

# **Bike Rental Sharing Demand Prediction with Machine Learning**

by

Jaanvi Gour  
(2139439)

Under the guidance of  
Dr. Jiju Gillariose



A Project report submitted in partial fulfillment of the requirements for the award of the degree of Master of Science (Data Analytics) of CHRIST (Deemed to be University)

May – 2023

## **CERTIFICATE**

*This is to certify that the report titled **Bike Rental Sharing Demand Prediction** is a bona fide record of work done by **Jaanvi Gour (2139439)** of CHRIST (Deemed to be University), Bangalore, in partial fulfillment of the requirements of VI Trimester MSc (Data Analytics) during the academic year 2022-23.*

**Head of the Department**

**Project Guide**

**Valued-by**

1.

Name : Jaanvi Gour

Register Number : 2139493

2.

Date of Exam : 13-05-2023

### **ACKNOWLEDGEMENT**

I would also like to express my gratitude to our Vice-Chancellor, Dr. Fr. Paul Achandy, Pro-Vice Chancellor Dr. Fr. Joseph CC, our Head of Department (HOD), Dr. Deepthi Das, Coordinator Dr. Sivakumar R and the faculty for providing all the required facilities to accomplish the project.

I would like to convey our gratitude to our project guide, Dr. Jiju Gillariose, for giving us a constant source of inspiration and help in preparing the project, personally correcting our work and providing encouragement throughout the project.

Finally, I would like to thank my parents and friends, without them this assignment would not have been completed.

## **ABSTRACT**

Bike rental sharing has become an increasingly popular mode of transportation in urban areas, providing a convenient and affordable means of transportation for short to medium distance commutes. However, the success of bike rental sharing services largely depends on the accurate prediction of demand, which can be affected by a range of factors, including weather, time of day, day of the week, and seasonal trends.

In this college project, we propose a machine learning-based approach to predict bike rental sharing demand, leveraging historical rental data and a range of external factors that may influence demand. We aim to develop a model that can accurately predict demand for each hour of the day, allowing bike rental sharing providers to optimize their operations and always ensure a sufficient supply of bikes.

Our approach involves the use of a regression model that utilizes a range of input features, including time-related features, weather data, and holiday and event data, among others. We also experimented with different machine learning algorithms, including random forest and gradient boosting, to determine the optimal approach for demand prediction.

To evaluate the performance of our proposed model, we use a range of evaluation metrics, including mean squared error (MSE), root mean squared error (RMSE), and mean absolute error (MAE), among others. We also conduct a comparative analysis of our approach with other commonly used methods for demand prediction, including time-series forecasting techniques and traditional regression models.

The results of our experiments demonstrate that our machine learning-based approach outperforms other methods for demand prediction, achieving high levels of accuracy in predicting hourly demand. Our findings highlight the potential of machine learning algorithms for predicting bike rental sharing demand, and the importance of utilizing external factors to improve the accuracy of demand predictions.

Overall, our project provides valuable insights into the use of machine learning for predicting bike rental sharing demand, and its potential to improve the efficiency and effectiveness of bike rental sharing services in urban area.

## **TABLE OF CONTENTS**

ACKNOWLEDGEMENT.....	(i)
ABSTRACT.....	(ii)
LIST OF TABLES .....	(vi)
LIST OF FIGURES .....	(viii)
<b>1 INTRODUCTION .....</b>	<b>1</b>
1.1.OVERVIEW OF BIKE RENTAL SHARING DEMAND PREDICTION.....	1
1.2.PROBLEM DESCRIPTION.....	2
1.3. EXISTING SYSTEM.....	3
1.4. DRAWBACKS OF THE CURRENT SYSTEM.....	4
1.5. OBJECTIVES.....	5
1.6. SCOPE OF THE PROJECT.....	5
1.7. PURPOSE OF THE PROJECT.....	6
<b>2 LITERATURE REVIEW.....</b>	<b>8</b>
2.1. BIKE RENTING PLATFORMS.....	8
2.1.1. OVERVIEW.....	8
<b>3 DATA OVERVIEW .....</b>	<b>10</b>
<b>4 EXPLORATORY DATA ANALYSIS .....</b>	<b>11</b>
4.1 WEATHER.....	11
4.2 SEASON.....	11
4.3 WORKING DAY.....	12
4.4 HOLIDAY .....	13

4.5	TEMPERATURE .....	13
4.6	HOUR .....	14
4.7	MONTH.....	17
<b>5</b>	<b>CORRELATION ANALYSIS .....</b>	<b>18</b>
5.1	HEATMAP .....	18
5.2	REGRESSION PLOTS.....	19
<b>6</b>	<b>DATA CLEANING .....</b>	<b>20</b>
6.1	MISSING DATA FIELDS .....	20
6.2	HANDLING OUTLIERS .....	20
6.2.1	WEATHER = 'HEAVY SNOW/RAIN' OUTLIER.....	20
6.2.2	ZSCORE OUTLIERS.....	20
<b>7</b>	<b>FEATURE ENGINEERING .....</b>	<b>21</b>
<b>8</b>	<b>MODELLING.....</b>	<b>23</b>
8.1	MODELLING OVERVIEW.....	23
8.2	TRAIN/TEST SPLIT .....	24
8.2.1.	TRAINING SET.....	24
8.2.2.	TEST SET.....	24
8.3	REGRESSION ALGORITHMS SUMMARY .....	25
8.4	STACKING MODEL DETAILS .....	26
8.5	EVALUATION METRIC – RMSLE .....	29
8.6	LINEAR REGRESSION .....	30
8.7	RIDGE REGRESSION.....	31
8.8	LASSO REGRESSION .....	33

8.9	RANDOM FOREST (RF1) .....	37
8.10	RANDOM FOREST (RF2) .....	41
8.11	RANDOM FOREST (RF3) .....	44
8.12	GRADIENT BOOST (GB1) .....	47
8.13	GRADIENT BOOST (GB2) .....	49
8.14	ADABOOST .....	51
8.15	STACKING VIA LINEAR REGRESSION.....	53
8.16	STACKING VIA RANDOM FOREST .....	55
8.17	STACKING VIA GRADIENT BOOST .....	56
<b>9</b>	<b>SUMMARY &amp; CONCLUSIONS.....</b>	<b>58</b>
9.1	RMSLE AND TRAIN+TEST TIME.....	58
9.2	SUMMARY .....	60
9.2.1	DATA EXPLORATION CONCLUSIONS .....	61
9.2.2	MODELING CONCLUSIONS .....	62
9.2.3	LIMITATIONS AND SCOPE FOR MODEL IMPROVEMENTS ....	62
	REFERNCES .....	64

### **LIST OF TABLES**

Table No	Table Name	Page No.
8.3.1	List of models tried	26
8.5.1	RMSLE Metric	29
8.6.1	RMSLE Summary for Linear Regression	30
8.7.1	RMSLE Summary for Ridge Regression	32
8.8.1	Lasso Regression hyperparameter	34
8.8.2	RMSLE Summary for Lasso Regression	34
8.9.1	Hyperparameter tuning	39
8.9.2	RMSLE Summary for Random Forest Regression (RF1)	39
8.10.1	optimal hyperparameters for the RF2 model	42
8.10.2	RMSLE Summary for Random Forest Regression (RF2)	43
8.11.1	optimal hyperparameters for RF3 model	45
8.11.2	RMSLE Summary for Random Forest Regression (RF3)	45
8.12.1	optimal hyperparameters for GB1 model	47
8.12.2	RMSLE Summary for Gradient Boost (GB1)	48
8.13.1	optimal hyperparameters for GB2 model	49
8.13.2	RMSLE Summary for Gradient Boost (GB2)	50
8.14.1	optimal hyperparameters obtained via GridSearchCV	51
8.14.2	RMSLE Summary for Adaboost (AB)	51
8.15.1	RMSLE Summary for Stacking via Linear Regression	53
8.16.1	optimal hyperparameters for this stacked Random Forest model	55
8.16.2	RMSLE Summary for Stacking via Random Forest	55
8.17.1	Optimal hyperparameters for this stacked Gradient Boost model.	56
8.17.2	RMSLE Summary for Stacking via Gradient Boost	57



## LIST OF FIGURES

Figure No	Figure Name	Page No.
4.1	Average bike rental across Weather conditions	11
4.2	Average bike rental across different Seasons	12
4.3	Average bike rentals on Working and Non-Working days	12
4.4	Average bike rentals on Holidays and Non-Holidays	13
4.5	Average bike rentals across different Temperature	14
4.6	Hourly average bike rentals for Working day or non- Working day	15
4.7	Hourly average bike rentals for Casual and Registered Users	15
4.8	Hourly average bike rentals for different days of the week	16
4.9	Average Monthly bike rental count	17
5.1	Heatmap of the correlation between all the numerical features	18
5.2	Regression Plots of 'Count' vs. a) Temperature, b) Humidity and c) Windspeed	19
6.1	Missing value graph	20
8.1	Overview of the Modelling procedure	23
8.2.1	Splitting the provided data into training and testing set.	24
8.4.1	Step1 to Step5 in the stacking algorithm	27
8.4.2	Prediction from individual base models used as features to build new model	28
8.4.3	Modelling via stacking – Overall Summary	28
8.6.1	Linear Regression: Actual vs. Predicted Bike rental count for between 2012-8-15 and 2012-8-19	31
8.7.1	Ridge Regression: Actual vs. Predicted Bike rental count for between 2012-8-15 and 2012-8-19	33

## **1.INTRODUCTION**

### **1.1.OVERVIEW OF BIKE RENTAL SHARING DEMAND PREDICTION**

The bike rental industry has seen significant growth in recent years, driven by an increased focus on sustainability, health and fitness, and urban transportation challenges. Bike rental services typically provide customers with access to a fleet of bicycles for a set period of time, ranging from a few hours to several days.

There are several types of bike rental services available, including traditional bike rental shops, docked bike sharing systems, and dockless bike sharing systems. Traditional bike rental shops typically require customers to visit a physical location to rent a bike, while docked bike sharing systems offer bikes at designated stations where users can pick up and drop off the bikes. Dockless bike sharing systems do not have designated stations, and instead allow users to pick up and drop off bikes anywhere within a specified service area.

The bike rental industry has seen significant growth in urban areas, where cycling can be an efficient and affordable mode of transportation. In addition, bike rental services have become popular for tourists, who use them to explore cities and attractions. Bike rental services have also become popular for health and fitness purposes, as many people use them for exercise and outdoor activities.

One of the challenges facing the bike rental industry is the need for infrastructure, such as bike lanes and designated parking areas, to support safe and efficient bike use. In addition, some bike rental services have faced issues with vandalism and theft, as well as challenges related to maintenance and repair of the bikes.

Overall, the bike rental industry is expected to continue growing as cities focus on sustainability and alternative modes of transportation, and as more

people seek out healthy and active lifestyles.

## **1.2.PROBLEM DESCRIPTION**

The problem with bike rental demand prediction using machine learning is to accurately predict the number of bikes that will be rented at a given time based on historical and real-time data. Accurate demand prediction is crucial for bike rental companies as it helps them optimize their fleet management, staffing, and resources, and ensure that they can meet customer demand.

The challenge with demand prediction in the bike rental industry is that there are many factors that can affect bike rental demand, including weather, seasonality, time of day, day of the week, holidays, and special events. In addition, demand can vary based on location and customer demographics.

Traditional methods of demand prediction, such as regression analysis, can be limited in their ability to capture the complexity and variability of demand factors in the bike rental industry. Machine learning algorithms offer the potential to capture more nuanced relationships between demand factors and bike rental demand, and to incorporate real-time data to improve accuracy.

However, there are also challenges associated with using machine learning for demand prediction, including the need for large and diverse data sets, the potential for overfitting, and the complexity of some algorithms. In addition, machine learning models must be regularly updated and refined to ensure accuracy over time.

Overall, the problem of bike rental demand prediction using machine learning is an important one for the bike rental industry, as accurate demand prediction can help companies improve efficiency and profitability while also providing better service to customers.

## **1.3.EXISTING SYSTEM**

The existing system of bike rental demand prediction typically relies on traditional statistical methods such as time series analysis and regression analysis. These methods use historical data on bike rental demand to identify patterns and trends, and to make predictions about future demand based on those patterns.

In addition, some bike rental companies may use basic forecasting models such as moving averages and exponential smoothing to predict demand. These models are based on historical demand data and assume that future demand will be similar to past demand, with some adjustment for seasonality and trend.

However, these traditional methods have some limitations in capturing the complexity of bike rental demand. They may not be able to account for factors such as weather, events, and customer demographics that can have a significant impact on demand. Furthermore, they may not be able to incorporate real-time data and may require significant manual effort to update and refine over time.

To address these limitations, some bike rental companies are beginning to explore the use of machine learning algorithms for demand prediction. Machine learning algorithms can learn from historical and real-time data to identify patterns and relationships between demand factors and bike rental demand and can incorporate a wider range of factors than traditional methods.

For example, some bike rental companies are using supervised learning algorithms such as linear regression, decision trees, and neural networks to predict demand based on weather data, day of the week, time of day, and other factors. Other companies are using unsupervised learning algorithms such as clustering and association rule mining to identify patterns and relationships in the data that may not be immediately apparent.

Overall, while traditional methods of bike rental demand prediction are still widely used, the growing availability of machine learning algorithms and tools is creating new opportunities to improve the accuracy and efficiency of

demand prediction in the bike rental industry.

#### **1.4. DRAWBACKS OF THE CURRENT SYSTEM**

The existing system of bike rental demand prediction using traditional statistical methods has several drawbacks that limit its accuracy and effectiveness. Some of these drawbacks include:

- i. Limited ability to capture complex relationships: Traditional statistical methods such as regression analysis rely on assumptions about the relationships between variables, which may not capture the complex interactions between demand factors in the bike rental industry.
- ii. Difficulty in incorporating real-time data: Traditional methods often rely on historical data to make predictions, which can be limiting when it comes to incorporating real-time data such as weather conditions, events, and social media trends that can affect bike rental demand.
- iii. Lack of scalability: Traditional methods may not be easily scalable as the volume of data increases, making it difficult to process and analyze large data sets.
- iv. Lack of flexibility: Traditional methods may not be easily adaptable to changing demand patterns or new variables that may impact bike rental demand.
- v. Potential for bias: Traditional methods may introduce bias if the historical data used to train the model reflects biases or changes in demand patterns that are no longer relevant.

Overall, these drawbacks limit the ability of traditional methods to accurately predict bike rental demand and highlight the need for more advanced machine learning techniques that can better capture the complex relationships and dynamics of the bike rental industry.

#### **1.5.OBJECTIVES**

The main objective of bike rental demand prediction using machine learning is to develop accurate and reliable models that can predict bike rental demand in real-time. The key objectives of this approach include:

- i. accuracy: Machine learning algorithms have the potential to capture more nuance relationships between demand factors and bike rental demand, and to incorporate real-time data to improve accuracy.
- ii. Increase efficiency: Accurate demand prediction can help bike rental companies optimize their fleet management, staffing, and resources, and ensure that they can meet customer demand, leading to increased efficiency and profitability.
- iii. Adaptability: Machine learning models can be more adaptable to changing demand patterns or new variables that may impact bike rental demand, ensuring that the models remain accurate and relevant over time.
- iv. Scalability: Machine learning models can be easily scalable to process and analyze large volumes of data, making it possible to generate accurate predictions even with large data sets.
- v. Reduction of bias: Machine learning algorithms can help reduce bias by identifying pattern and trends in the data that may not be immediately apparent using traditional statistical methods.

Overall, the objective of bike rental demand prediction with machine learning is to develop more accurate, efficient, and adaptable models that can help bike rental companies optimize their operations, improve customer service, and maximize profitability.

## **1.6. SCOPE OF THE PROJECT**

The scope of bike rental demand prediction with machine learning is broad, and can be applied to various aspects of bike rental operations. Some of the key areas where machine learning can be applied include:

- i. Fleet management: Bike rental companies can use machine learning algorithms to predict demand for different types of bikes, allowing them to optimize their fleet management by allocating resources to areas where demand is highest.

- ii. Pricing strategies: Machine learning models can help bike rental companies develop pricing strategies that reflect changing demand patterns, allowing them to maximize revenue while remaining competitive.
- iii. Staffing: Accurate demand prediction can help bike rental companies adjust staffing levels to meet changing demand, ensuring that there are enough staff members available to provide high-quality customer service.
- iv. Marketing and promotions: Machine learning can be used to identify patterns in customer behavior and preferences, allowing bike rental companies to develop targeted marketing campaigns and promotions that are more likely to be effective.
- v. Expansion planning: Bike rental companies can use machine learning models to predict demand in new locations or markets, allowing them to make informed decisions about expansion and growth opportunities.

Overall, the scope of bike rental demand prediction with machine learning is vast, and has the potential to transform the way that bike rental companies operate by providing them with more accurate and reliable predictions of demand.

### **1.7.PURPOSE OF THE PROJECT**

The purpose of bike rental demand prediction with machine learning is to provide bike rental companies with more accurate and reliable predictions of demand. By using advanced machine learning algorithms, companies can analyze large volumes of data from various sources, including historical rental data, weather forecasts, events calendars, social media trends, and more, to generate accurate predictions of future demand.

The primary purpose of bike rental demand prediction is to help companies optimize their operations and resources to meet customer demand, while also maximizing profitability. Accurate demand predictions can help companies make informed decisions about fleet management, staffing, pricing strategies, marketing and promotions, and expansion planning. By using machine learning algorithms, companies can develop more accurate and reliable demand prediction models that can be adapted to changing conditions and evolving market trends.

Ultimately, the purpose of bike rental demand prediction with machine learning is

---

to help bike rental companies improve customer satisfaction, reduce costs, and increase revenue by providing the right bikes and services in the right locations at the right times.



---

## 2. LITERATURE REVIEW

### 2.1. BIKE RENTING PLATFORMS

#### 2.1.1. OVERVIEW

Bike renting platforms have become increasingly popular in recent years, as more people seek out convenient and eco-friendly modes of transportation. The following is a brief literature review of some key studies and research papers related to bike renting platforms:

"Bike-sharing: History, impacts, models of provision, and future" by Susan Shaheen and Stacey Guzman (2011): This paper provides an overview of the history and evolution of bike-sharing systems around the world and explores the impacts of bike-sharing on mobility, sustainability, and public health. The authors also discuss various models of bike-sharing provision, including public and private partnerships, and identify key challenges and opportunities for future development.

"An analysis of bike-sharing usage and rebalancing: Evidence from Barcelona" by Juan Carlos Martín, Juan Montero, and Beatriz López-García (2017): This study examines the usage patterns and rebalancing strategies of a bike-sharing system in Barcelona, Spain. The authors analyze data on bike trips and usage patterns and explore the effectiveness of different rebalancing strategies for ensuring that bikes are available in high-demand areas.

"The adoption of bike-sharing in European cities: Cross-national comparisons" by Ralph Buehler and John Pucher (2012): This paper provides a comparative analysis of the adoption and implementation of bike-sharing systems in several European cities, including Paris, Barcelona, and Vienna. The authors explore the factors that have contributed to the success of bike-sharing in different cities, including infrastructure, policy support, and public awareness.

"Bike sharing and the economy of cities" by Susan Shaheen, Elliot Martin, and Nelson Chan (2012): This study examines the economic impacts of bike-sharing

---

systems in various cities around the world, including their effects on tourism, job creation, and local businesses. The authors also explore the potential for bike-sharing systems to reduce transportation costs and increase economic opportunities for low-income residents.

"Designing bike-sharing systems for sustainable mobility: Lessons learned from a comparative case study" by Heleen Buldeo Rai, Klaus Niederhuber, and Silvia Scarpetta (2017): This paper presents a comparative case study of bike-sharing systems in four European cities and explores the factors that have contributed to their success or failure. The authors identify key design principles for sustainable and effective bike-sharing systems, including integration with public transport networks, user-friendly technology, and effective pricing strategies.

Overall, these studies highlight the potential benefits of bike renting platforms for improving urban mobility, sustainability, and economic development, and provide insights into the key factors that contribute to the success of bike-sharing systems in different contexts.

an intuitive approach of the business problem without even thinking about the available data.

---

### 3. DATA OVERVIEW

The data set is obtained from Kaggle with the following column labels.

- datetime - hourly date + timestamp
- season - 1 = spring, 2 = summer, 3 = fall, 4 = winter
- holiday - whether the day is considered a holiday
- working day - whether the day is neither a weekend nor holiday
- weather -
  - 1: Clear, Few clouds, partly cloudy, partly cloudy
  - 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
  - 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
  - 4: Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog
- temp - temperature in Celsius
- atemp - "feels like" temperature in Celsius
- humidity - relative humidity
- windspeed - wind speed
- casual - number of non-registered user rentals initiated
- registered - number of registered user rentals initiated
- count - number of total rentals

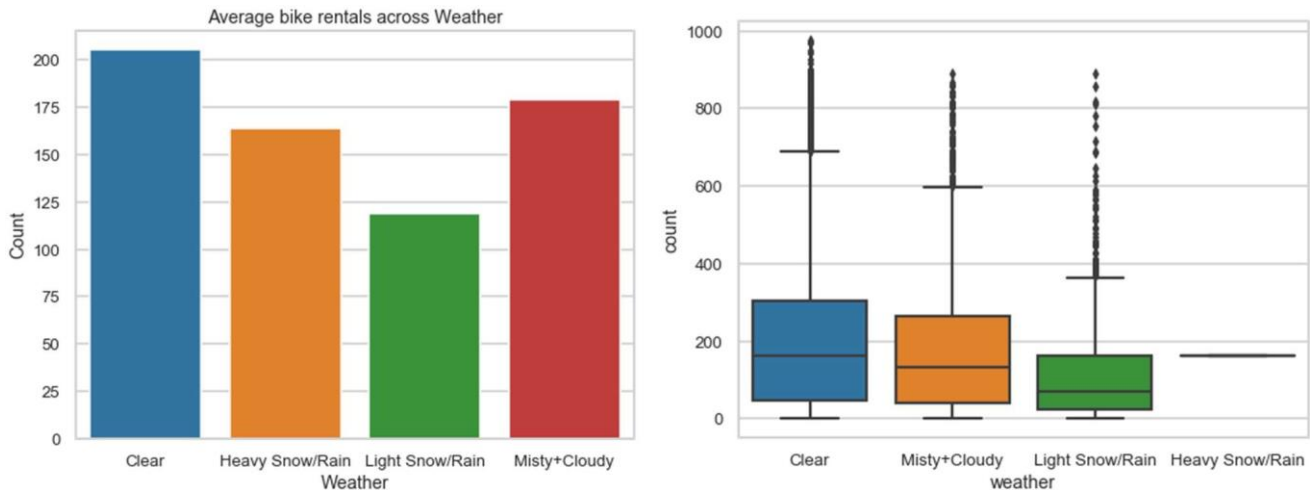
In this project, we use the 8 of the above columns as the feature set ['datetime', 'season', 'holiday', 'working day', 'weather', 'temp', 'atemp', 'humidity', 'windspeed'] to predict the value of 'count'. The other two columns (casual and registered) comprise of the split-up of the target column 'count'. The provided data consists of 10886 observations with 11 column variables (excluding the datetime column - which has been used as an index)

## 4. EXPLORATORY DATA ANALYSIS

Before we start modeling, let us first get an idea on how the number of bike rentals depend on the various features provided to us one by one.

### 4.1. Weather

Below are the plots that show the average bike rental (count) for various weather conditions.



*Figure 4.1: Average bike rental across Weather conditions*

We observe higher bike rentals when the weather is clearer and sunnier. We also notice that there is a single instance where there were rentals under heavy rain/snow condition. Few possibilities arise:

- 4.1.1. Could be an error in logging.
- 4.1.2. Could be an outlier in the data set.
- 4.1.3. Maybe, observations were made at a time when the weather was good. But weather conditions logged sometime later in the same hour when the conditions were heavy rains/snow.

We will investigate this later in the Data Cleaning section (5.2.1) to explore more.

### 4.2. Season

Below are bar plots and box plots of bike rental counts across the 4 seasons.

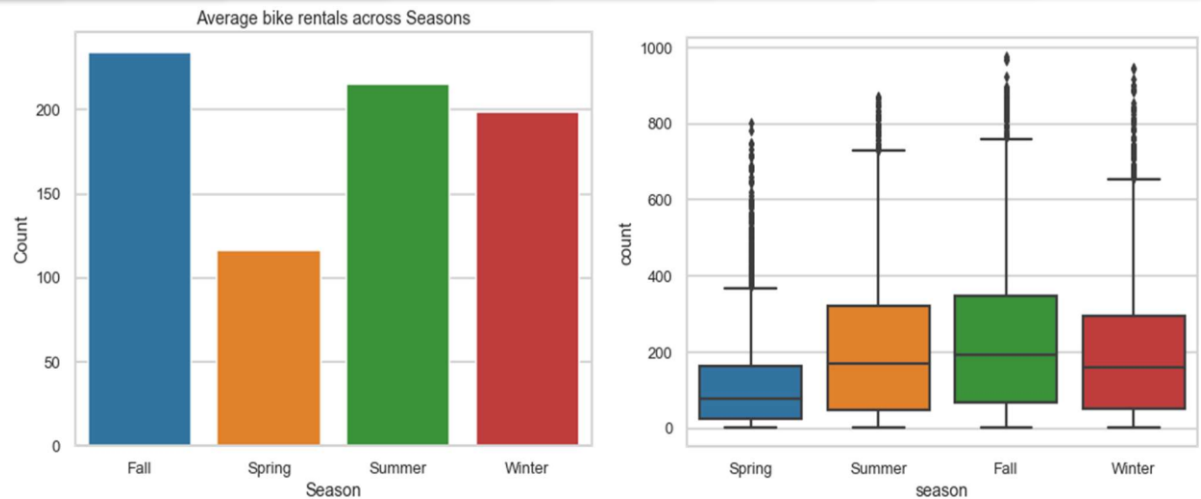


Figure 4.2: Average bike rental across different Seasons

Bike reservations are highest during the Summer (April to June) and Fall (July to September) season and least during the Spring season (January to March)

### 4.3. Working Day

Below are bar plots and box plots of average bike rental counts on working and non-working days.

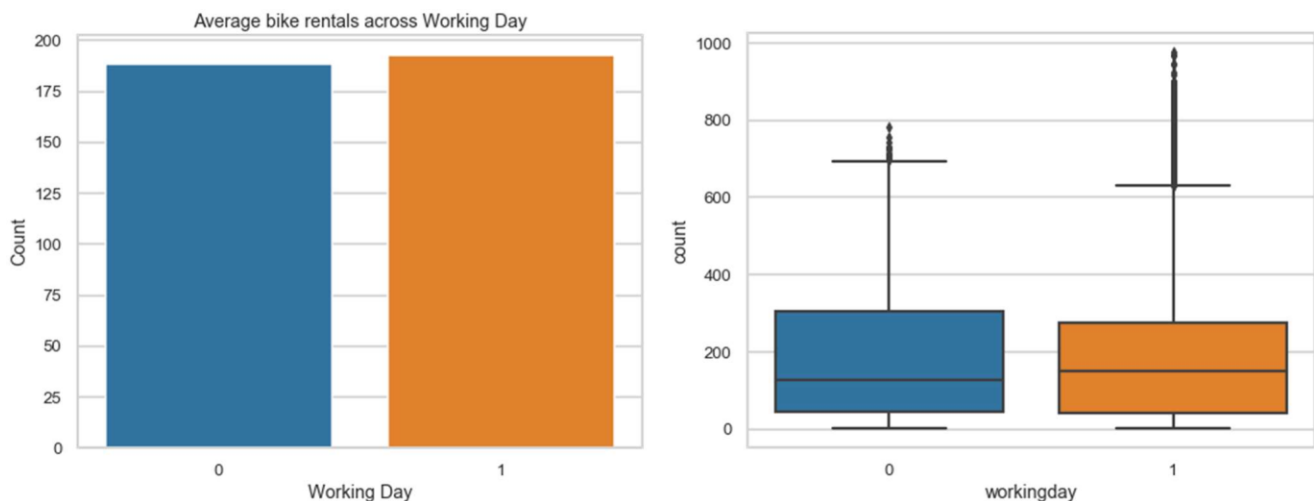


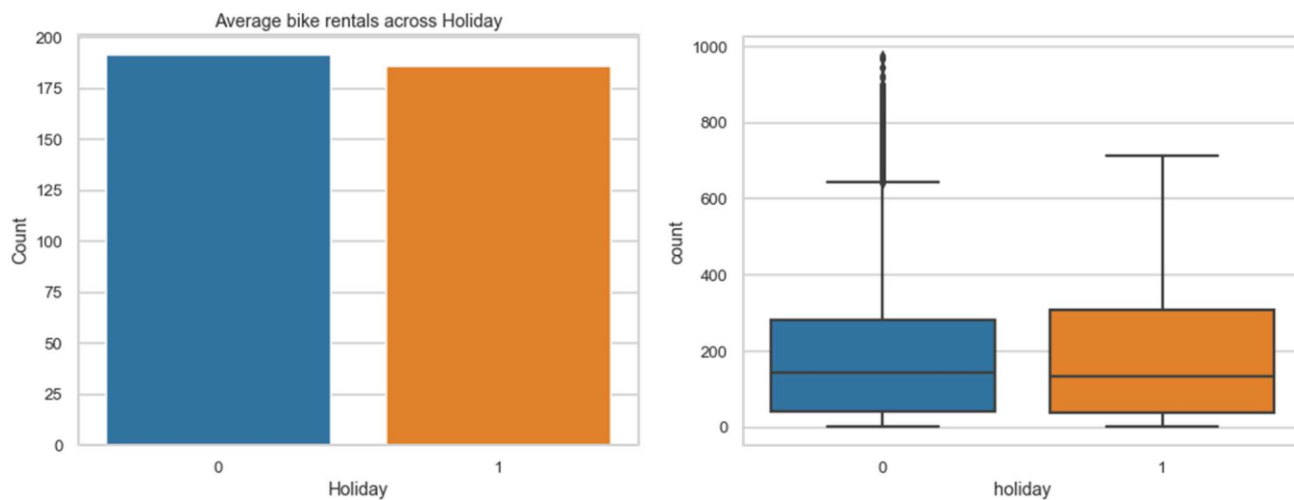
Figure 4.3: Average bike rentals on Working and Non-Working days

We can see that the mean count is similar. However, we see more outliers on working days.

#### 4.4. Holiday

Below are bar plots and box plots of average bike rental counts on holidays and non-holidays.

Holidays correspond to non-working days excluding weekends.

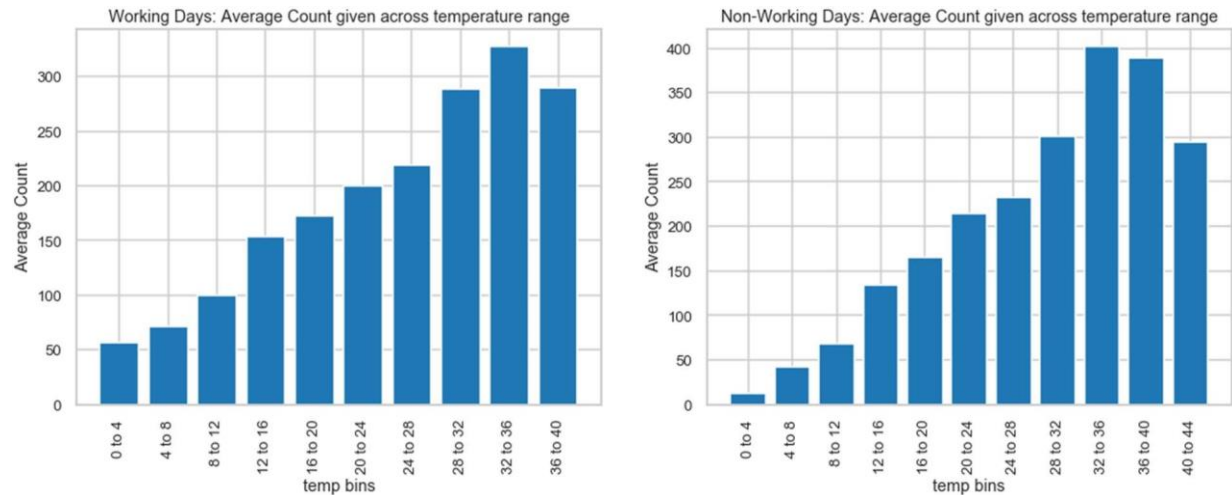


*Figure 4.4: Average bike rentals on Holidays and Non-Holidays*

Mean count seem to be similar with more outlier data points for non-holidays.

#### 4.5. Temperature

Below is a histogram plot for average bike rental count given a certain range of temperature for working and non-working days.



*Figure 4.5: Average bike rentals across different Temperature*

From the above histogram plot, we can see that there is a steady increase in the average bikes rented with temperature with a small decrease at the highest temperature bin. Temperature between 32 and 36 degrees Celsius seems to be the ideal temperature for biking.

#### 4.6.Hour

Now let us examine how the average bike count varies across the day (vs. hour) for the below three categories

- 4.6.1. Working day and non-Working day in Figure 3-6
- 4.6.2. Casual and Registered Users in Figure 3-7
- 4.6.3. Different days of the week in Figure 3-8

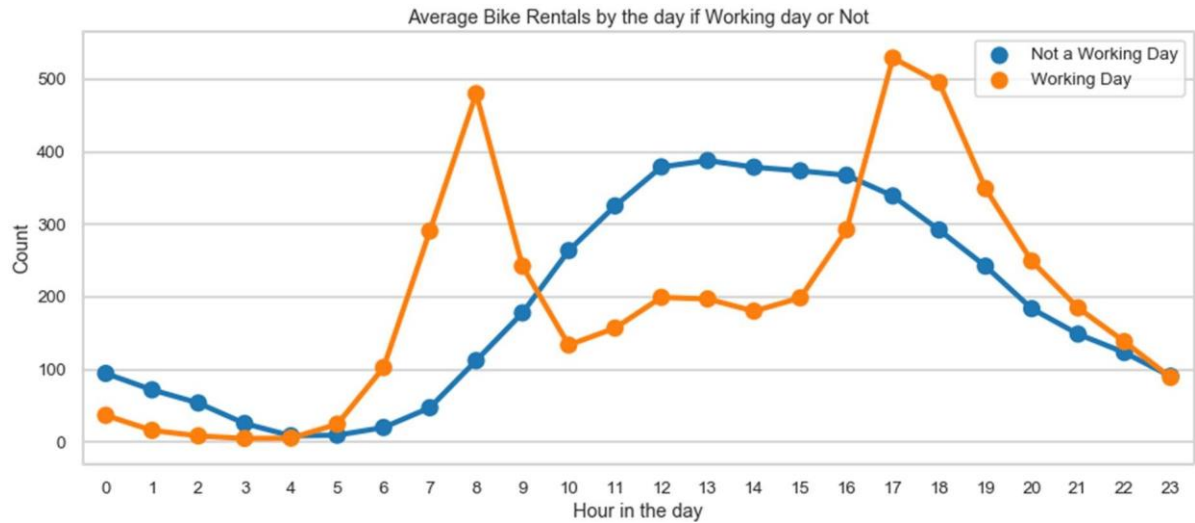


Figure 4.6: Hourly average bike rentals for Working day or non- Working day

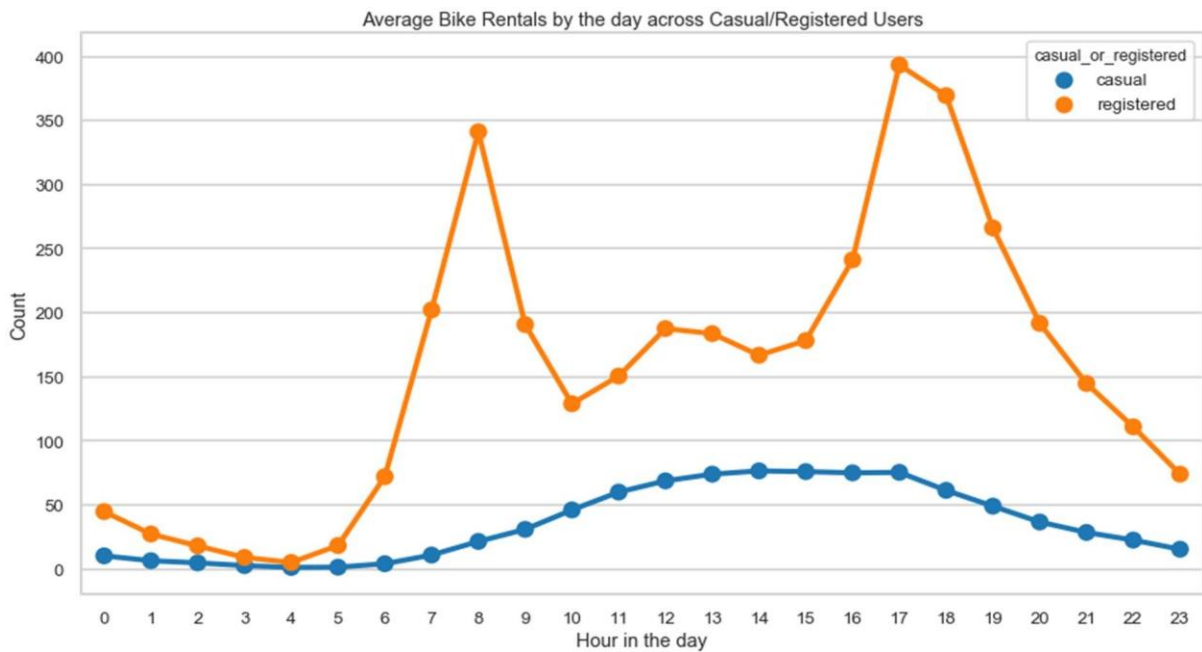


Figure 4.7: Hourly average bike rentals for Casual and Registered Users



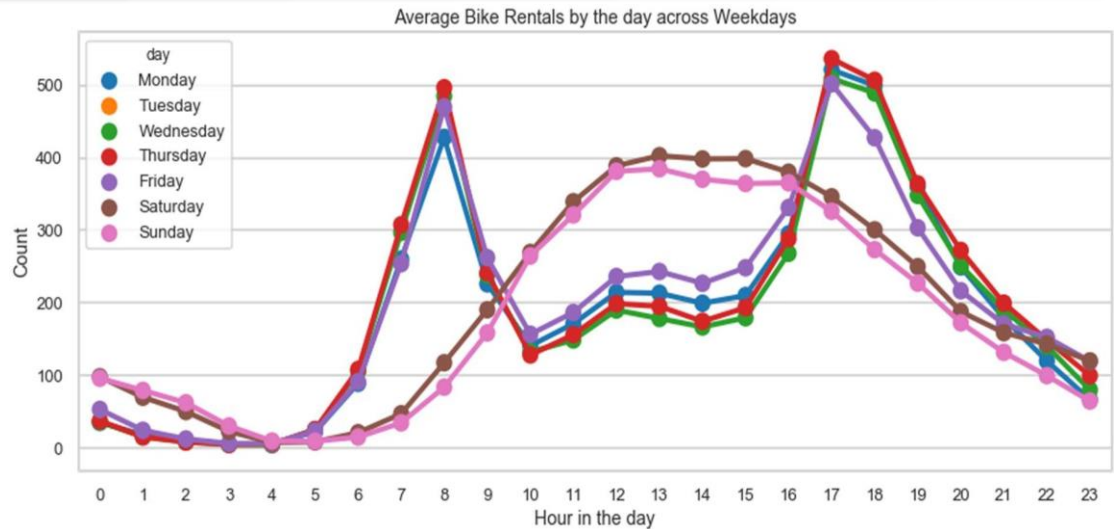


Figure 4.8: Hourly average bike rentals for different days of the week

There are few interesting observations that we can make from the above 3 figures.

4.6.4. Firstly, we see that there are 2 patterns in the bike rentals during the day.

4.6.4.1. Working Day: First pattern where there is a peak in the rentals at around 8am and another at around 5pm. These correspond to working local bikers who typically are *registered* and go to work on working day which are *Monday to Friday*

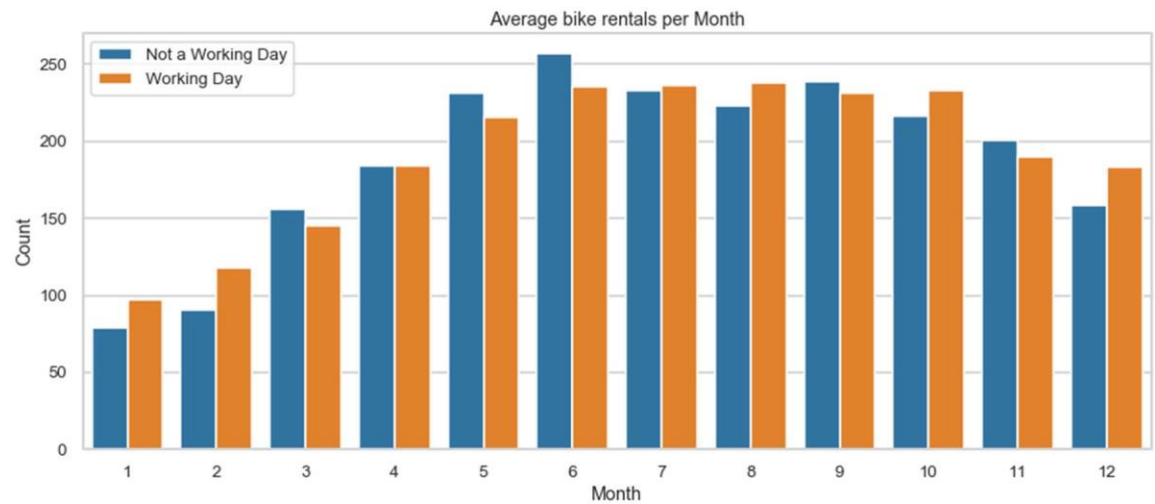
4.6.4.2. Non-Working Day: Second pattern where there is a steady increase in bike rental during the day reaching the peak at around noon and a gradual decrease in the count in the evening after 4pm. These correspond to probably tourists who typically are *casual* users who rent/drop off bikes uniformly during the day and tour the city of Washington on non- working days which typically are *Saturday* and *Sunday*

4.6.5. The average number of the casual users are in general lower than the average number for non- working days and weekends. This seem to indicate that there are several registered users too, who follow the same non-working day rental trend during weekends and non-working day.

---

### 4.7.Month

Below plot contains the average bike count over each month of a calendar year.



*Figure 4.9: Average Monthly bike rental count*

The above figure is highly correlated with the seasons bar plot since seasons plot effectively is the average count for 3 of these months. We can see that the most rentals are in the months of June and May while the lowest are on January and February.

## 5. CORRELATION ANALYSIS

### 5.1.Heatmap

Below is a heatmap plot of the correlation between all the numerical columns.

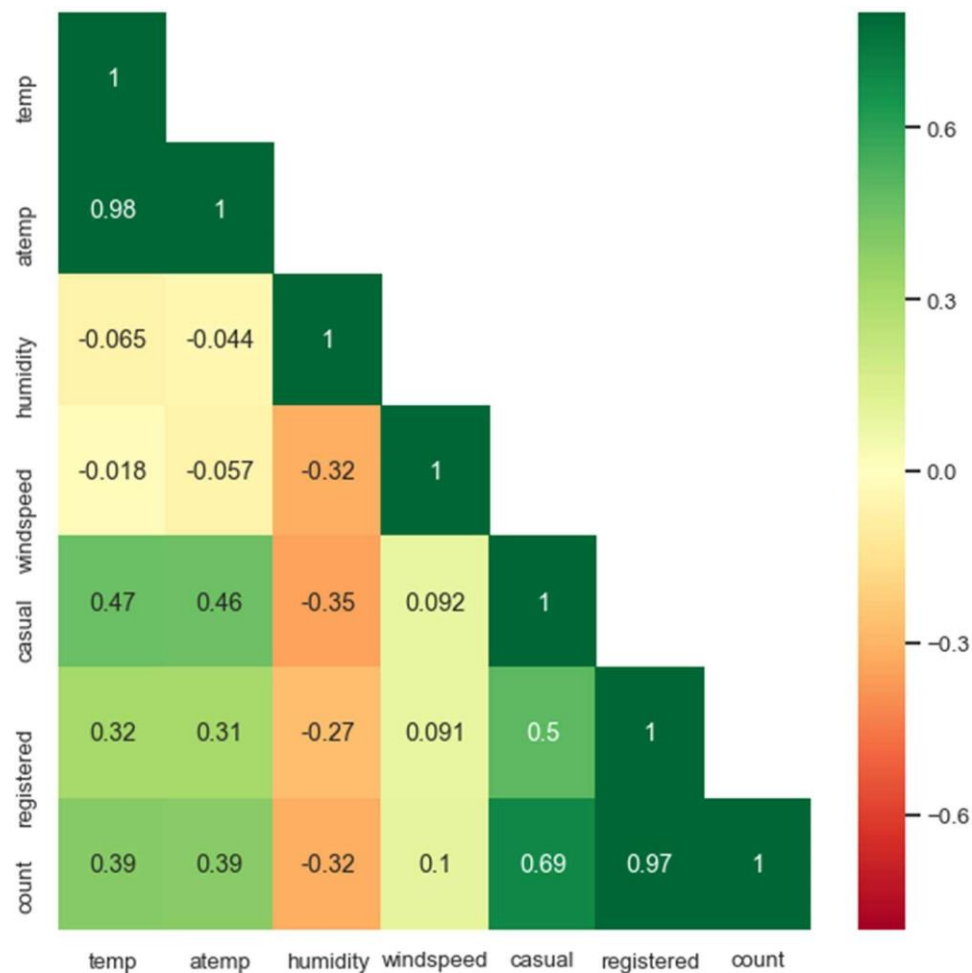


Figure 5.1: Heatmap of the correlation between all the numerical features

We can infer the following from the above heatmap.

- 5.1.1. *temp* (true temperature) and *atemp* (feels like temperature) are highly correlated, as one would expect.
- 5.1.2. *count* is highly correlated with *casual* and *registered* as expected since  $\text{count} = \text{casual} + \text{registered}$

5.1.3. We see a positive correlation between *count* and *temperature*. Note that this is only true for the range of temperatures provided.

5.1.4. We see a negative correlation between *count* and *humidity*. The more humidity, the less people prefer to bike.

5.1.5. *Count* has a weak dependence on *windspeed*.

## 5.2. Regression Plots

Below are regression plots of the bike rental count vs. Temperature, Humidity and Windspeed, respectively.

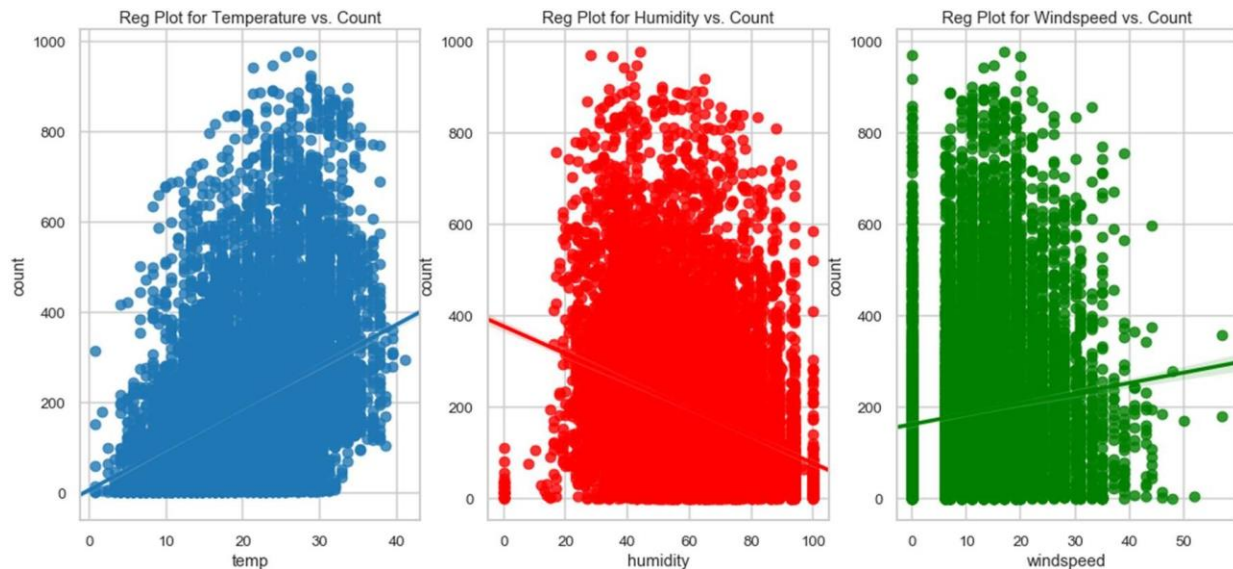


Figure 5.2: Regression Plots of 'Count' vs. a) Temperature, b) Humidity and c) Windspeed

The above regression plots indicate a strong positive correlation of *count* with *temperature*, a weak positive correlation with *windspeed*, and a strong negative correlation with *humidity* (as we saw in the heatmap plots too). We can also see several instances with *windspeed* = 0. These are probably missing values. Since *windspeed* has a very low dependence and has several missing (or erroneous) data points, let us exclude this from our model.

## 6. DATA CLEANING

### 6.1. Missing Data Fields

As seen from the below figure, the provided data set has no missing data fields for any of the features.

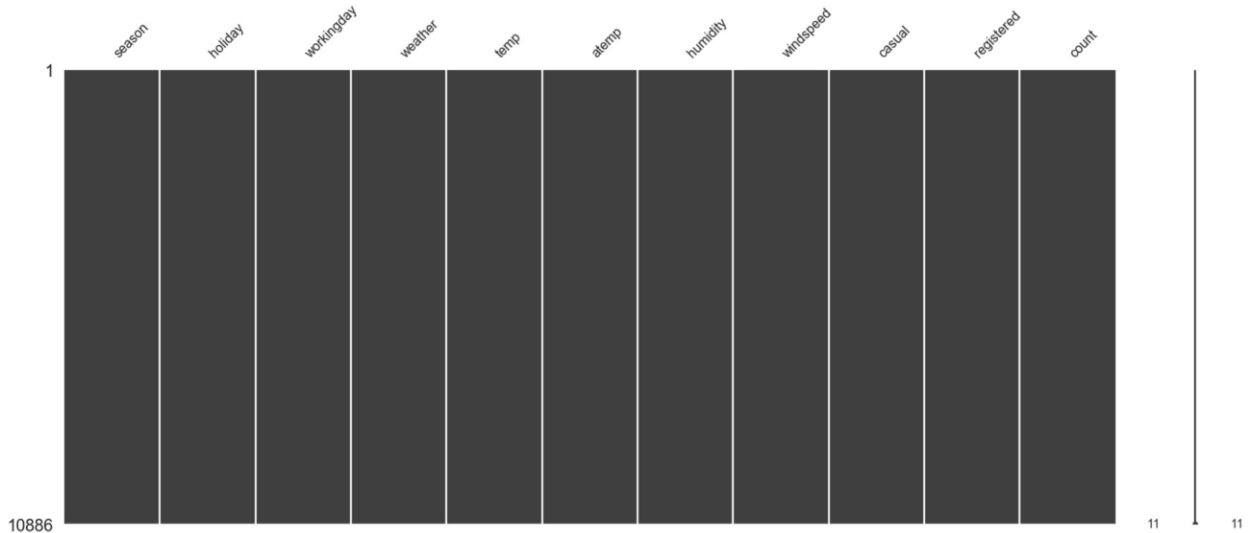


Figure 6.1: Missing value graph

### 6.2. Handling Outliers

#### 6.2.1. Weather = 'Heavy Snow/Rain' outlier

The provided data set has just one instance of 'Heavy Snow/Rain' observation and wouldn't be very helpful in training any of our model. Hence, we replace that instance with 'Light Snow/Rain.'

#### 6.2.2. Z-score outliers

We prune out all observations with z-score  $> 4$ , i.e., data with more than 4 standard deviations away from the mean. These correspond to observations which have probability of  $\sim 6 \times 10^{-5}$ . There were 15 such observations and all these outliers occur mostly early in the morning or late at night. These could be due to some late-night shows or events. We prune these outliers out and use the remaining 10871 observations to train our model.

---

## 7. FEATURE ENGINEERING

The provided data in its raw form wasn't directly used as an input to the model. Several feature engineering was carried out where few features were modified, few were dropped, and few were added. Below is a summary of the feature engineering carried out with the provided data set.

1. The *datetime* column which contained the date-time stamp in 'yyyy-mm-dd hh:mm:ss' format was split into individual [*'month'*, *'date'*, *'day'*, *'hour'*] categorical columns.
2. Drop *season* column: This is because *month* column has a direct mapping with *season* (Winter: January to March, Summer: April to June, Fall: July to September and Spring: October to December). Hence, we retain *month* column due to its higher cardinality.
3. Drop *holiday* and *day* columns: The working day column had information about holiday embedded in it. *working day* = weekday and not a holiday. Since we noticed that there were two kinds of bikerental behaviors - during working days and not a working day, we will retain only the *working day* column and drop *'day'* and *'holiday'* column.
4. Drop *date* column: Intuitively, there should be no dependency on date. Hence drop this column.
5. Drop *casual* and *registered* columns: These are individual components of the 'to be predicted' column (*count*). These are not provided in the Kaggle Test Data too. Hence drop these columns.
6. Drop *windspeed* column: Very poorly correlated with *count* and has several missing/erroneous data. Hence drop this column.
7. Drop *atemp* column: *temp* and *atemp* are very highly correlated and essentially indicate the something. Hence retain only the *temp* column
8. One Hot Encoding of categorical feature set
  - a. *Weather*: Split weather column to *weather\_1*, *weather\_2* and *weather\_3* (recall that we had relabeled all the weather = 4 data

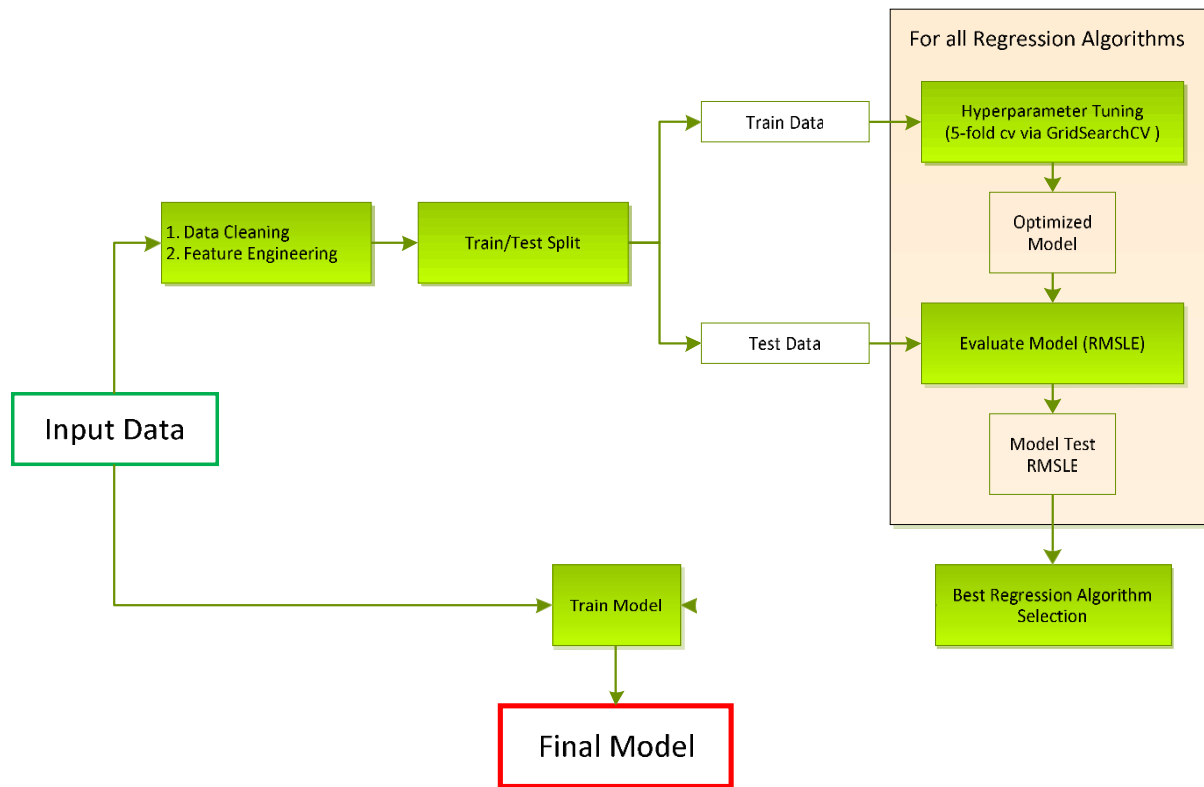
---

points to *weather* = 3 due to its sparseness). Drop *weather\_3* and *weather* columns since it is a function of the rest of the retained *weather* columns.

- b. *Month*: Split *month* column to *month\_1*, *month\_2*, ..., *month\_12*. Drop *month\_12* and *month* columns since they are a function of the rest of the retained *month* columns.
- c. *Hour*: Split *hour* column to *hour\_0*, *hour\_1*, ..., *hour\_23*. Drop *hour\_23* and *hour* columns since they are a function of the rest of the retained *hour* columns.

## 8. MODELLING

### 8.1. Modeling Overview



*Figure 8.1: Overview of the Modelling procedure*

Figure 7-1 depicts the overall procedure followed to obtain the Final Model. The provided data is first cleaned and transformed using Feature Engineering. We then split the data into Train set (for Hyperparameter tuning) and Test set (for Model Evaluation). Using RMSLE as our evaluation metric, we compare various models and select the regression algorithm based on the lowest RMSLE on the Test data. The final model used for submission is then obtained by again training the selected Regression Algorithm on the entire Input Data set.



## 8.2. Train/Test Split

Below figure summarizes the method to split the provided data into training and testing data set. Kaggle has held out data from 20th to the end of the month (for every month) as test set (no labels, i.e., count values have been provided for those data). We follow a similar approach to split the provided labelled data set (consisting of 1<sup>st</sup> to 19<sup>th</sup> of every month) into two sets.

### 8.2.1. Training set

8.2.1.1. This contains data from the 1<sup>st</sup> to 15<sup>th</sup> of every month.

8.2.1.2. This set is used to train various model and obtain the best set of hyperparameters for these models. We use GridSearchCV to tune the hyperparameters using this training set.

### 8.2.2. Test set

8.2.2.1. This contains data from the 16<sup>th</sup> to 19<sup>th</sup> of every month.

8.2.2.2. This is used to evaluate all our models. The model with the best test score is finally chosen for submission.

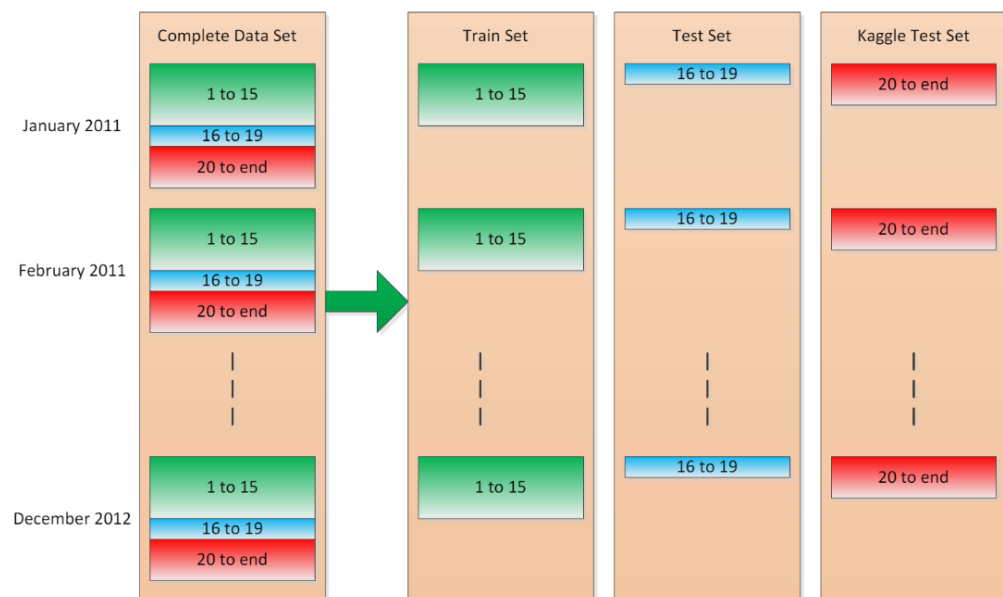


Figure 8.2.1: Splitting the provided data into training and testing set.

### 8.3. Regression Algorithms Summary

Below is a list of models tried out and a brief description of the algorithm.

Category	Modeling Algorithm	Details
<b>Linear</b>	Linear Regression	Two separate models for Working and Non-Working days. One Hot encoded Features
<b>Linear</b>	Ridge Regression	Two separate models for Working and Non-Working days. One Hot encoded Features
<b>Linear</b>	Lasso Regression	Two separate models for Working and Non-Working days. One Hot encoded Features
<b>Ensemble</b>	Random Forest-1	Single Model for Working and Non-Working days. No One Hot Encoding
<b>Ensemble</b>	Random Forest-2	Two separate models for Working and Non-Working days. One Hot encoded Features
<b>Ensemble</b>	Random Forest-3	Two separate models for Working and Non-Working days. No One Hot Encoding
<b>Ensemble</b>	Gradient Boost-1	Two separate models for Working and Non-Working days. One Hot encoded Features
<b>Ensemble</b>	Gradient Boost-2	Two separate models for Working and Non-Working days. No One Hot Encoding
<b>Ensemble</b>		Two separate models for Working and Non-Working days. No One Hot Encoding

<b>Ensemble</b>	Ada Boost	
<b>Stacking</b>	Linear Regression	Use all the Linear+Ensemble model prediction as the Feature set to make final predictions
<b>Stacking</b>	Random Forest	Use all the Linear+Ensemble model prediction as the Feature set to make final predictions
<b>Stacking</b>	Gradient Boost	Use all the Linear+Ensemble model prediction as the Feature set to make final predictions

*Table 8.3.1: List of models tried*

#### 8.4. Stacking Model Details

Stacking is an ensemble learning technique that uses predictions from multiple models (for example Linear Regression, Random Forest, Gradient Boost) to build a new model. This model is used for making predictions on the test set. Below is a stepwise explanation for a simple stacked ensemble:

- 8.4.1. The train set is split into 5 parts – data from a) (1<sup>st</sup> to 3<sup>rd</sup>) of every month, b) (4<sup>th</sup> to 6<sup>th</sup>) of every month, c) (7<sup>th</sup> to 9<sup>th</sup>) of every month, d) (10<sup>th</sup> to 12<sup>th</sup>) of every month and e) (13<sup>th</sup> to 15<sup>th</sup>) of every month.
- 8.4.2. A base model (E.g., Linear Regression) is fitted on 4 parts and predictions are made for the 5<sup>th</sup> part. In Figure 7-3, LR-Model 1 corresponds to Linear Regression Model obtained by training all except the 1<sup>st</sup> part. This is done for each part of the train set (LR-Model 1 to LR-Model 5)
- 8.4.3. The base model (Linear Regression Model) is then fitted on the whole train data set (LR-Model)
- Using this model (LR-Model), predictions are made on the Test Set (containing 16<sup>th</sup> to 19<sup>th</sup> of every month)
- 8.4.4. Steps 2 to 4 are repeated for another base model (say Random

Forest) resulting in another set of predictions for the train set and test set.

Figure 7-3 depicts Steps 1 through 5.

8.4.5. All the predictions from train and test set are compiled into a new table as shown in Figure 7-4. The predictions from the train set are used as features to build a new model.

8.4.6. This model is used to make final predictions on the test prediction set.

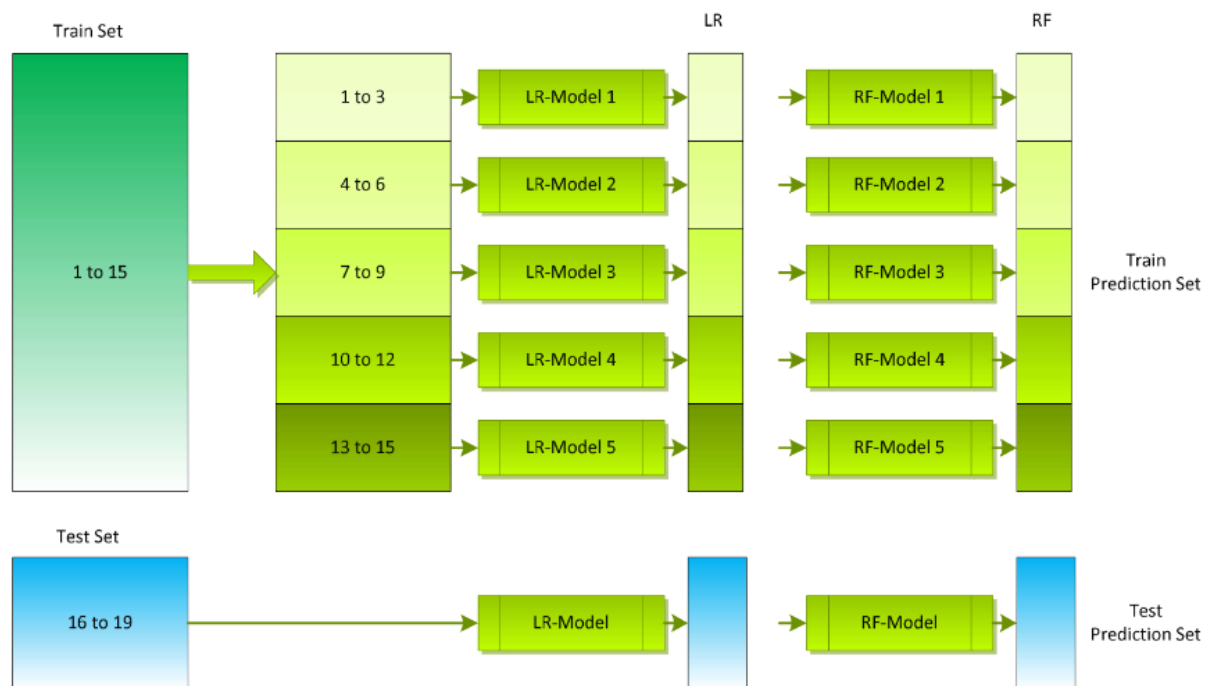


Figure 8.4.1: Step1 to Step5 in the stacking algorithm

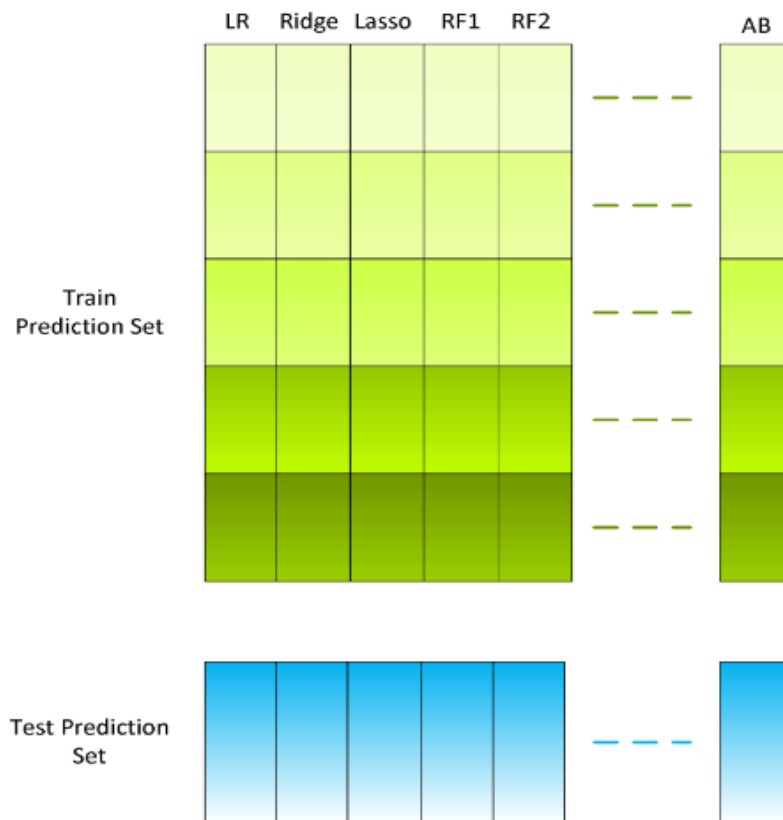


Figure 8.4.2: Prediction from individual base models used as features to build new model

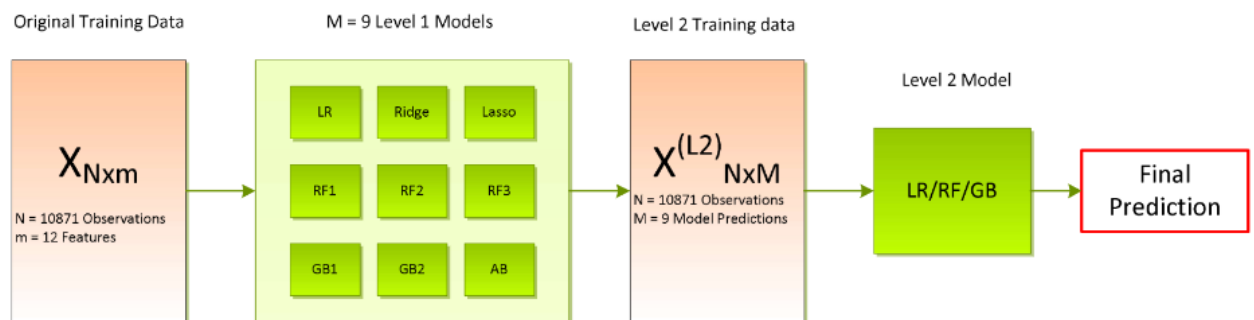


Figure 8.4.3: Modelling via stacking – Overall Summary

The overall summary of the Model that we developed via Stacking is shown in Figure 7-5. Predictions from 9 Individual Models are used as features for the Level 2 Model (for which Linear Regression, Random Forest and Gradient Boost were tried out) to obtain the final prediction.

### 8.5.Evaluation Metric – RMSLE

Root Mean Square Log Error (RMSLE) is chosen as the evaluation metric by Kaggle and is used in our modelling too. The RMSLE is calculated as

$$\sqrt{\frac{1}{n} \sum_i^n (\log(p_i + 1) - \log(a_i + 1))^2}$$

Where

- $n$  is the number of hours in the test set
- $p_i$  is the predicted count
- $a_i$  is the actual count
- $\log(x)$  is the natural logarithm

As a result, each of our models are trained using  $\log(\text{count})$  as the target columns (instead of count). This would result in the predicted values being returned in the log domain. Thus, we further transform the predicted output via  $\exp(y_{\text{pred}}) - 1$ .

For each model, we will evaluate the following metrics for 3 portions of data: {Working Days, Non- Working Days, All combined}

RMSLE Metric	Model Training Data Set	Model Testing Data Set
<b>Train</b>	Train Set	Train Set
<b>CV Test</b>	(4/5) th of the Train Set	The remaining (1/5) th of the Train Set
<b>Test</b>	Train Set	Test Set

*Table 8.5.1: RMSLE Metric*

CV Test RMSLE is obtained as a by-product of the above Stacking Model building. It is essentially the RMSLE computed on the Train prediction and Test prediction set for each of the models.

## 8.6.Linear Regression

For Linear Regression, we obtain two separate models – one for working and the other for non-working days. We also use the entire set of features obtained via OneHotEncoding.

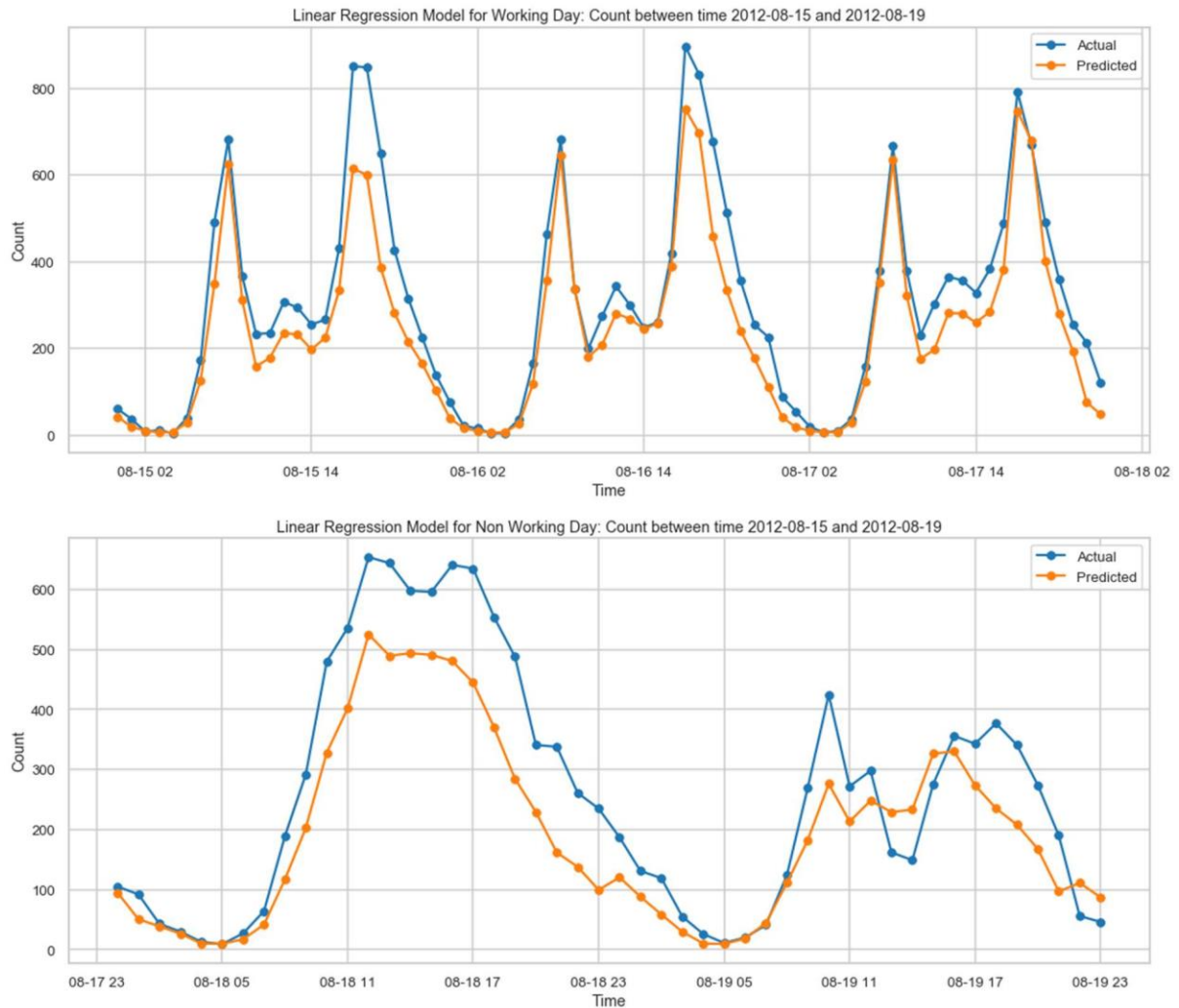
The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.418155	0.432817	<b>0.422797</b>
<b>CV Test</b>	0.430658	0.474808	<b>0.444944</b>
<b>Test</b>	0.390881	0.495181	<b>0.428666</b>
<b>Total Train + Test Time</b>			<b>0.197 sec</b>

*Table 8.6.1: RMSLE Summary for Linear Regression*

RMSLE on the test data = 0.43 and training data = 0.42 are almost the same. Linear Regression model is definitely not an overfit model.

Below plots the true *count* value vs. the model predicted value for a few days in the test data set. We can see that the predicted *count* closely tracks the true value.



*Figure 8.6.1: Linear Regression: Actual vs. Predicted Bike rental count for between 2012-8-15 and 2012-8-19*

### 8.7. Ridge Regression

For Ridge Regression too, we obtain two separate models – one for working and the other for non-working days. We also use the entire set of features obtained via OneHotEncoding. We tune the hyperparameter using the Train Set using GridSearchCV with 5-fold cross validation.



Hyperparameters	Working Day Model	Non-Working Day Model
$\alpha$	10	10

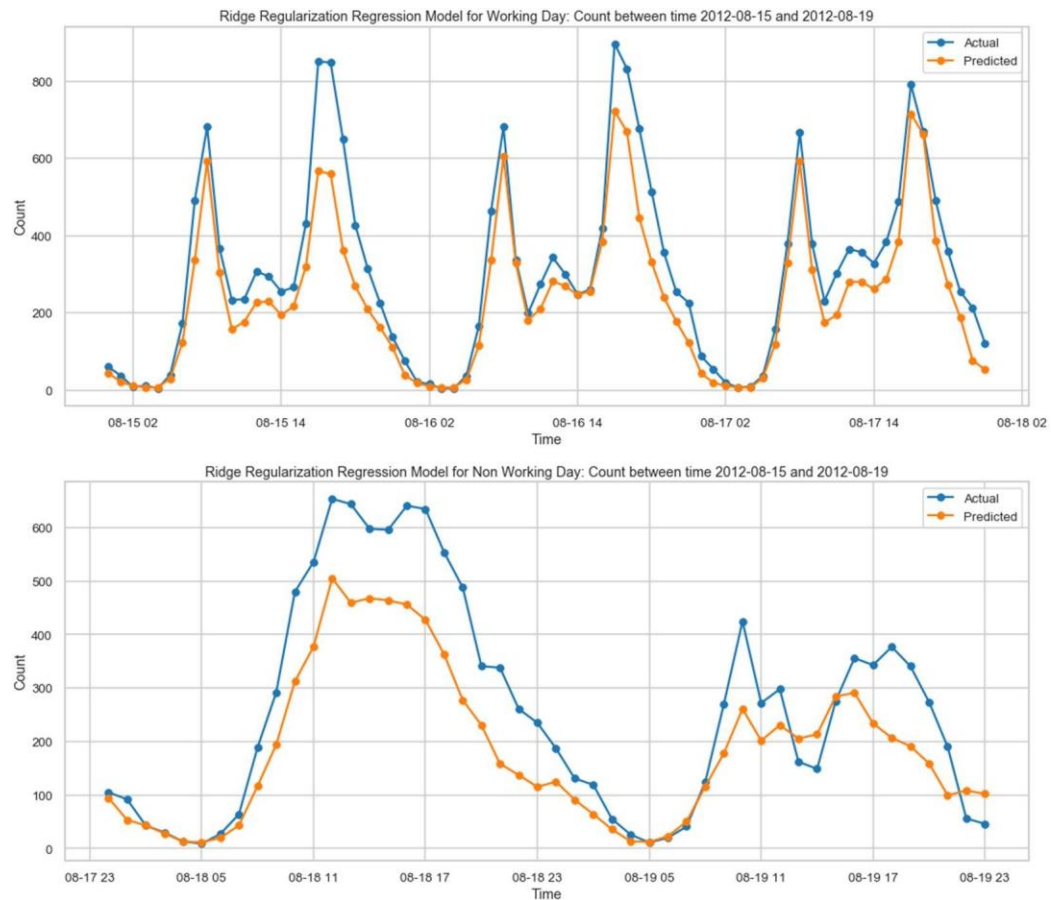
The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.423189	0.448141	<b>0.431152</b>
<b>CV Test</b>	0.438216	0.492361	<b>0.45585</b>
<b>Test</b>	0.39423	0.516847	<b>0.439148</b>
<b>Total Train + Test Time</b>			<b>1.3 sec</b>

*Table 8.7.1: RMSLE Summary for Ridge Regression*

Ridge Regression Model has a higher RMSLE compared to the Linear Regression Model.

Below plot the actual 'count' value vs. the model predicted value for a few days in the test data set.



*Figure 8.7.1: Ridge Regression: Actual vs. Predicted Bike rental count for between 2012-8-15 and 2012-8-19*

## 8.8. Lasso Regression

As with the previous two Linear models, we obtain two separate models – one for working and the other for non-working days and use the entire set of features obtained via OneHotEncoding for Lasso RegressionModel too

We tune the hyperparameter using the Train Set using GridSearchCV with 5-fold cross validation. Below table summarizes the optimal hyperparameters for Lasso Regression

Hyperparameters	Working Day Model	Non-Working Day Model
$\alpha$	0.1	0.5

*Table 8.8.1: Lasso Regression hyperparameter*

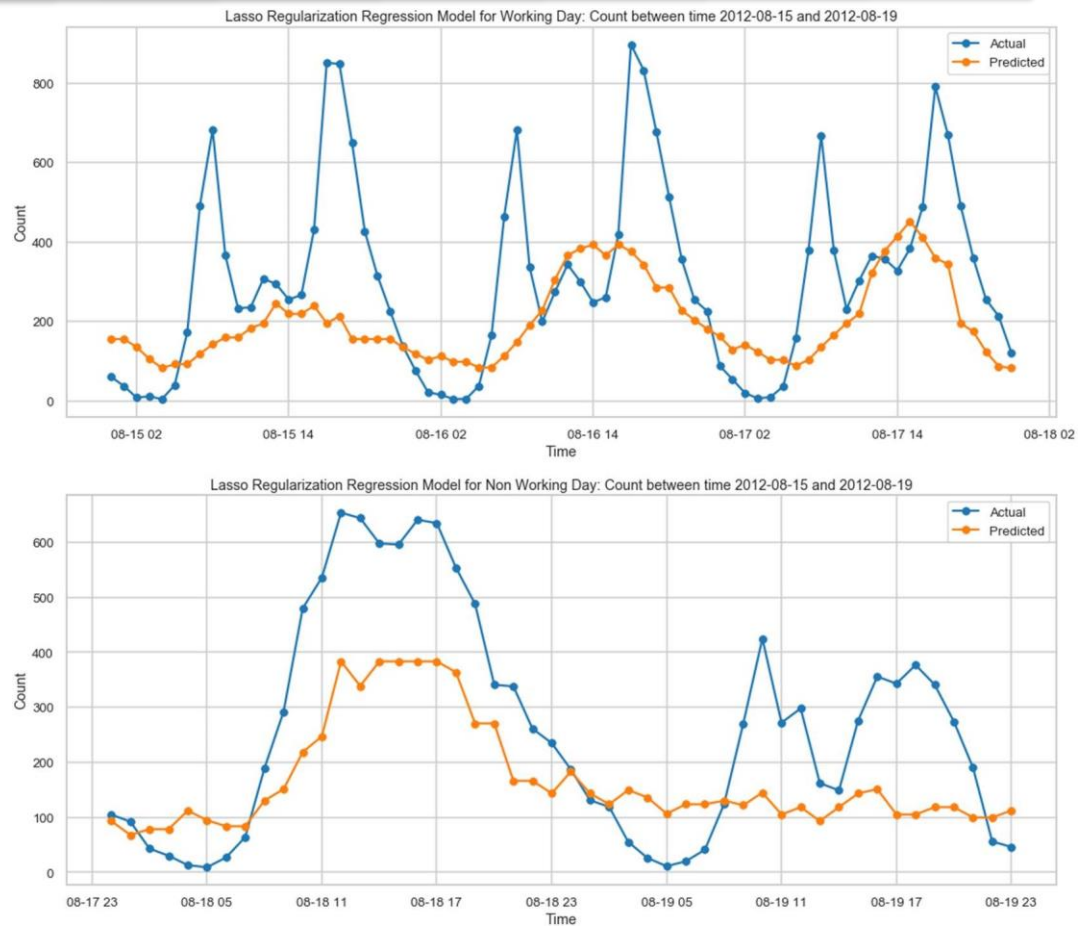
The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	1.313353	1.065252	<b>1.241062</b>
<b>CV Test</b>	1.31459	1.075364	<b>1.244685</b>
<b>Test</b>	1.285346	1.117959	<b>1.231794</b>
<b>Total Train + Test Time</b>			<b>1.3 sec</b>

*Table 8.8.2: RMSLE Summary for Lasso Regression*

Clearly, the Lasso Regression Model performs very bad. Lasso Model has a RMSLE of 1.23 for the test data compared to 0.43 for Linear Regression

Below plot the actual 'count' value vs. the model predicted value for a few days in the test data set. We can see that the predicted value is way off the true value.



*Figure 8.8.1: Lasso Regression: Actual vs. Predicted Bike rental count for between 2012-8-15 and 2012-8-19*

### **Feature Importance – Linear Models**

Feature importance for the Linear Models are obtained from the Coefficients assigned to each of the feature. Higher the magnitude, more is the importance. Figure 7-9 compares the feature coefficient for each of the 3 Linear models.

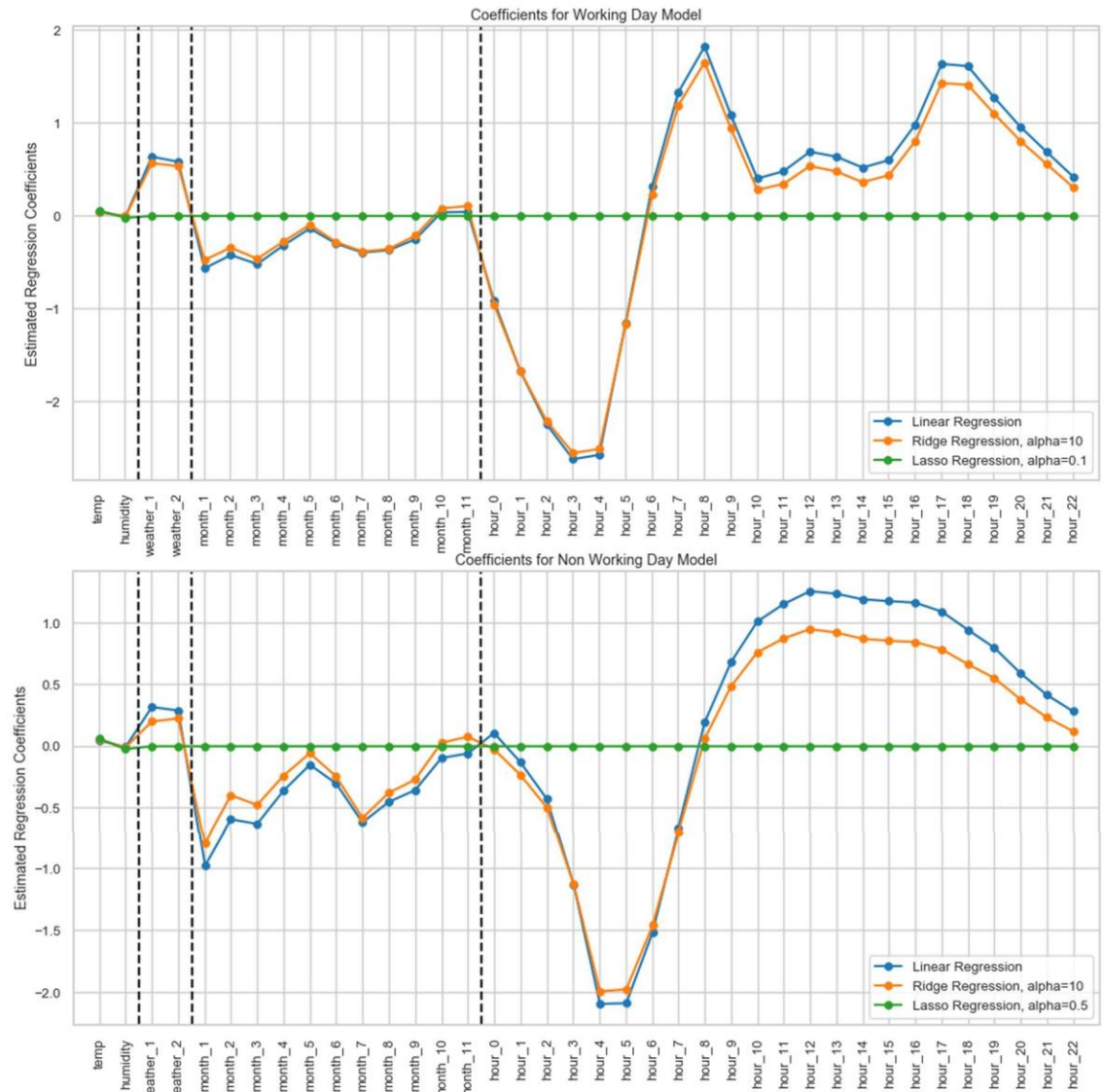


Figure 8.8.2: Feature Coefficient Comparison for the Linear Models (Linear Regression, Ridge and Lasso)

Below are few observations from the

- Linear Regression coefficients are most aggressive (highest magnitude). Ridge Regression coefficients closely follows Linear Regression. And Lasso nullifies most of the coefficients.

- 
- Lasso model provides a non-zero value only for temp and humidity for working day model
  - From the Linear and Ridge coefficient plots, we can see that maximum dependency on the bike rental count are the hour in the day (highest coefficient magnitudes). Negative values of coefficient for early morning hours and positive values with greater magnitude during peak hours of the respective models.
  - Working day coefficient has highest positive coefficient values at 8am and 5pm
  - Non-working day coefficients has a single peak across hours ~ 12 noon (as observed from our earlier plots)
  - Weather\_1 and weather\_2 have positive coefficient value, indicating a negative bias for weather\_3 (light snow/rain). Note that if weather=3, then weather\_1 and weather\_2 = 0. Hence effective weather\_3 coefficient = 0
  - Though the coefficient value for temp and humidity is low compared to the others, the range of these values are higher too.
  - The absolute value of the coefficients is higher for months 5, 6, 10, 11, 12 indicating higher bike rentals during that month.

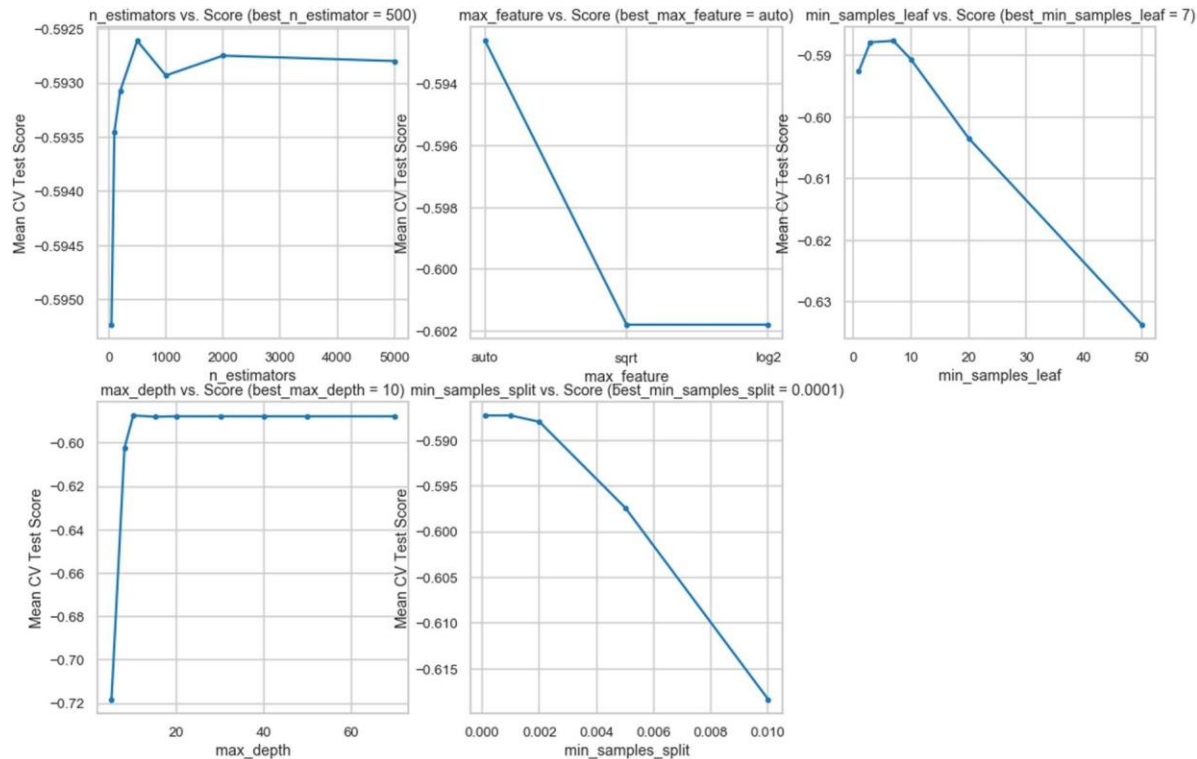
### 8.9. Random Forest (RF1)

Since Random Forest Regression predicts based on decision tree, we expect it to handle categorical features including hours, months, working and non-working day gracefully. In RF1, we use a single model for both Working and Non-Working days and do not use One Hot Encoding to transform the categorical data.

We tune the hyperparameter using the Train Set using GridSearchCV with 5-fold cross validation. Below is the procedure adopted to tune hyperparameter for Random Forest Model

- 8.9.1. First obtain `n_estimators` using default values of the remaining parameters
- 8.9.2. Tune for the `max_features` using the best `n_estimators`
- 8.9.3. Tune for `min_samples_leaf` using the best `n_estimators` and `max_features`
- 8.9.4. Tune for `max_depth` using the best `n_estimators`, `max_features` and `min_samples_leaf`
- 8.9.5. Tune for `min_samples_split` using the best `n_estimators`, `max_features`, `min_samples_leaf` and `max_depth`

Below figures plot the variation of the `cv_score` (which reflects CV Test score)



*Figure 8.9.1: Hyperparameter tuning (one parameter at a time) for Random Forest (`n_estimators`, `max_features`, `min_samples_leaf`, `max_depth` and `min_samples_split`)*

The below table summarizes the optimal hyperparameters for RF1 model

Hyperparameters	Optimal Value
<b>n_estimators</b>	500
<b>max_features</b>	'auto'
<b>min_samples_leaf</b>	7
<b>max_depth</b>	10
<b>min_samples_split</b>	1

*Table 8.9.1: Hyperparameter tuning*

The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.353823	0.387612	<b>0.364732</b>
<b>CV Test</b>	0.425339	0.482957	<b>0.444173</b>
<b>Test</b>	0.393007	0.48925	<b>0.427676</b>

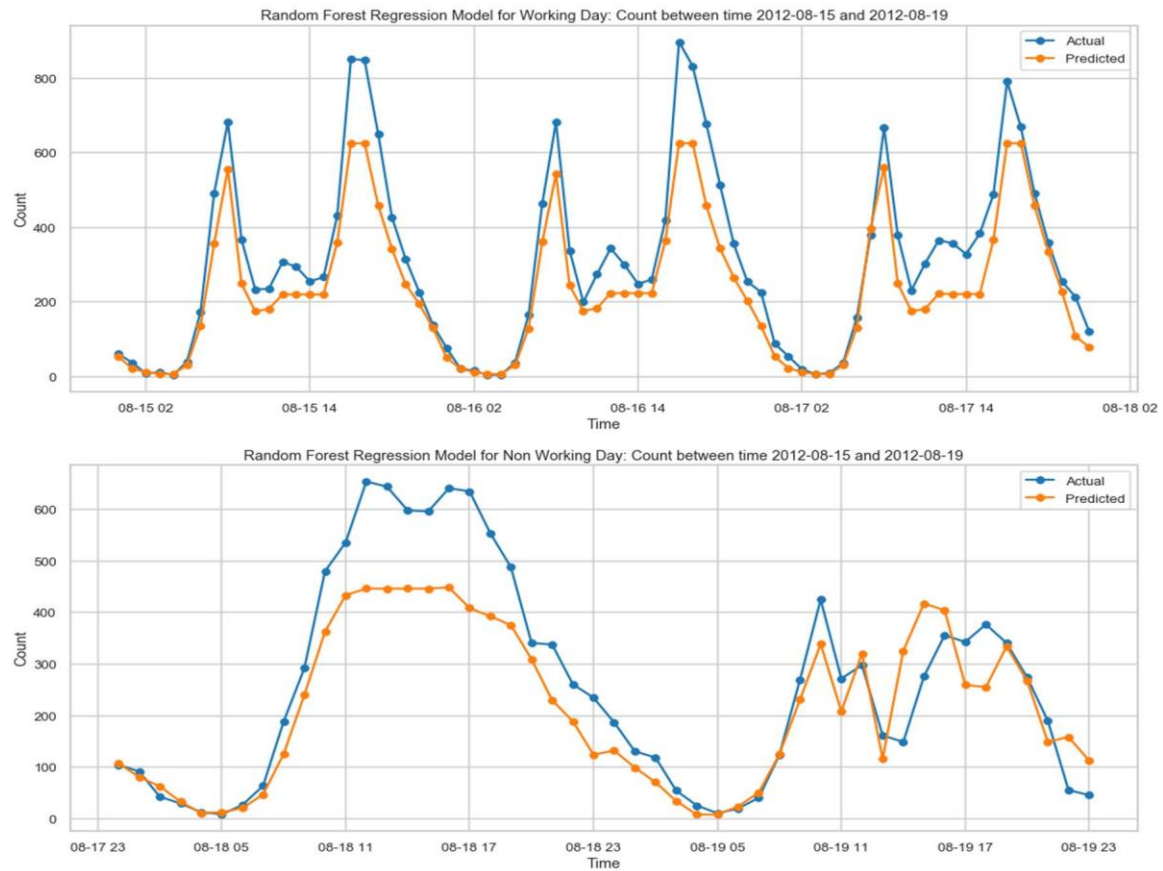
<b>Total Train + Test Time</b>	<b>5.48 sec</b>
--------------------------------	-----------------

*Table 8.9.2: RMSLE Summary for Random Forest Regression (RF1)*

RF1 model performs on par with Linear Regression Model. Also, the train and test RMSLE are similar, which indicates that it is not overfit.

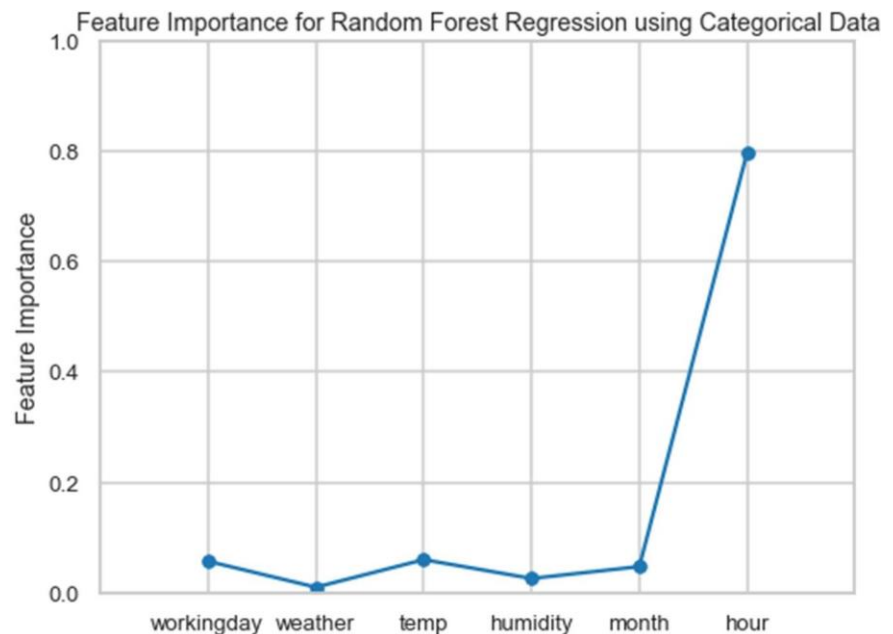


Below plot the actual 'count' value vs. the model predicted value for a few days



in the test data set

*Figure 8.9.2: Random Forest (RF1): Actual vs. Predicted Bike rental count for between 2012-8-15 and 2012-8-19*



*Figure 8.9.3: Feature Importance for Random Forest Model, RF1*

As expected 'hour' feature holds maximum importance. We saw spikes and dips in count value depending on the hour of the day. We also notice that *workingday* feature has marginal importance. In our next model, we will use two separate models for Working and Non-Working days.

#### **8.10. Random Forest (RF2)**

In this Random Forest model (RF2), we use all the features transformed via OneHotEncoding and obtain two separate models for Working and Non-Working days.

Hyperparameter tuning is done as explained in RF1. Figure 7-13 shows the *cv\_score* obtained during hyperparameter tuning via GridSearchCV.

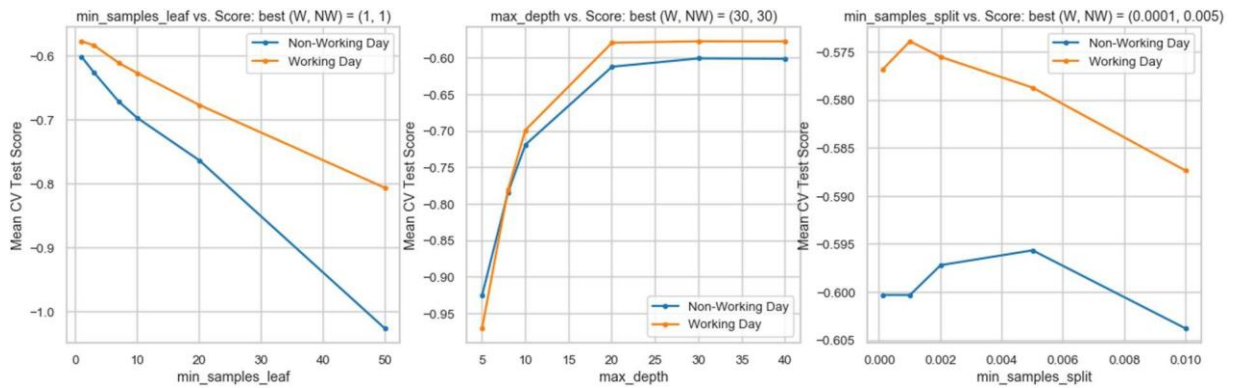


Figure 8.10.1: Hyperparameter tuning for Random Forest (*min\_samples\_leaf*, *max\_depth* and *min\_samples\_split*)

The below table summarizes the optimal hyperparameters for the RF2 model

Hyperparameters	Working Day Model	Non-Working Day Model
<b>n_estimators</b>	2000	200
<b>max_features</b>	'sqrt'	'log2'
<b>min_samples_leaf</b>	1	1
<b>max_depth</b>	30	30
<b>min_samples_split</b>	0.0001	0.005

Table 8.10.1: optimal hyperparameters for the RF2 model

The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.181652	0.378901	<b>0.259992</b>
<b>CV Test</b>			

<b>Test</b>	0.437254	0.53261	<b>0.46918</b>
	0.40938	0.541029	<b>0.457731</b>
<b>Total Train + Test Time</b>			<b>22.4 sec</b>

Table 8.10.2: RMSLE Summary for Random Forest Regression (RF2)

From the table, we see a dip in prediction relative to our earlier Random Forest model (which used Categorical Features and single model for working and non-working days). Average Test RMSLE increases from 0.43 to 0.46. Also note that the Train RMSLE is reasonably lower than Test RMSLE. This possibly is a result of on over-fit model

Below is a feature importance plot for RF2 model.

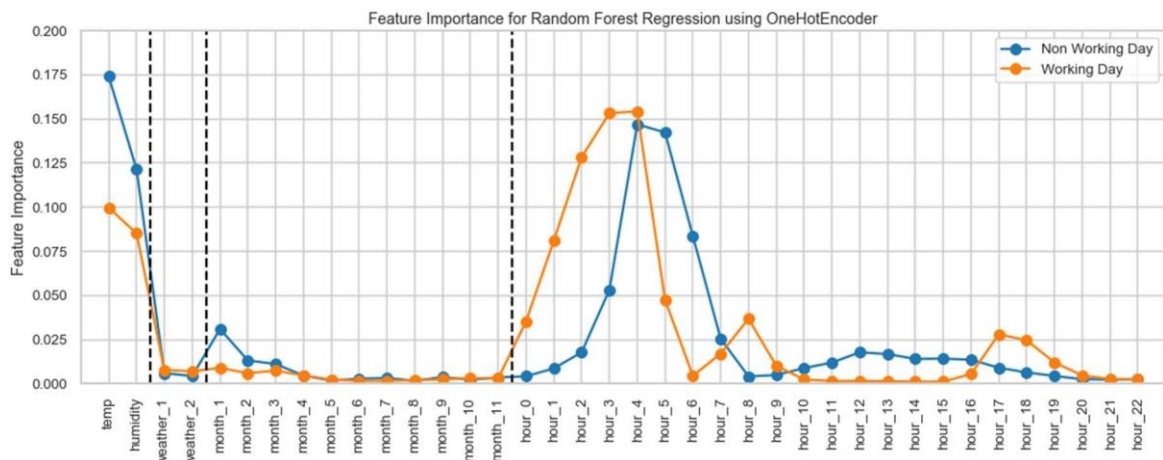


Figure 8.10.2: Feature Importance plot for Random Forest Model, RF2

From the feature importance plot Figure 7-14 we can see that the features for which the average *count* is significantly different from the generic mean *count* (late night or peak hours), tends to hold higher weightage. For example

- Non-working day hour\_x feature set: *hour\_3* to *hour\_6* have high importance as the average count at these hours are very low compared to the daily average. Not splitting based on these features correctly would result in higher MSE. The peak

hours for non-working days are at 1pm and 2pm; correspondingly, we do see relatively higher importance for *hour\_12* and *hour\_13*

- Working day hour\_x feature set: *hour\_1* to *hour\_4* have high importance for similar reasons explained in the previous bullet. *hour\_8*, *hour\_17*, *hour\_18* too hold reasonably high importance since they are the peak hours and tend to have higher than usual count values\

### 8.11. Random Forest (RF3)

In this model, we will not use OneHotEncoding to transform the categorical features but have two separate models to make independent predictions for Working and Non-Working days.

Below are hyperparameter tuning plots obtained via GridSearchCV with cv=5

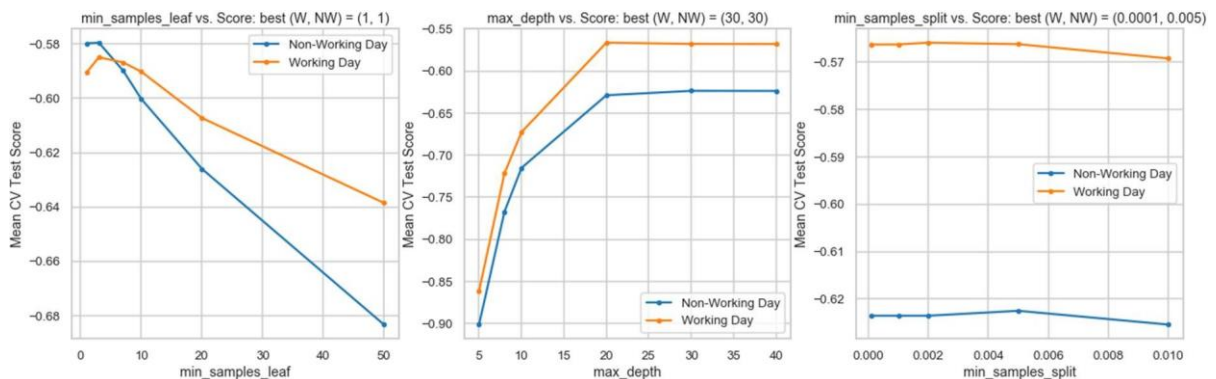


Figure 8.11.1: Hyperparameter tuning for Random Forest (*min\_samples\_leaf*, *max\_depth* and *min\_samples\_split*)

The below table summarizes the optimal hyperparameters for RF3 model

Hyperparameters	Working Day Model	Non-Working Day Model
<b>n_estimators</b>	500	500
<b>max_features</b>	'auto'	'sqrt'
<b>min_samples_leaf</b>	3	3

<b>max_depth</b>	20	30
<b>min_samples_split</b>	0.002	0.005

*Table 8.11.1: optimal hyperparameters for RF3 model*

The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.290989	0.346424	<b>0.309405</b>
<b>CV Test</b>	0.418383	0.496303	<b>0.444237</b>
<b>Test</b>	0.391217	0.524161	<b>0.44026</b>
<b>Total Train + Test Time</b>			<b>6.54 sec</b>

*Table 8.11.2: RMSLE Summary for Random Forest Regression (RF3)*

Average Test RMSLE = 0.44 which is still a bit higher than RF1 Random Forest model. Additionally, the average training RMSLE = 0.31 indicates that this model too might be a bit overfit

Below is the feature importance plot for RF3 model. As RF1, *temp* and *hour* feature hold maximum importance in the regression

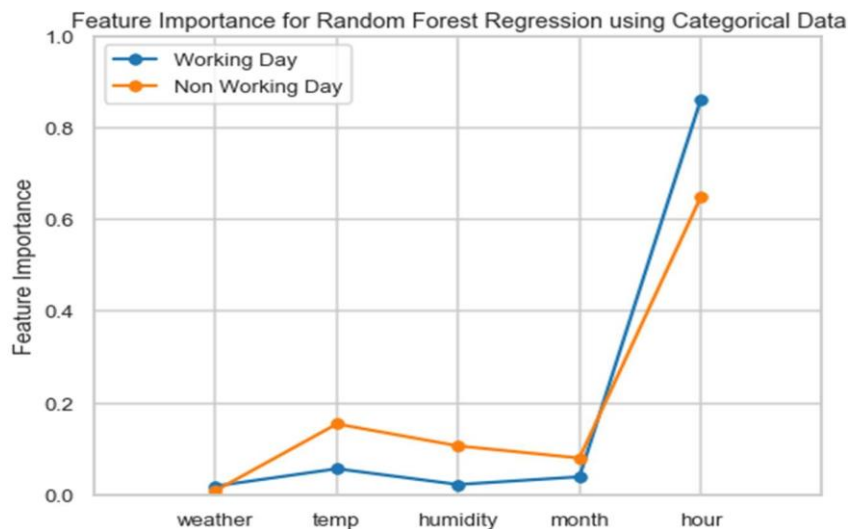


Figure 8.11.2: Feature Importance plot for Random Forest Model, RF3

Random Forest uses an ensemble of 500 decision trees in our earlier model. Let us visualize one of them and see if it makes sense. The entire tree would have several leaf nodes and would be difficult to analyze. Instead, let us limit the `max_depth` to 3 and plot one decision tree (as shown in Figure 7-17)

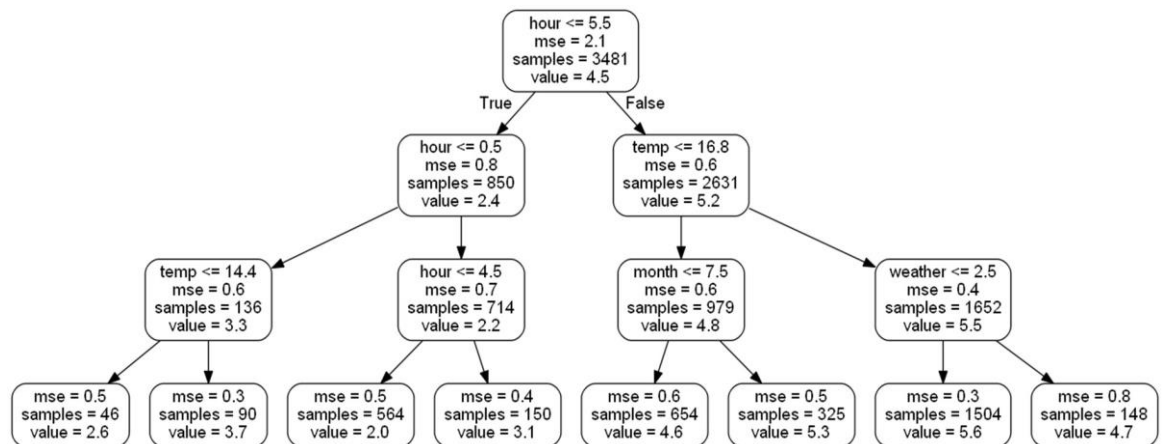


Figure 8.11.3: One of the Decision Trees used in our Random Forest Models

The first split is made based on *hour* feature indicating that it is very important (since that first split at 5.5 hours minimizes the MSE). Note that the value indicated at the leaf node corresponds to the predicted value (which is  $\log(1+\text{count})$ ).

## 8.12. Gradient Boost (GB1)

Gradient Boosting is another powerful regression technique which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. Here, we use a two Gradient Boost models to predict the bike rental count for Working and Non-Working days. We will also use the feature set obtained via OneHotEncoder.

We adopt the below procedure to tune hyperparameter for the Gradient Boost Model

- 8.12.1. Pick a large `n_estimator` = 3000
- 8.12.2. Tune `max_depth`, `learning_rate`, `min_samples_leaf`, and `max_features` via grid search.
- 8.12.3. Increase `n_estimators` even more (5000) and tune `learning_rate` again holding the other parameters fixed.

The below table summarizes the optimal hyperparameters for GB1 model

Hyperparameters	Working Day Model	Non-Working Day Model
<b>n_estimators</b>	3000	3000
<b>learning_rate</b>	0.01	0.005
<b>max_features</b>	1	1
<b>min_samples_leaf</b>	5	5
<b>max_depth</b>	6	6

*Table 8.12.1: optimal hyperparameters for GB1 model*



The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.272817	0.299252	<b>0.281356</b>
<b>CV Test</b>	0.425073	0.559165	<b>0.471153</b>
<b>Test</b>	0.407705	0.550037	<b>0.460327</b>
<b>Total Train + Test Time</b>			<b>42.9 sec</b>

Table 8.12.2: RMSLE Summary for Gradient Boost (GB1)

GB1 has a higher Test RMSLE compared to RF1 model. Also, there is a significant difference between the Train (0.3) and Test RMSLE (0.55), indicating the possibility of an overfit.

Below is the feature importance plot for GB1 model. Unlike the Random Forest model, all the hour\_x features are given almost equal weightage.

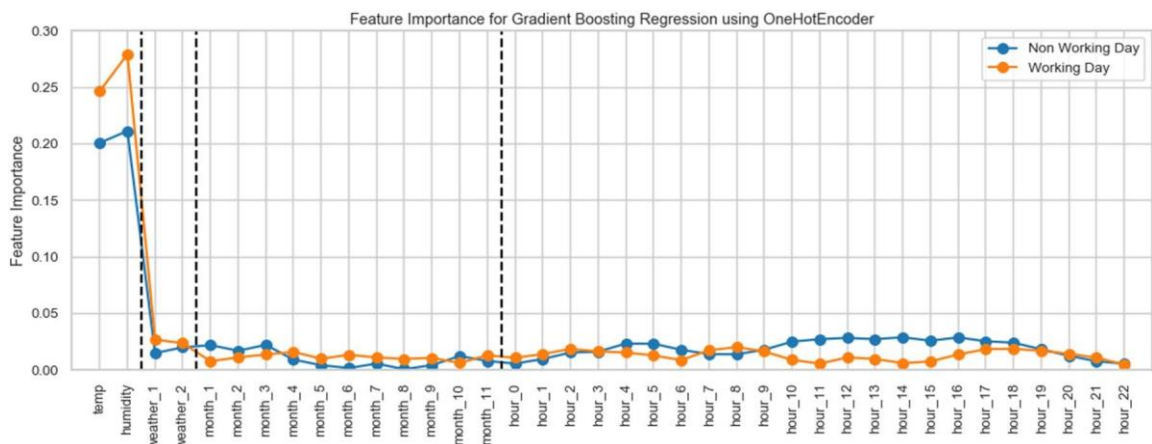


Figure 8.12.1: Feature Importance plot for Gradient Boost Model, GB1

### 8.13. Gradient Boost (GB2)

Let us now try out another variation of the Gradient Boost model. For GB2, we will have 2 models – one for Working days and another for Non-Working days. Also, we will not use OneHotEncoding on categorical features.

The below table summarizes the optimal hyperparameters for GB2 model

Hyperparameters	Working Day Model	Non-Working Day Model
<b>n_estimators</b>	3000	3000
<b>learning_rate</b>	0.01	0.005
<b>max_features</b>	4	4
<b>min_samples_leaf</b>	1	1
<b>max_depth</b>	21	21

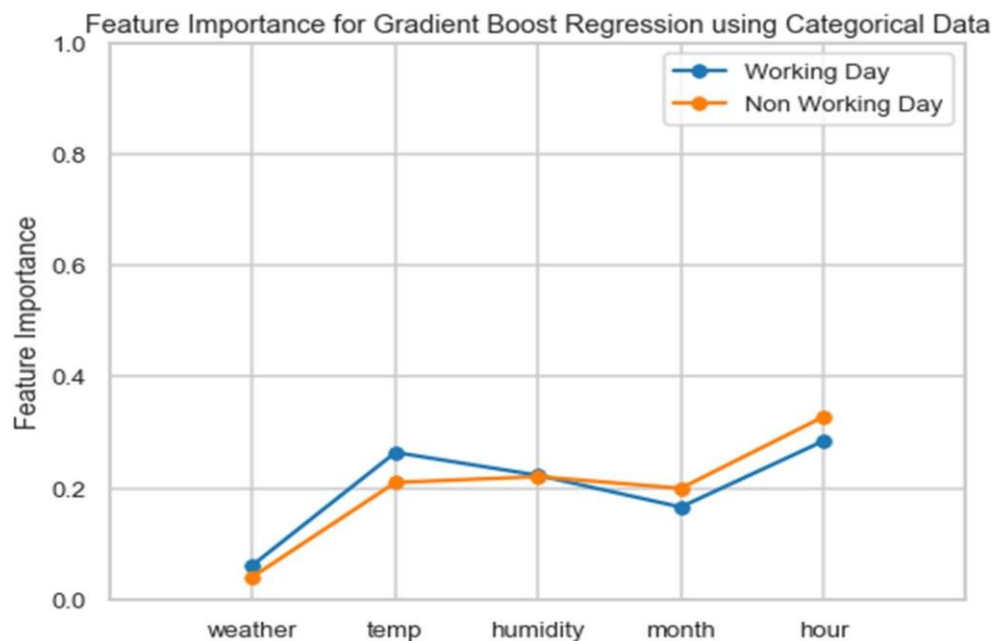
*Table 8.13.1: optimal hyperparameters for GB2 model*

The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.312417	0.334973	<b>0.319646</b>
<b>CV Test</b>	0.404077	0.501422	<b>0.436876</b>
<b>Test</b>	0.387938	0.493577	<b>0.426262</b>
<b>Total Train + Test Time</b>			<b>10.1 sec</b>

*Table 8.13.2: RMSLE Summary for Gradient Boost (GB2)*

GB2 results in a low test RMSLE (0.43) but a significantly lower train RMSLE (0.32) seem to suggest that this model too might be overfit.



*Figure 8.13.1: Feature Importance plot for Gradient Boost Model, GB2*

### 8.14. Adaboost

Adaboost is another popular boosting algorithm where we convert the predictions of several weak learning algorithms into a strong one by weighted combining. Let us u

We have two hyperparameters for Adaboost and the below table summarizes the optimal hyperparameters obtained via GridSearchCV with cv=5

Hyperparameters	Working Day Model	Non-Working Day Model
<b>n_estimators</b>	5000	5000
<b>learning_rate</b>	0.001	0.001

Table 8.14.1 : optimal hyperparameters obtained via GridSearchCV

The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.604791	0.558883	<b>0.590809</b>
<b>CV Test</b>	0.621749	0.5989	<b>0.614691</b>
<b>Test</b>	0.579941	0.651514	<b>0.604868</b>
<b>Total Train + Test Time</b>			<b>59.2</b>

Table 8.14.2: RMSLE Summary for Adaboost (AB)

AdaBoost model has higher RMSLE and doesn't perform as well as the other 2 Ensemble methods (RF and GB). Below is a feature importance plot for the Adaboost model which prioritizes like other ensemble models (RF and GB)

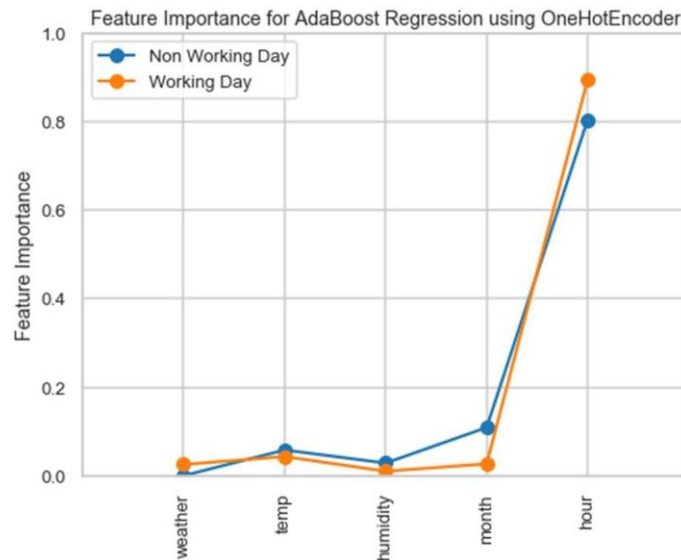


Figure 8.14.1: Feature Importance plot for Adaboost Model

Below, we plot the estimator errors and the weights given to each of the individual 5000 estimators used for our Adaboost model.

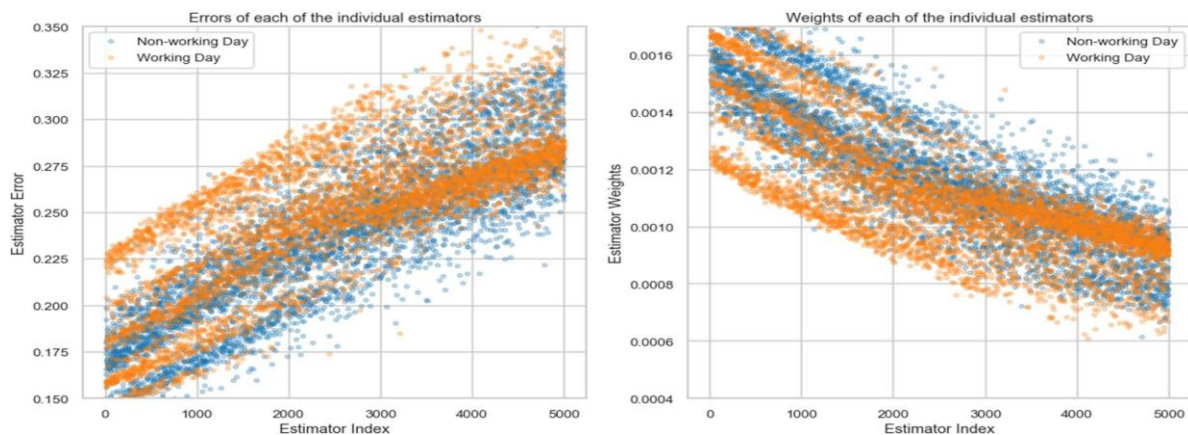


Figure 8.14.2: Estimator Error and Estimator Weights for each of the 5000 individual estimators for Working and Non-Working day Adaboost models

From the above plot, it AdaBoost might not be a good model for this problem. Every subsequent estimator seems to result in an increase in the error (and as a result, we give it a lower weight for our final combined estimator model)

### 8.15. Stacking via Linear Regression

Model stacking is an efficient ensemble method in which the predictions, generated by using various machine learning algorithms, are used as inputs in a second-layer learning algorithm. This second-layer algorithm is trained to optimally combine the model predictions to form a new set of predictions. Next few models will apply the stacking concept by using the predictions obtained from above individual models (7.6 to 7.14) as meta-features to build a new ensemble model.

First let us use Linear Regression as the second layer model - this would estimate the weight for each of the individual model predictions by minimizing the least square errors.

The RMSLE and the model train+test time summary for this stacked Linear Regression model is given in the below table.

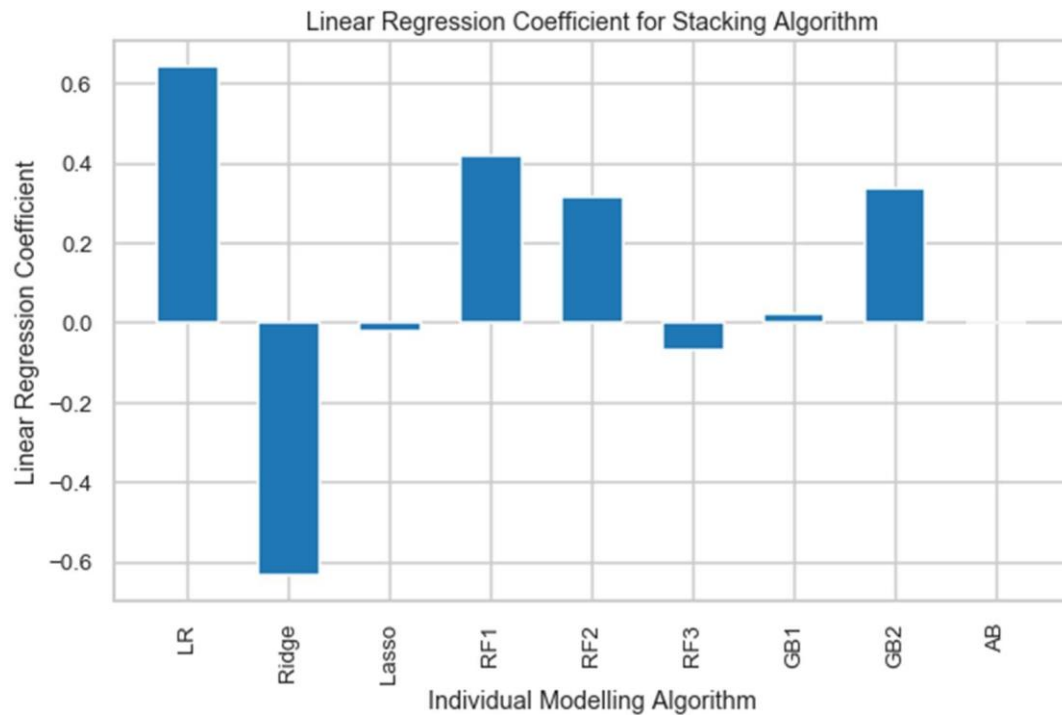
	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.436751	0.499553	<b>0.457331</b>
<b>Test</b>	0.402243	0.48662	<b>0.432355</b>
<b>Total Train + Test Time</b>			<b>584 secs</b>

*Table 8.15.1: RMSLE Summary for Stacking via Linear Regression*

Stacking using Linear Regression results in poorer test performance (0.46 RMSLE) compared to few ensemble individual models (e.g. RF1 with 0.43 RMSLE).

Also note that the total train+test time for stacking models are usually very high. This is because, we must train+test every individual model first, to obtain the predicted value which then becomes the feature set for the stacked model.

Below we plot the Coefficients obtained for our Level 2 Linear Regression model.



*Figure 8.15.1: Linear Regression Coefficients for Stacking Algorithm*

- We notice that the Lasso and GB1 have been given very low values indicating that those models were probably not very accurate.
- LR and Ridge have equal opposite signed coefficients. Since these two are highly correlated (as we saw by their coefficients in Figure 7-9: Feature Coefficient Comparison for the Linear Models (Linear Regression, Ridge and Lasso))
- So, to sum up, this Linear Regression seem to primarily use the predicted values from RF1, RF2 and GB2.

### 8.16. Stacking via Random Forest

Let us use Random Forest to model our stacking model. The below table summarizes the optimal hyperparameters for this stacked Random Forest model

Hyperparameters	Optimal Value
<b>n_estimators</b>	1000
<b>max_features</b>	'sqrt'
<b>min_samples_leaf</b>	10
<b>max_depth</b>	5

*Table 8.16.1: optimal hyperparameters for this stacked Random Forest model*

The RMSLE and the model train+test time summary for the model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.422611	0.472015	<b>0.438669</b>
<b>Test</b>	0.3899	0.494256	<b>0.427714</b>
<b>Total Train + Test Time</b>			<b>593 secs</b>

*Table 8.16.2: RMSLE Summary for Stacking via Random Forest*

Stacking using Linear Regression results in poorer test performance (0.46 RMSLE) compared to fewensemble individual models (e.g. Random Forest 3.5.1 with 0.43 RMSLE)

Below is feature importance plot for our stacked Random Forest model.



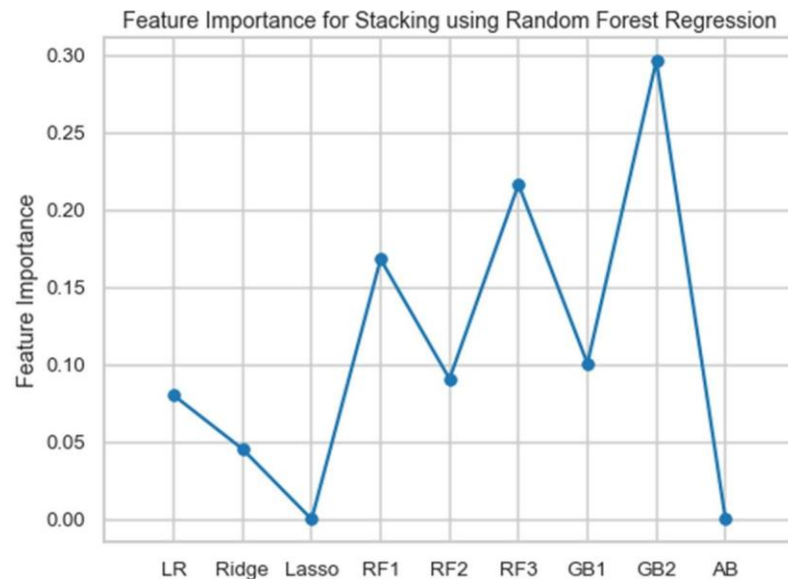


Figure 8.16.1: Feature importance plot for Stacking using Random Forest Regression

Lasso and Adaboost model have zero importance (which agrees with our RMSLE analysis). Predictions from GB2, RF3 and RF1 seem to be of the most importance.

### 8.17. Stacking via Gradient Boost

Finally, let us use Gradient Boost to model our stacking model. The below table summarizes the optimal hyperparameters for this stacked Gradient Boost model.

Hyperparameters	Optimal Value
<b>n_estimators</b>	3000
<b>learning_rate</b>	0.01
<b>max_features</b>	0.3
<b>min_samples_leaf</b>	100
<b>max_depth</b>	4

Table 8.17.1: Optimal hyperparameters for this stacked Gradient Boost model.

The RMSLE and the model train+test time summary for this stacked Gradient Boost model is given in the below table

	Working Day	Non-Working Day	All Combined
<b>Train</b>	0.408437	0.458788	<b>0.424835</b>
<b>Test</b>	0.393135	0.494927	<b>0.429937</b>
<b>Total Train + Test Time</b>			<b>592 secs</b>

Table 8.17.2: RMSLE Summary for Stacking via Gradient Boost

Below is feature importance plot for our stacked Gradient Boost model.

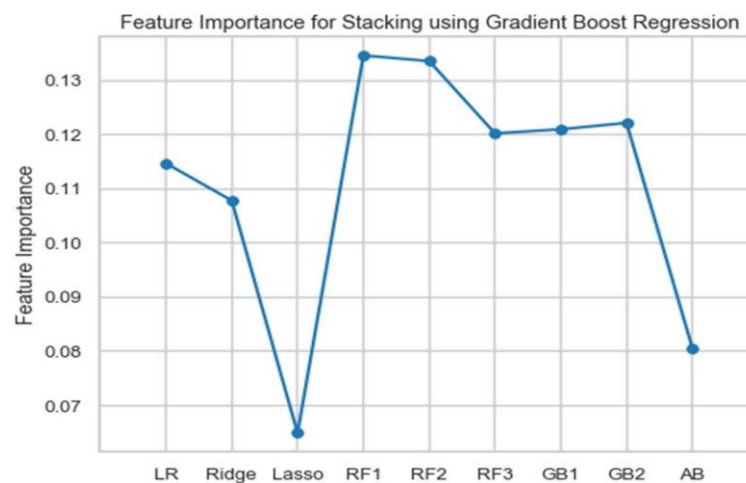


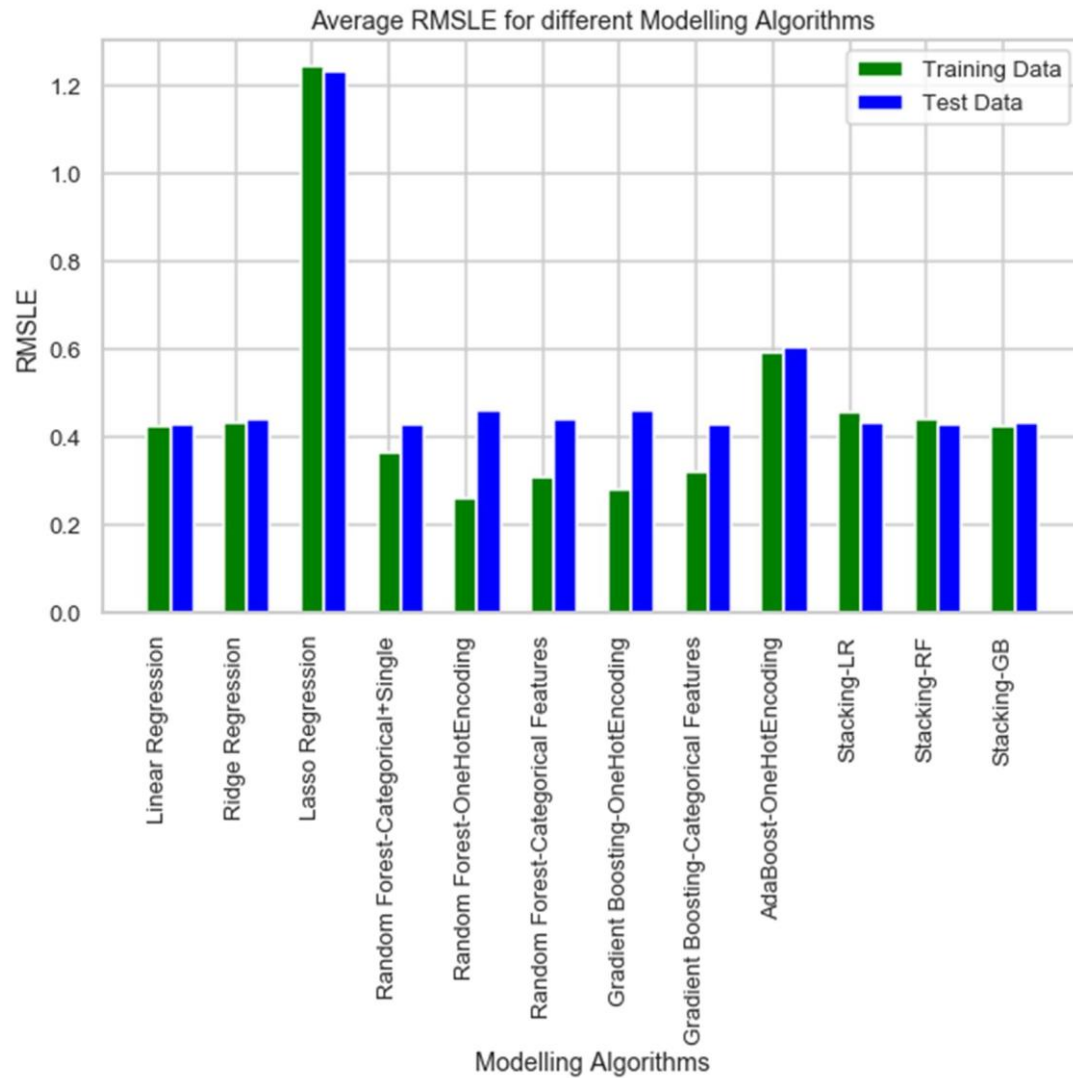
Figure 8.17.1: Feature importance plot for Stacking using Gradient Boost Regression

As with the Random Forest stacked model, Lasso and Adaboost gets the lowest importance. All the remaining models are given more or less equal priority here.

## 9. SUMMARY AND CONCLUSION

### 9.1.RMSLE and Train+Test time

We summarize the average Test RMSLE and Train+Test time for every model



we used in the below barcharts

*Figure 9.1.1: Test RMSLE for all the Modeling algorithms*

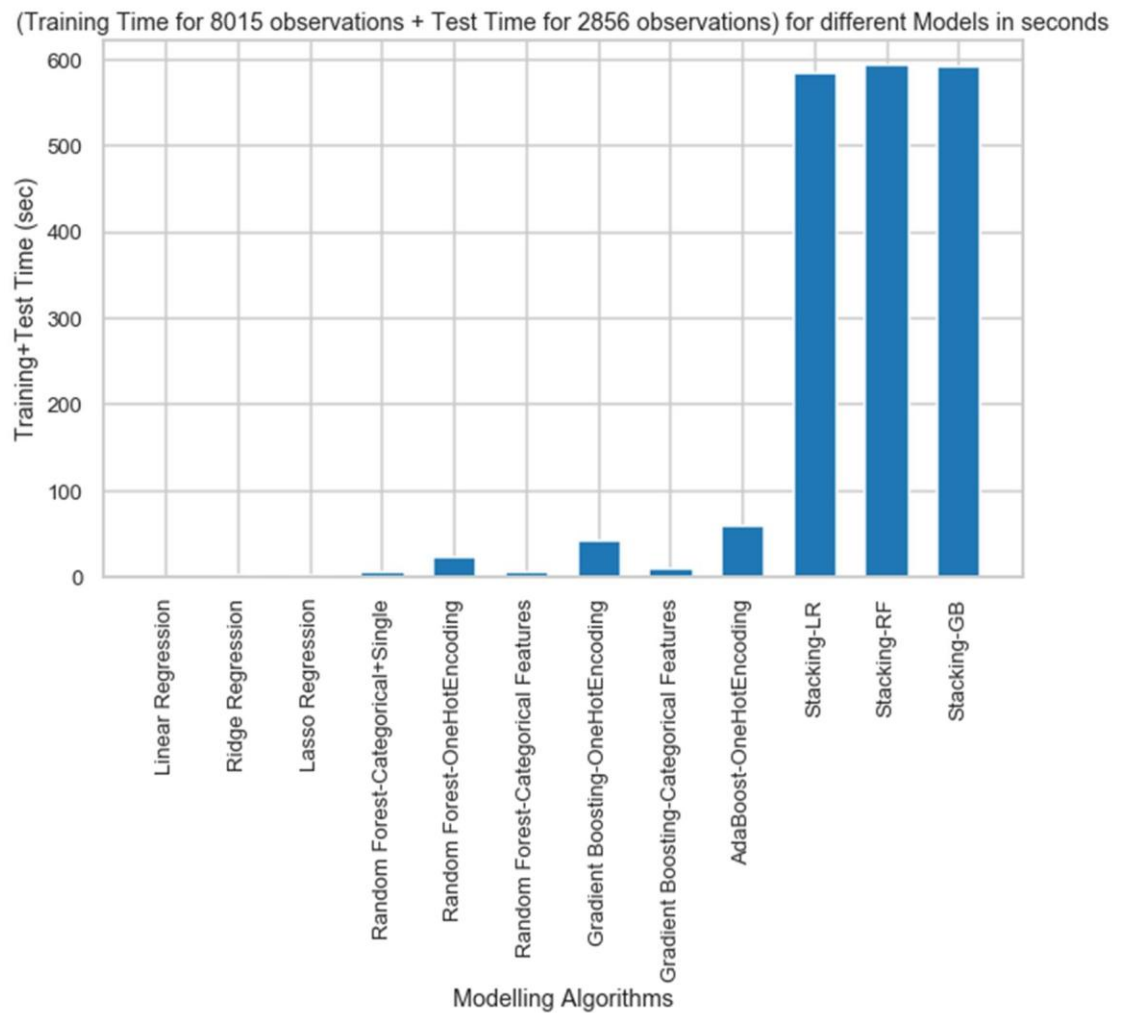


Figure 9.1.2: Train+Test time for all the considered Models

Based on the above results, let us use RF1, i.e., Single Model Random Forest Regressor Model (with categorical feature) based on the below couple of reasons

- 9.1.1. One of the lowest Average Test RMSLE
- 9.1.2. Single model which works for both working and non-working days
- 9.1.3. Not an Overfit model (Train Average RMSLE is closer to Test Average RMSLE)
- 9.1.4. Lesser to train due to lesser number of features

---

Other reasonable choice would be the Linear Regression Model.

## 9.2. Summary

In this report, we train a model to predict the number of bike rentals at any hour of the year given the weather conditions. The data set was obtained from the Capital Bikeshare program in Washington, D.C. which contained the historical bike usage pattern with weather data spanning two years.

First, we do Exploratory Data Analysis on the data set. We look for missing data values (none were found) and outliers and appropriately modify them. We also perform correlation analysis to extract out the important and relevant feature set and later perform feature engineering to modify few existing columns and drop out irrelevant ones.

We then look at several popular individual models from simple ones like Linear Regressor and Regularization Models (Ridge and Lasso) to more complicated ensemble ones like Random Forest, Gradient Boost and Adaboost. Additionally, few options for model formulation were tried - 1. A single unified model for working and non-working days, 2. Two separate models for working and non-working days, 3. Using OneHotEncoding to get Binary Vector representation of Categorical Features and 4. Using Categorical Features as provided. Finally, we also tried stacking algorithms where the predictions from the level 1 individual models were used as meta-features into a second level model (Linear Regressor, RandomForest and Gradient Boost) to further enhance the predicting capabilities.

The labeled data set provided comprised of the first 19 days of each month while the Kaggle Test Data comprised of rental information from 20th to the end of the month. Hyperparameters were tuned using GridSearchCV cross validation

using 5 folds on part of the provided training data set (first 14 days of each month). The remaining data (15th to 19th of each month) were used as hold out set to test our model performance. Of all the methods and models, we found Random Forest Ensemble method using a Single Model and Categorical Feature set an ideal choice with train and test scores of 0.36 and 0.427, respectively. The training and test time for the chosen model are very reasonable (~23 seconds for total train+test observation size = 10871). The chosen model yields a score of 0.49 on Kaggle Test Data

### 9.2.1. Data Exploration Conclusions

In this project, we explored several factors that influenced bike rental count. Below is a quick summary of exploratory data analysis

9.2.1.1. Working or Non-working Day We see 2 rental patterns across the day in bike rentals count - first for a Working Day where the rental count high at peak office hours (8am and 5pm) and the second for a Non-working day where rental count is more or less uniform across the day with a peak at around noon.

9.2.1.2. Hour of the day: Bike rental count is mostly correlated with the time of the day. As indicated above, the count reaches a high point during peak hours on a working day and is mostly uniform during the day on a non-working day

9.2.1.3. Casual and Registered Users: While most casual users are likely to be tourists whose rental count is high during non-working days, most registered users are most likely city natives whose rental count is high during working days

9.2.1.4. Temperature: People generally prefer to bike at moderate to high temperatures. We see highest rental counts between 32 to 36 degrees Celsius

9.2.1.5. Season: We see highest number bike rentals in Fall (July to

September) and Summer (April to June) Seasons and the lowest in Spring (January to March) season

9.2.1.6. Weather: As one would expect, we see highest number of bike rentals on a clear day and the lowest on a snowy or rainy day

9.2.1.7. Humidity: With increasing humidity, we see decrease in the number of bike rental count.

### **9.2.2. Modeling Conclusions**

We used 6 Regression Models to predict the bike rental count at any hour of the day - Linear Regression, Ridge, Lasso, Random Forest, Gradient Boost and Adaboost. Using the predictions made by these level 1 individual models as features, we trained 3 level 2 stacking algorithms (Linear Regression, Random Forest and Gradient Boost) to make more refined predictions. Below is a summary of the model performances

9.2.2.1. Of all the models, we found a simple Random Forest Model providing the best/lowest RMSLE score.

9.2.2.2. Stacking individual models didn't provide any improvement over the best individual model

9.2.2.3. Having separate models for Working days and Non-working days didn't provide any improvement in prediction accuracy

9.2.2.4. 'Hour' of the day holds most importance among all the features for prediction

### **9.2.3. Limitations and Scope for Model Improvements**

Below are few limitations in this analysis and ideas to improve model prediction accuracy

9.2.3.1. Since casual + registered = total count, we just predicted the total bike rental count by ignoring the casual and registered user information. Another (possibly better) method would be to train separate models

---

for casual and registered users and add the two

9.2.3.2. Windspeed wasn't used due to very low correlation. This might have been due to several instances where  $\text{windspeed} = 0$ . One possible method could be to first estimate those windspeed and then use it as a feature to estimate count.

9.2.3.3. One limitation in the provided training data set is that it lacked data with extreme weather condition data ( $\text{weather} = 4$ ). Hence, we had to modify it to  $\text{weather} = 3$



---

**REFERENCE**

- i. Shaheen, S. A., & Guzman, S. (2011). Bike-sharing: History, impacts, models of provision, and future. *Journal of Public Transportation*, 14(4), 1-21.
- ii. Martín, J. C., Montero, J., & López-García, B. (2017). An analysis of bike-sharing usage and rebalancing: Evidence from Barcelona. *Transportation Research Part A: Policy and Practice*, 103, 370-384.
- iii. Buehler, R., & Pucher, J. (2012). The adoption of bike-sharing in European cities: Cross-national comparisons. *Transportation Research Part A: Policy and Practice*, 46(4), 984-999.
- iv. Shaheen, S. A., Martin, E., & Chan, N. (2012). Bike sharing and the economy of cities. *Transportation Research Record: Journal of the Transportation Research Board*, 2314(1), 25-36.
- v. Buldeo Rai, H., Niederhuber, K., & Scarpetta, S. (2017). Designing bike-sharing systems for sustainable mobility: Lessons learned from a comparative case study. *Transportation Research Part A: Policy and Practice*, 97, 320-336.