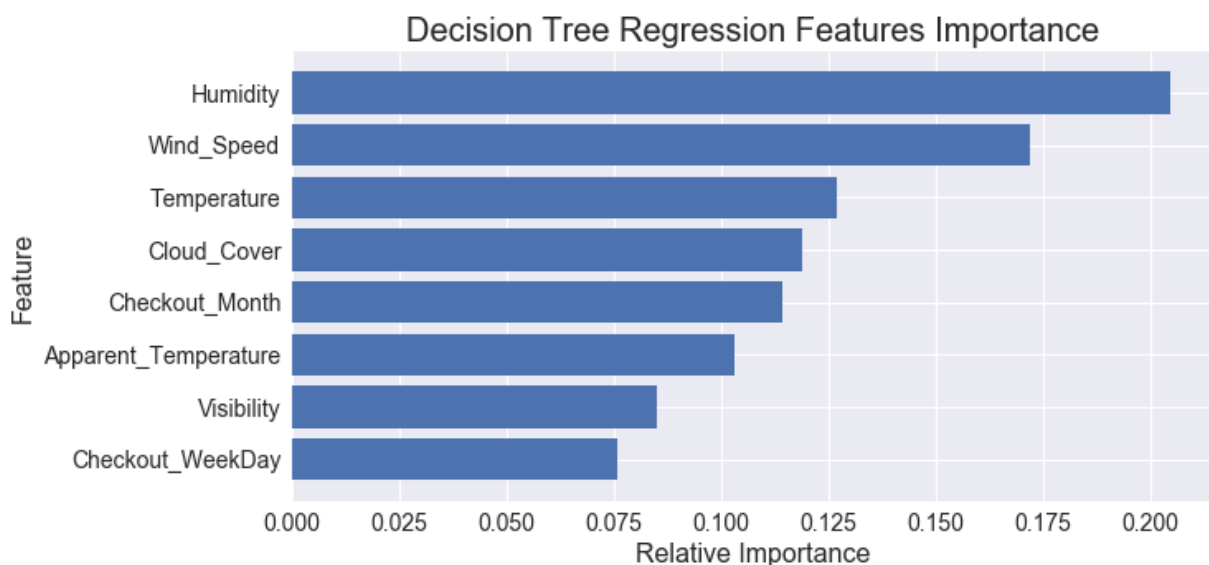
The Trip\_Duration and Trip\_Distance features were removed from the original list of features. These reduced number of features were then used in the regression models to see if the train and test accuracy could be improved.

## Decision Tree Regression

Training Set Score: 1.000

Testing Set Score: 0.998

	Decision Tree Regression
R Squared	0.997807



## Linear Regression

Training Set Score: 0.106

Testing Set Score: 0.102

	Linear Regression
R Squared	0.102364

## Lasso Regression

Training Set Score: 0.106

Testing Set Score: 0.102

	Lasso Regression
R Squared	0.102387

## Ridge Regression

Training Set Score: 0.106

Testing Set Score: 0.102

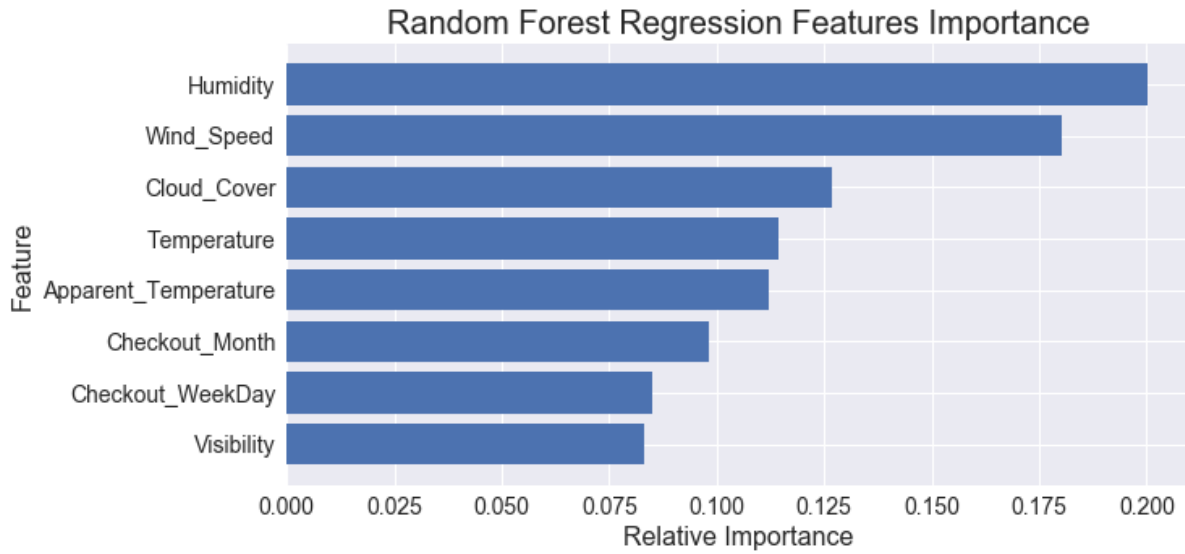
	Ridge Regression
R Squared	0.102387

## Random Forest Regression

Training Set Score: 1.000

Testing Set Score: 0.998

	Random Forest Regression
R Squared	0.997989

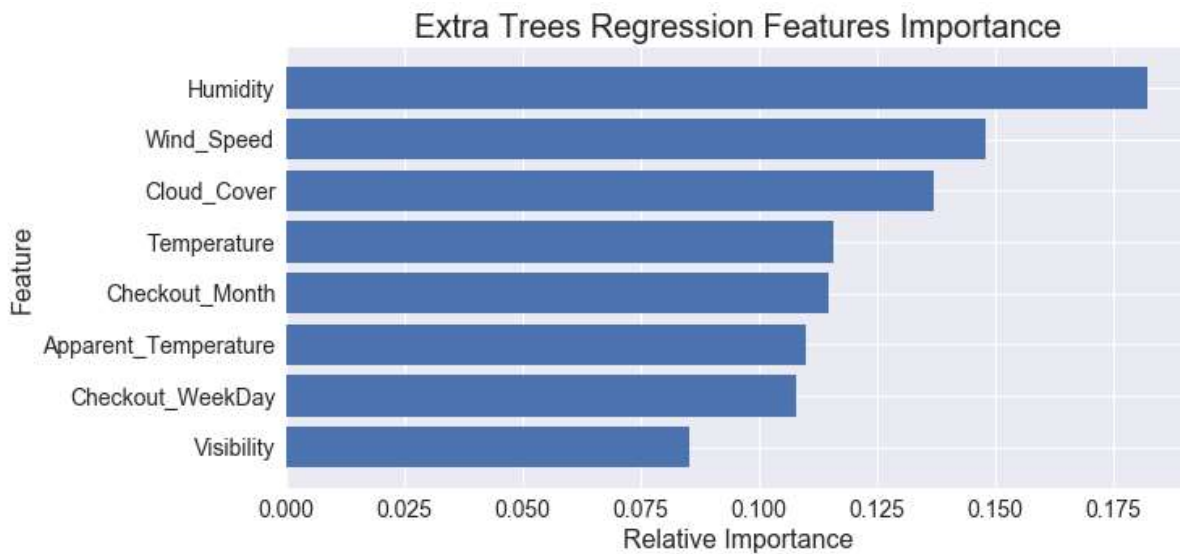


## Extra Trees Regression

Training Set Score: 1.000

Testing Set Score: 0.999

	Extra Trees Regression
R Squared	0.999011



## Nearest Neighbors Regression

Training Set Score: 1.000

Testing Set Score: 0.997

	Nearest Neighbors Regression
R Squared	0.997435

## Bayesian Ridge Regression

Training Set Score: 0.106

Testing Set Score: 0.102

	Bayesian Ridge Regression
R Squared	0.102365

## Regression Modeling Summary – Selected Features

	Decision Tree	Linear Regression	Lasso Regression	Ridge Regression	Random Forest Regression	Extra Trees Regression	Nearest Neighbors Regression	Bayesian Ridge Regression
R Squared	0.997807	0.102364	0.102387	0.102387	0.997989	0.999011	0.997435	0.102365

## Part 3: Classification Modeling

In this section various classification models were used to test and train the Trips data that was merged with the weather data to try to predict the checkout time based on weather conditions.

The following classification models were used in this study:

- Decision Tree Classification
- Linear (Logistic) Classification
- Random Forest Classification
- Extra Trees Classification
- Naïve Bayes Classification
- Nearest Neighbors Classification

	Decision Tree Classification	Logistic Classification	Random Forest Classification	Extra Trees Classification	Naive Bayes Classification	Nearest Neighbors Classification
Accuracy	0.998865	0.171393	0.999125	0.999287	0.133939	0.822490
F1 (macro)	0.996534	0.091614	0.997270	0.998073	0.058921	0.784425
F1 (micro)	0.998865	0.171393	0.999125	0.999287	0.133939	0.822490
Precision (macro)	0.997204	0.128753	0.997803	0.998689	0.110323	0.796228
Precision (micro)	0.998865	0.171393	0.999125	0.999287	0.133939	0.822490
Recall (macro)	0.995891	0.117052	0.996747	0.997472	0.104769	0.775852
Recall (micro)	0.998865	0.171393	0.999125	0.999287	0.133939	0.822490



	Decision Tree Classification	Logistic Classification	Random Forest Classification	Extra Trees Classification	Naive Bayes Classification	Nearest Neighbors Classification
Accuracy	0.999330	0.170276	0.999293	0.999280	0.148475	0.999287
F1 (macro)	0.998256	0.090561	0.997880	0.997823	0.088862	0.997960
F1 (micro)	0.999330	0.170276	0.999293	0.999280	0.148475	0.999287
Precision (macro)	0.998903	0.121817	0.998247	0.998240	0.117142	0.998502
Precision (micro)	0.999330	0.170276	0.999293	0.999280	0.148475	0.999287
Recall (macro)	0.997624	0.116308	0.997522	0.997415	0.115514	0.997431
Recall (micro)	0.999330	0.170276	0.999293	0.999280	0.148475	0.999287