

Analysis of Hubway Bike Rental Data

Subhabrata Mukherjee

Dec 2017

Introduction

The objective of this analysis:

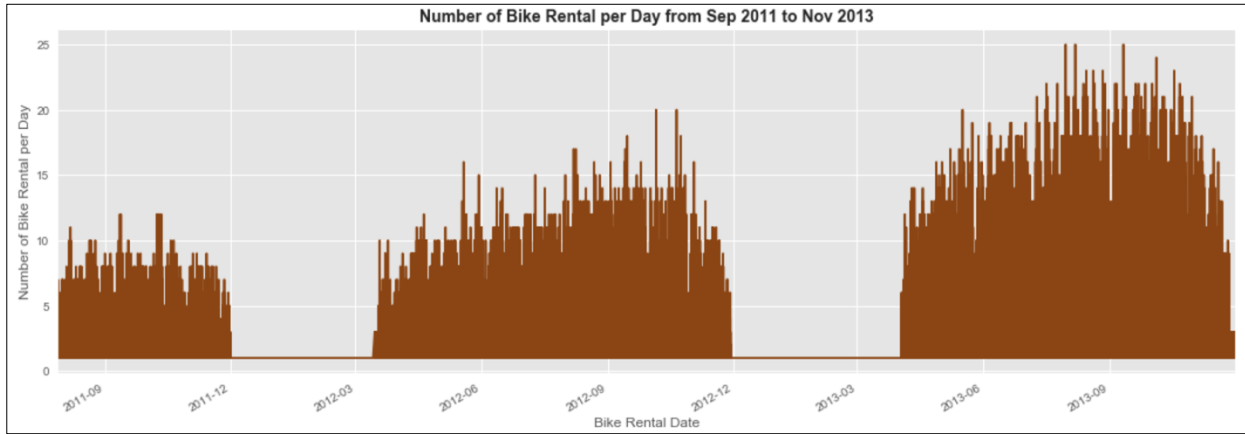
- ▶ Analyze bike rental data to know the demand for bike on weekdays
- ▶ Identify most prominent time of the day having highest demand for bike
- ▶ Create predictive model of Hubway bike demand using machine learning

Data sets

- ▶ Hubway Trip Data from 2011 to 2013
 - ▶ Date, time, origin and destination stations, plus the bike number
- ▶ Weather Data from 2012 to 2013
 - ▶ Daily temperature, humidity, wind speed

DATA EXPLORATION

Daily Bike Rental by Year



- ❑ The bike share program is gaining popularity over time.
- ❑ Weather is an important factor in daily demand of bike rental

PREDICTIVE MODELS

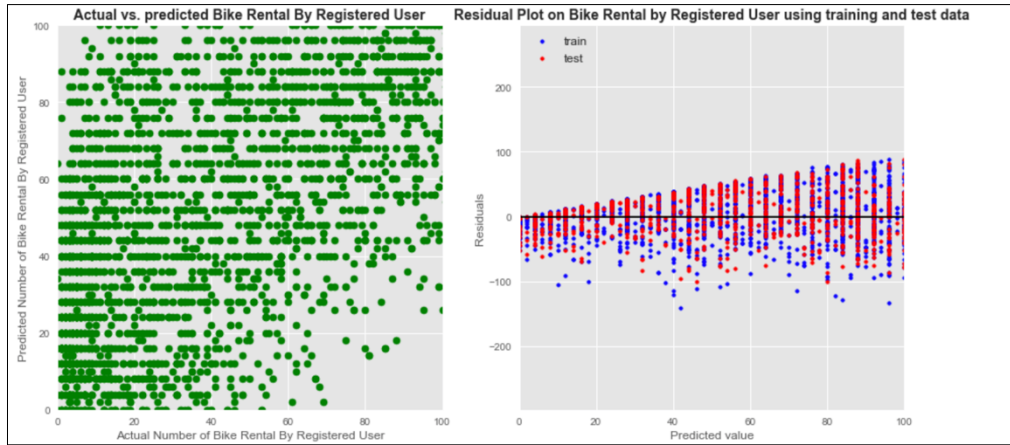
Linear Regression



Selected 12 features	
TEMP	ATEMP
HUMIDITY	SPEED
RAIN	SNOW
CASUAL	REGISTERED

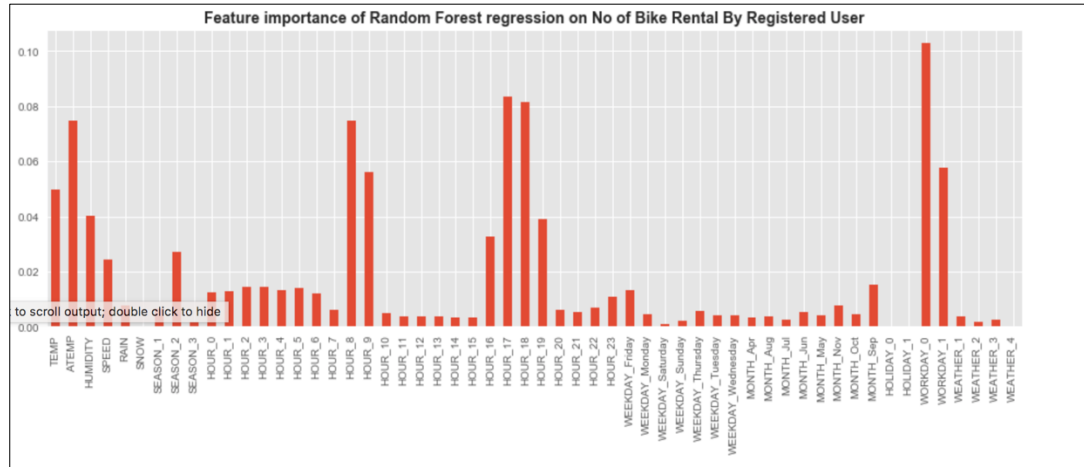
Correlation matrix shows some relationships but none strong enough to skew the model

Linear Regression - Registered User



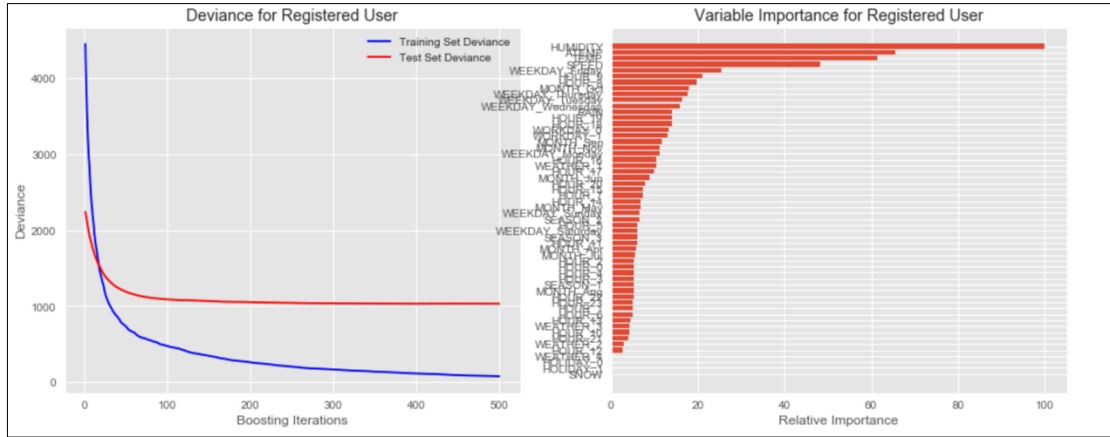
- ❑ Model explains 60% of the variance in the data - a moderate value for fit
- ❑ Training data has a MSE = 3720.20 and test data has a MSE = 3776.83
- ❑ Linear regression is not a good choice for this data set

Random Tree Regression - Registered User



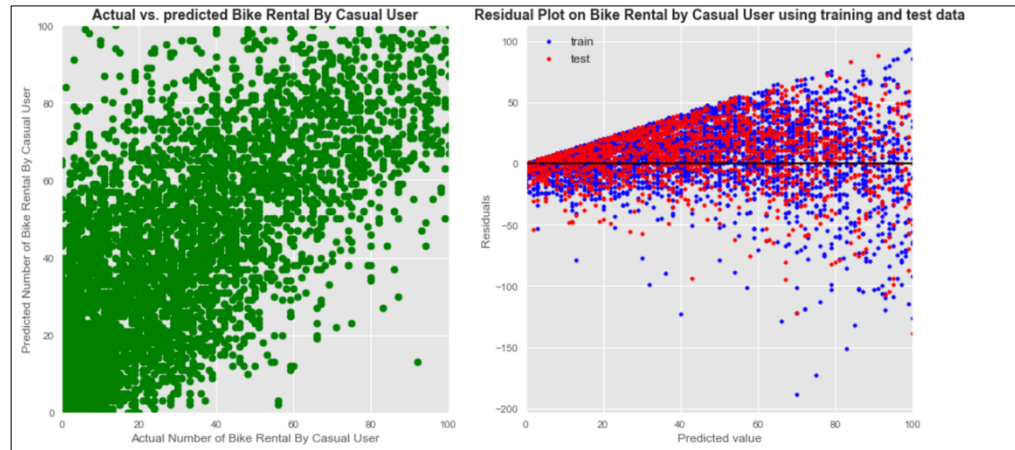
- ❑ Model explains 81% of data's variance but does not generalize well to the test set
- ❑ WORKDAY and HOUR are most important features during training

Gradient Boosting Regression - Registered User



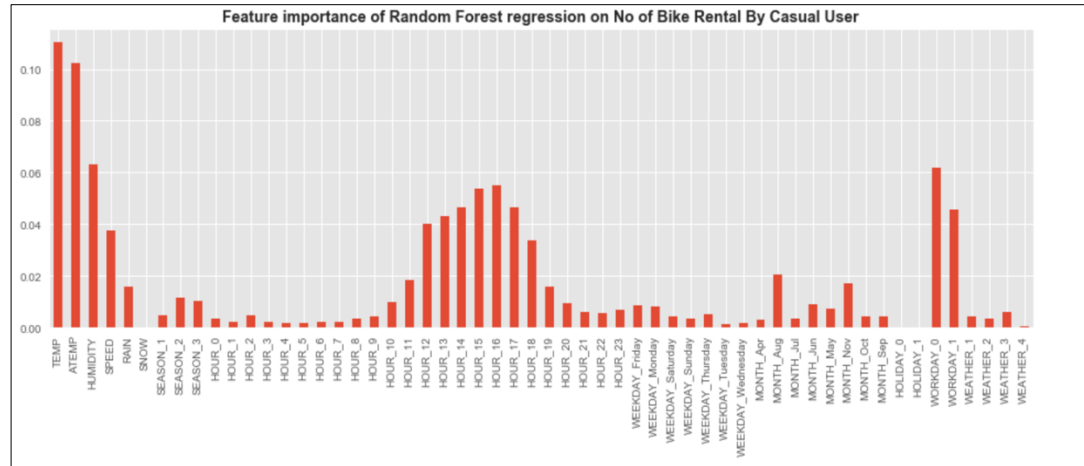
- ❑ Model explains 85% of data's variance but does not generalize well to the test set
- ❑ Temperature and Humidity are most important features during training

Linear Regression - Casual User



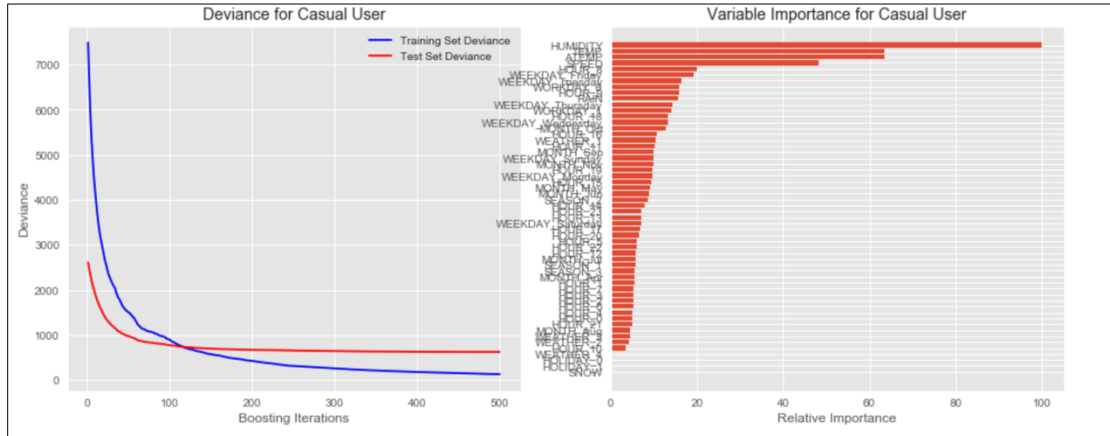
- ❑ Model explains 58% of the variance in the data - a moderate value for fit
- ❑ Training data has a MSE = 897.87 and test data has a MSE = 830.70
- ❑ Linear regression is not a good choice for this data set

Random Tree Regression - Casual User



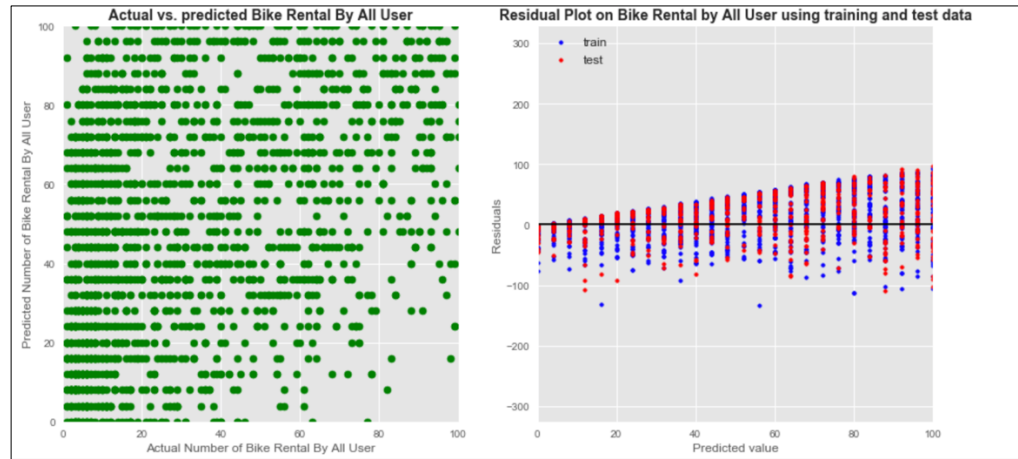
- ❑ Model explains 76% of data's variance but does not generalize well to the test set
- ❑ Temperature and Humidity are most important features during training

Gradient Boosting Regression - Casual User



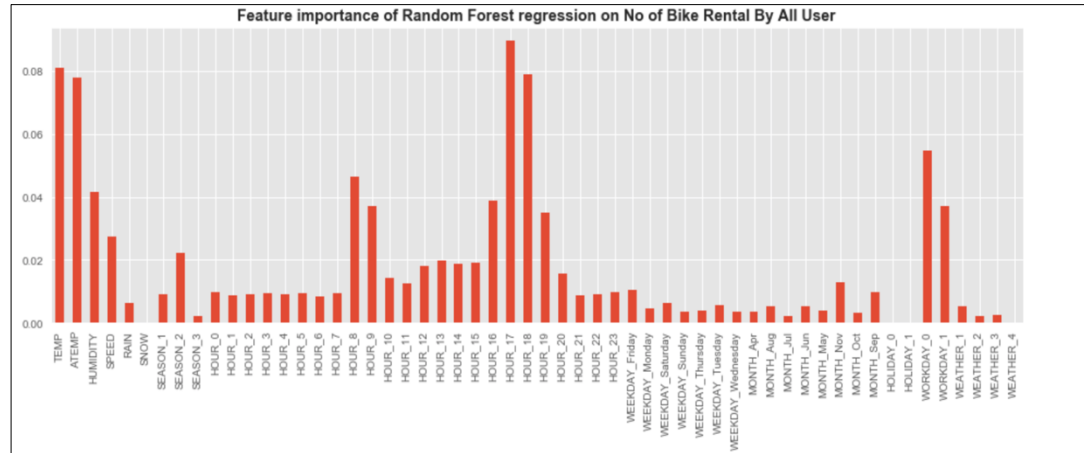
- ❑ Model explains 84% of data's variance but does not generalize well to the test set
- ❑ Temperature and Humidity are most important features during training

Linear Regression - All User



- ❑ Model explains 63% of the variance in the data - a moderate value for fit
- ❑ Training data has a MSE = 5756.49 and test data has a MSE = 5714.19
- ❑ Linear regression is not a good choice for this data set

Random Tree Regression - All User



- ❑ Model explains 81% of data's variance but does not generalize well to the test set
- ❑ Temperature and Humidity are most important features during training

Deviance for all User

Training Set Deviance (blue line) and Test Set Deviance (red line) are plotted against Boosting Iterations (0 to 500). The Training Set Deviance starts high (around 7500) and decreases rapidly, stabilizing around 1000 after 100 iterations. The Test Set Deviance starts lower (around 2500) and decreases more slowly, stabilizing around 700 after 100 iterations.

Variable Importance for all User

Relative Importance of features is shown. The most important feature is HUMIDITY, followed by WIND_SPEED, WIND_DIRECTION, and WIND_DIRECTION_16.

-
- An abstract graphic featuring overlapping triangles in various shades of green and yellow. A thin grey line runs diagonally across the composition. The text 'nce but', 'st set', and 'st' is visible on the left side.

Recommendations and future work

► For Registered User:

- Conduct a study to determine the root cause if there is any relationship between a starting station and ending station and then predict the demand for such combination

► For Casual User:

- Investigate if the number of casual user can be converted in to registered user so that the demand can be prediction with more accuracy .

Future opportunities to continue this work include further training of the regression model on similar city with more accurate dataset to determine if the demand for bike rental holds good.