

Capstone Project-2

AI

Bike Sharing Demand Prediction

(Supervised Machine Learning regression)

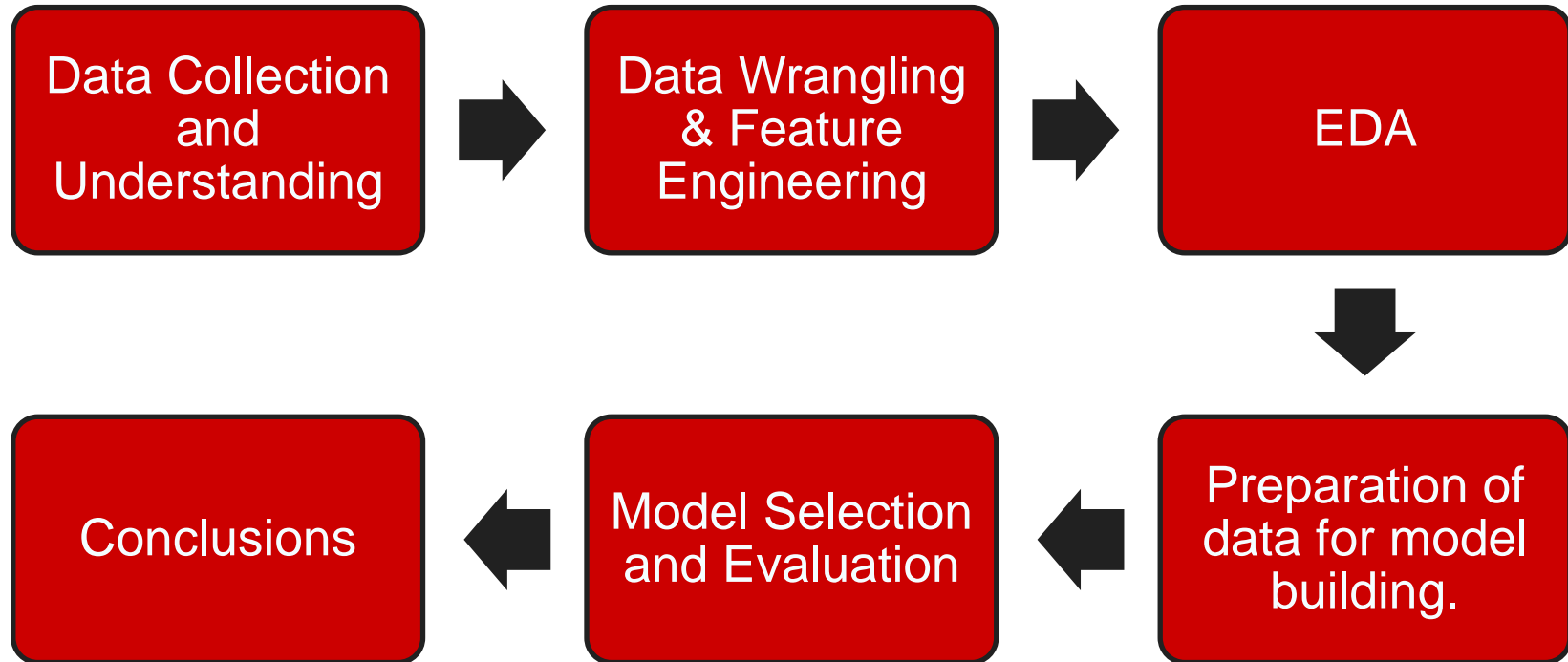
BY

Akash S. Kawade

❖ Problem Statement:

- **Currently Rental bikes are introduced in many urban cities for the enhancement of mobility comfort. The client is Seoul Bike, which participates in a bike share program in Seoul, South Korea. An accurate prediction of bike count is critical to the success of the Seoul bike share program. It is important to make the rental bike available and accessible to the public at the right time as it lessens the waiting time. Eventually, providing the city with a stable supply of rental bikes becomes a major concern.**
- **The final aim of this project is the prediction of bike count required at each hour for the stable supply of rental bikes.**

➤ So we will divide our work flow into following steps.



❖ Data Collection and Understanding:

- We had a Seoul Bike Data for our analysis and model building
- The dataset contains weather information (Temperature, Humidity, Wind speed, Visibility, Dew point, Solar radiation, Snowfall, Rainfall), the number of bikes rented per hour and date information.
- In this we had total 8760 observations and 14 features including target variable.

Data Description:

Date : year-month-day.

Hour - Hour of the day.

Temperature-Temperature in Celsius.

Humidity - %.

Wind speed - m/s.

Visibility - m.

Dew point temperature - Celsius.

Solar radiation - MJ/m².

Rainfall - mm.

Snowfall - cm.

Seasons - Winter, Spring, Summer, Autumn.

Holiday - Holiday/No holiday.

Functional Day - NoFunc(Non Functional Hours), Fun(Functional hours).

Rented Bike count - Count of bikes rented at each hour (Target Variable i.e Y variable).

❖ Data Wrangling and Feature Engineering:

As we know we had 8760 observations and 14 features.

➤ **Categorical Features:** Seasons, Holiday and Functioning day.

➤ **Numerical Columns:**

Date, Hour, Temperature, Humidity, Wind speed, Visibility, Dew point temperature, Solar radiation, Rainfall, Snowfall, Rented Bike count .

Rename Columns: We renamed columns because they had units mentioned in brackets and was difficult to copy the feature name while working.

```
bike_df.columns
```

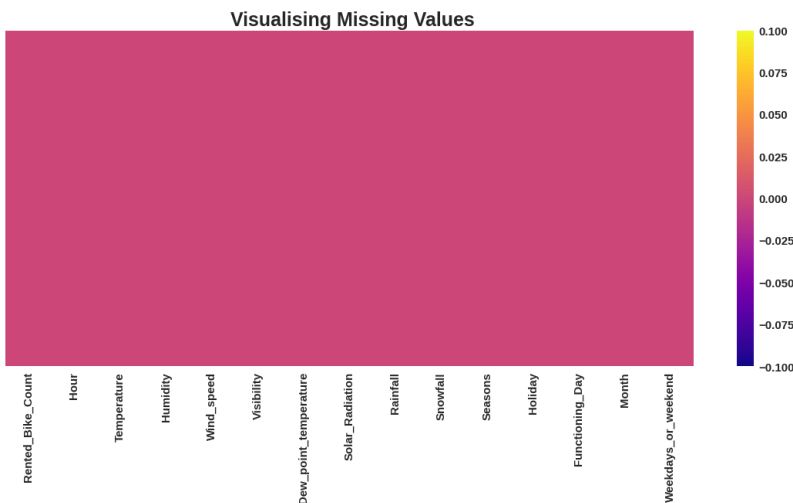
```
Index(['Date', 'Rented Bike Count', 'Hour', 'Temperature(°C)', 'Humidity(%)',  
      'Wind speed (m/s)', 'Visibility (10m)', 'Dew point temperature(°C)',  
      'Solar Radiation (MJ/m2)', 'Rainfall(mm)', 'Snowfall (cm)', 'Seasons',  
      'Holiday', 'Functioning Day'],  
      dtype='object')
```



```
#Since the variables having units with name, renaming columns for better variable analysis.  
bike_df.rename(columns={'Rented Bike Count':'Rented_Bike_Count', 'Temperature(°C)':'Temperature',  
                        'Visibility (10m)':'Visibility', 'Dew point temperature(°C)':'Dew_point_t',  
                        'Rainfall(mm)':'Rainfall', 'Snowfall (cm)':'Snowfall', 'Functioning Day':
```

❖ Data Wrangling and Feature Engineering:

- We had zero null values in our dataset.
- Zero Duplicate entries found.
- We changed the data type of Date column from 'object' to 'datetime64[ns]'. This was done for feature engineering.
- We Created two new columns with the help of Date column 'Month' and 'Day'. Which were further used for EDA. And later we dropped Date column.



```
# Change The datatype of Date columns to extract 'Month', 'Day', "year".
bike_df['Date']=bike_df['Date'].astype('datetime64[ns]')
```

```
# checking Duplicate rows in our BikeData.
duplicates=bike_df.duplicated().sum()
print(f"We have {duplicates} rows in our Bike Data.")
# No duplicate rows found
```

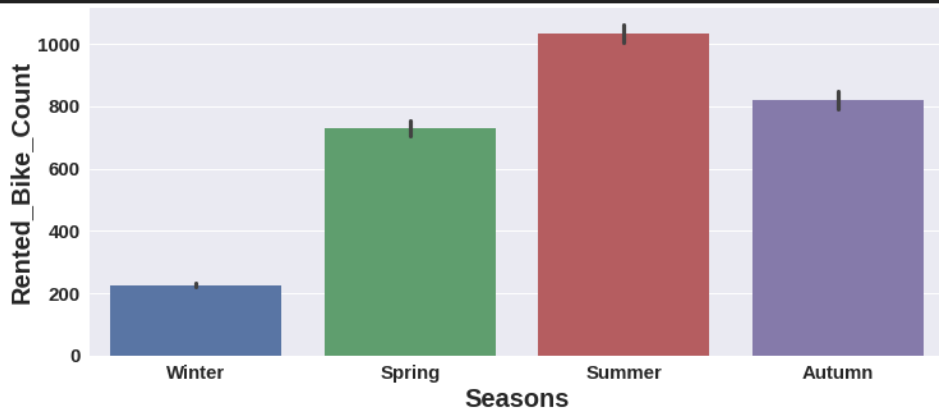
We have 0 rows in our Bike Data.

```
[ ] # Creating new columns 'Month', 'Year', 'Day'.
bike_df['Month']=bike_df['Date'].dt.month

bike_df['Day']=bike_df['Date'].dt.day_name()
```

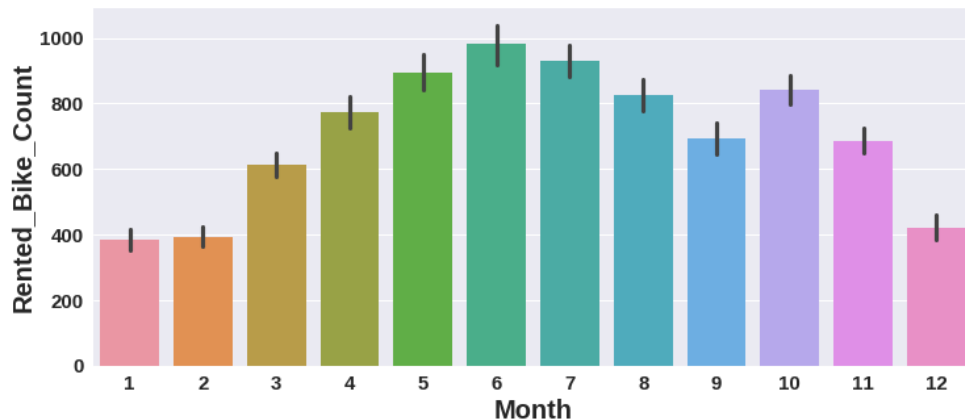


EDA (Exploratory Data Analysis):



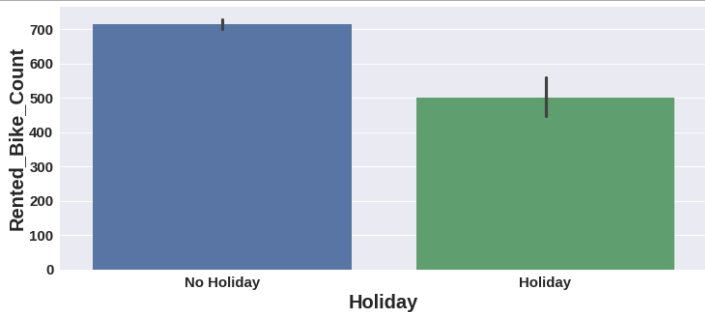
Relation of rented bike count with categorical features:

Summer season had the highest Bike Rent Count. People are more likely to take rented bikes in summer. Bike rentals in winter is very less compared to other seasons.



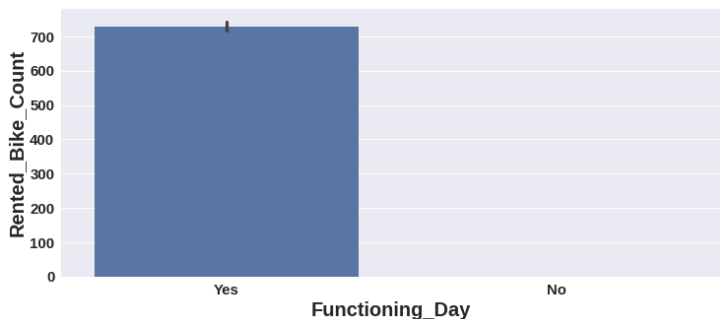
From March Bike Rent Count started increasing and it was highest in June.

❖ EDA (Exploratory Data Analysis):

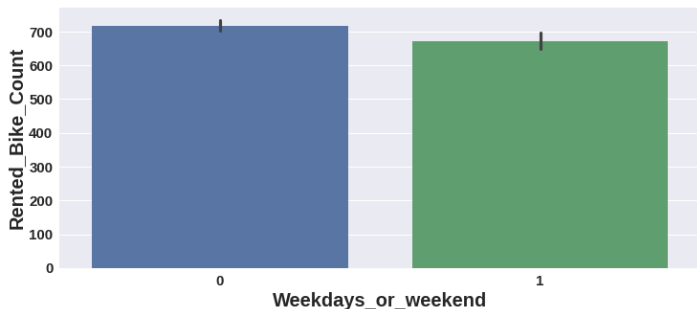


Conclusions:

High number of bikes were rented on No Holidays. Which is almost 700 bikes.



Zero Bikes were rented on no functioning day. More than 700 bikes rented on functioning day.

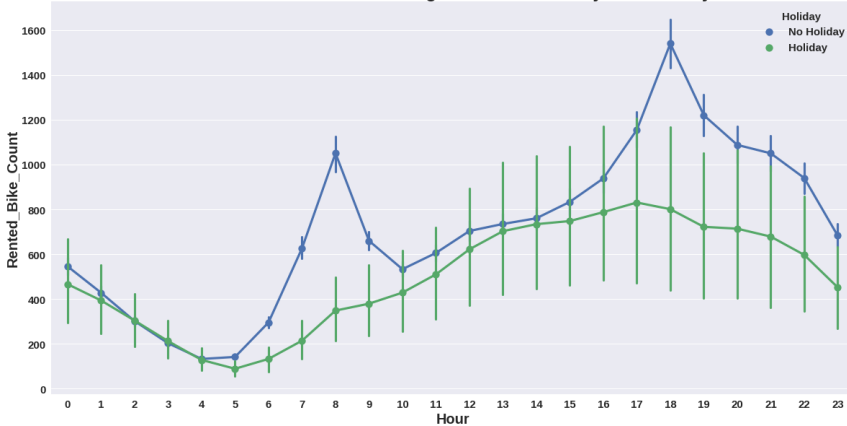


More than 700 bikes were rented on weekdays. On weekdays, almost 650 bikes were rented.

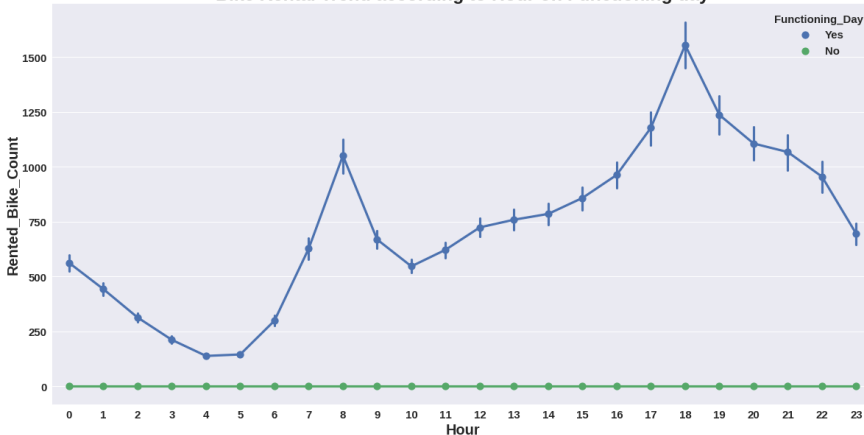
❖ EDA (Exploratory Data Analysis):

Bike Rent Trend according to hour in different scenarios.

Bike Rental Trend according to Hour on Holiday / No Holiday



Bike Rental Trend according to Hour on Functioning day



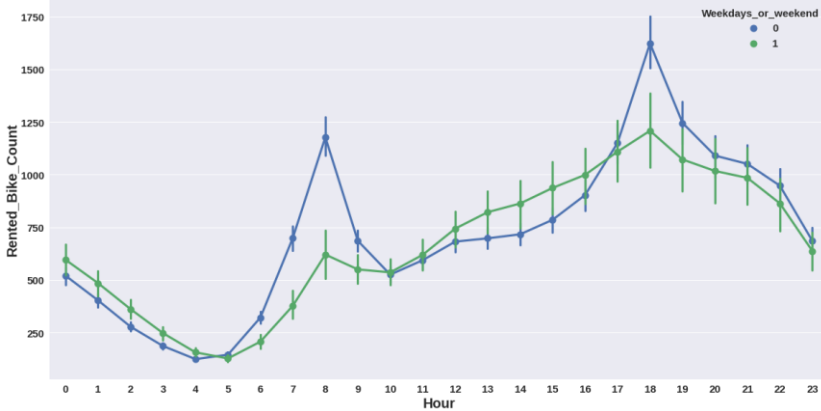
Observations:

- 1) Here we observed that, Bike rental trend according to hours is almost similar in all scenarios.
- 2) There is sudden peak between 6/7AM to 10 AM. Office /College going time could be the reason for this sudden peak on NO Holiday. But on Holiday the case is different, very less bike rentals happened.
- 3) Again there is peak between 4PM to 7 PM. may be its office leaving time for the above people.(NO Holiday).
- 4) Here the trend for functioning day is same as of No holiday. Only the difference is on No functioning day there were zero bike rentals.

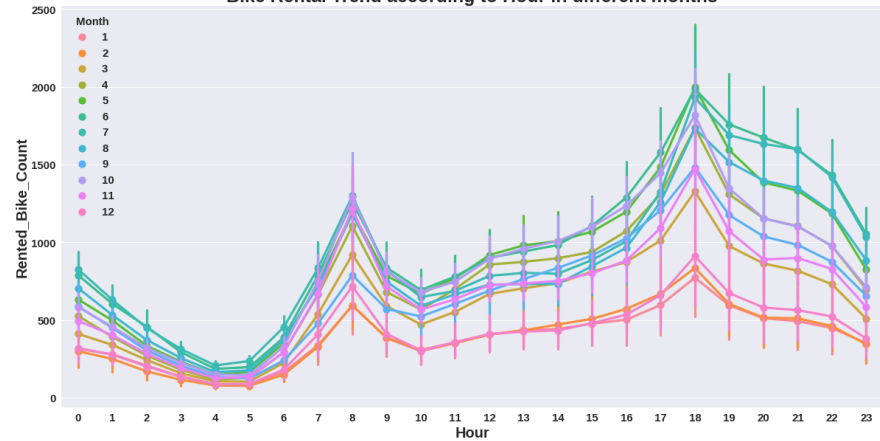
❖ EDA (Exploratory Data Analysis):

Bike Rent Trend according to hour in different scenarios.

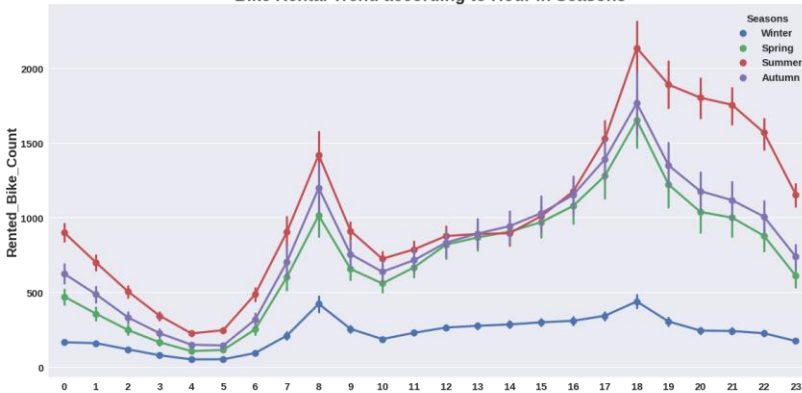
Bike Rental Trend according to Hour in Weekdays_or_weekend.



Bike Rental Trend according to Hour in different months

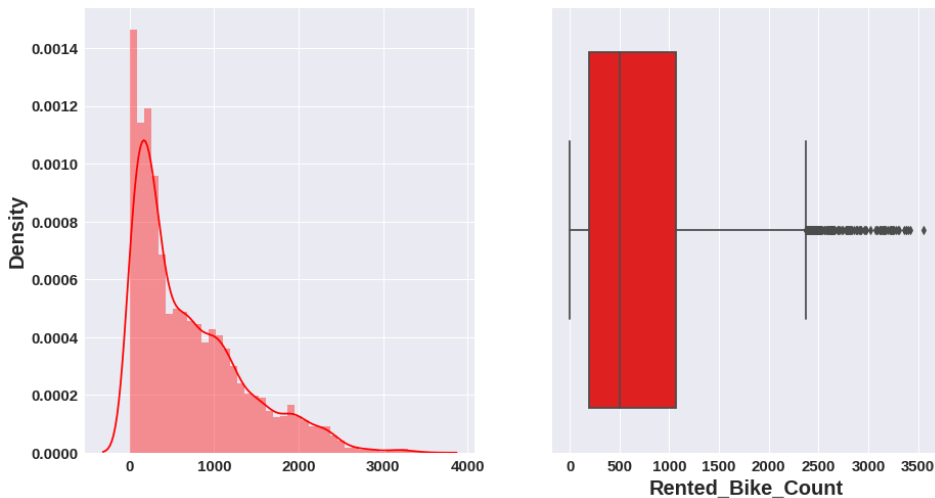


Bike Rental Trend according to Hour in Seasons

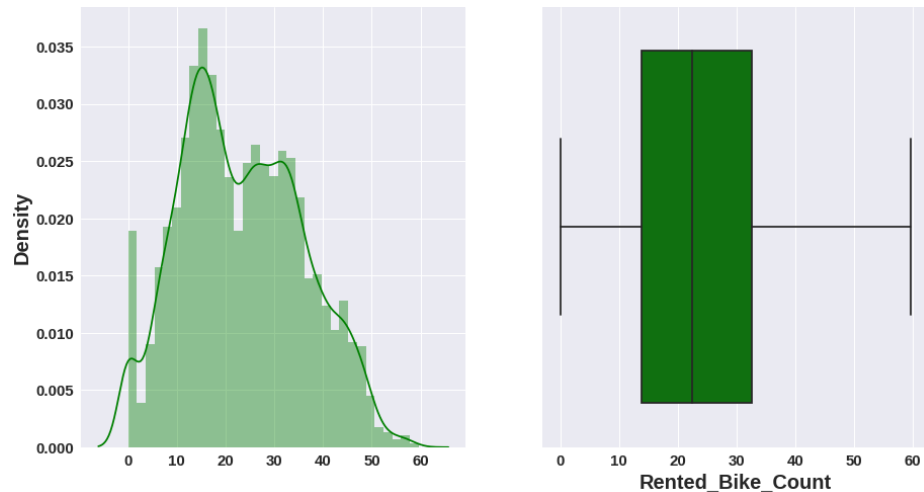


❖ EDA (Exploratory Data Analysis):

Distribution of target variable- Bike Rent Count

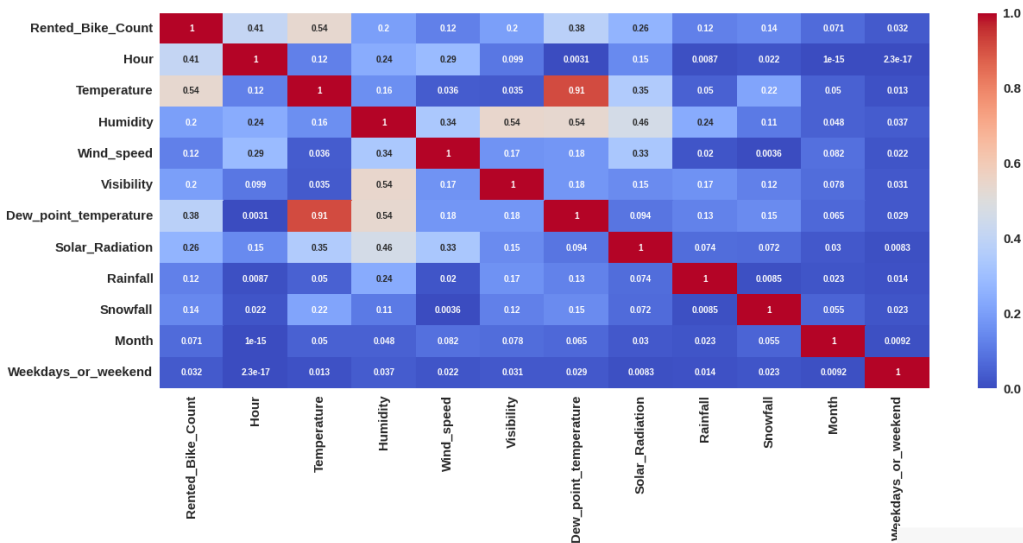


Distribution is rightly skewed and some outliers are observed.



To normalize the distribution we applied square root method. After normalization no outliers were found.

❖ Preparation of data for model building:



➤ With the heat map we dropped highly correlated variables. As we can see Temperature and Dew point temperature are 91 % correlated. So we dropped the Dew point temperature because it has very low correlation with our target variable as compared to temperature.

➤ Later by using variation inflation factor we dropped 'Visibility' and 'Humidity' features as they had VIF value more than 5.

➤ Next we created dummy variables for categorical Seasons column and did mapping with 0 and 1 for holiday and functioning column.

➤ Thus we prepared our data for model building.

```
[ ] # Createing dummy variables
df=pd.get_dummies(df,columns=['Seasons'],prefix='Seasons',drop_first=True)
```

```
[ ] # Labeling for holiday=1 and no holiday=0
df['Holiday']=df['Holiday'].map({'No Holiday':0, 'Holiday':1})
```

```
[ ] ## Labeling for Yes=1 and no No=0
df['Functioning_Day']=df['Functioning_Day'].map({'Yes':1, 'No':0})
```

❖ Model Selection and Evaluation:

As this is the regression problem we are trying to predict continuous value. For this we used following regression models.

- Linear Regression
- Lasso regression (regularized regression)
- Ridge Regression (regularized regression)
- Decision Tree regression.
- Random forest regression
- Gradient Boosting regression.

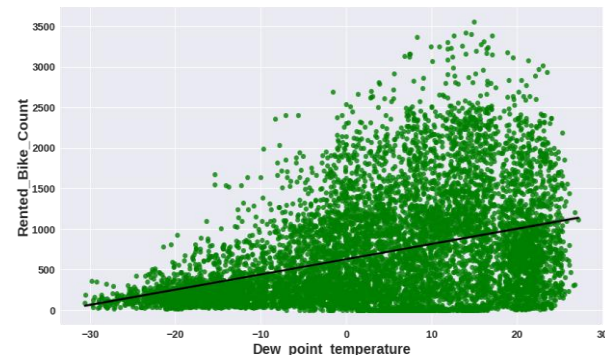
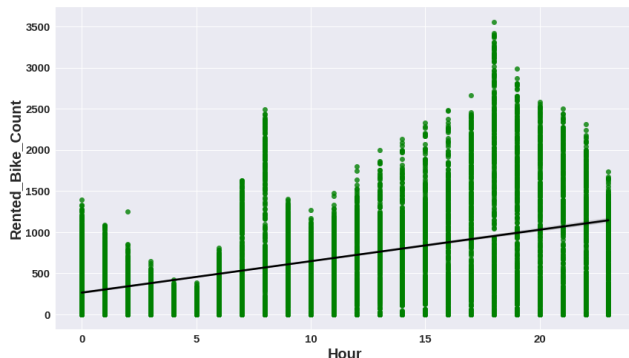
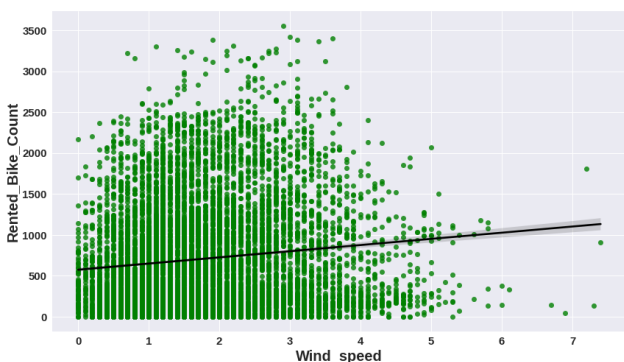
Assumptions of regression line:

1. The relation between the dependent and independent variables should be almost linear.
2. Mean of residuals should be zero or close to 0 as much as possible. It is done to check whether our line is actually the line of “best fit”.
3. There should be homoscedasticity or equal variance in a regression model. This assumption means that the variance around the regression line is the same for all values of the predictor variable (X).
4. There should not be multicollinearity in regression model. Multicollinearity generally occurs when there are high correlations between two or more independent variables.

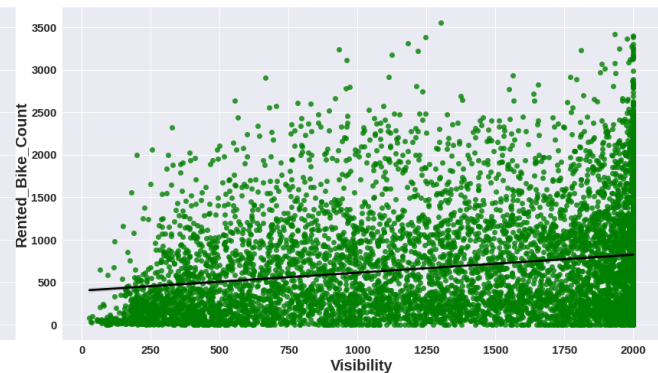
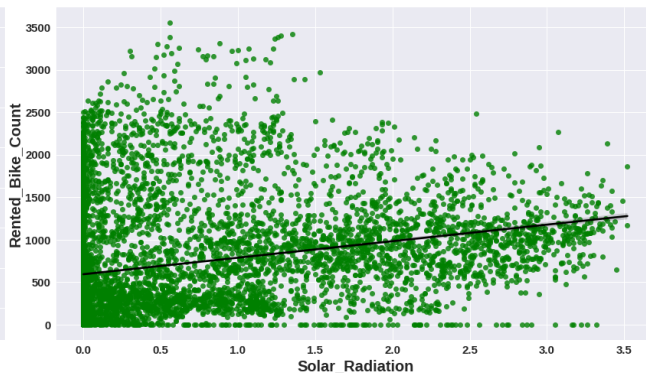
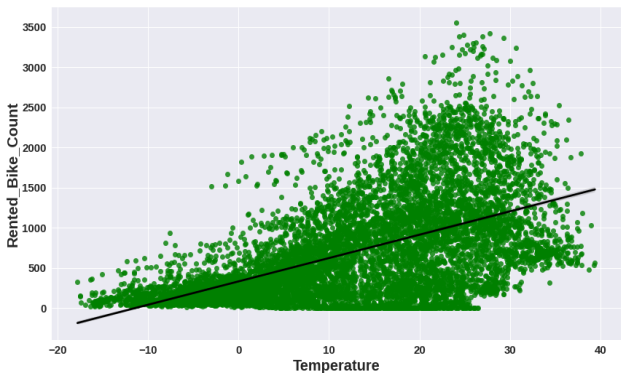
- Before and after applying these models we checked our regression assumptions by distribution of residuals, scatter plot of actual and predicted values, removing multi-collinearity among independent variables.



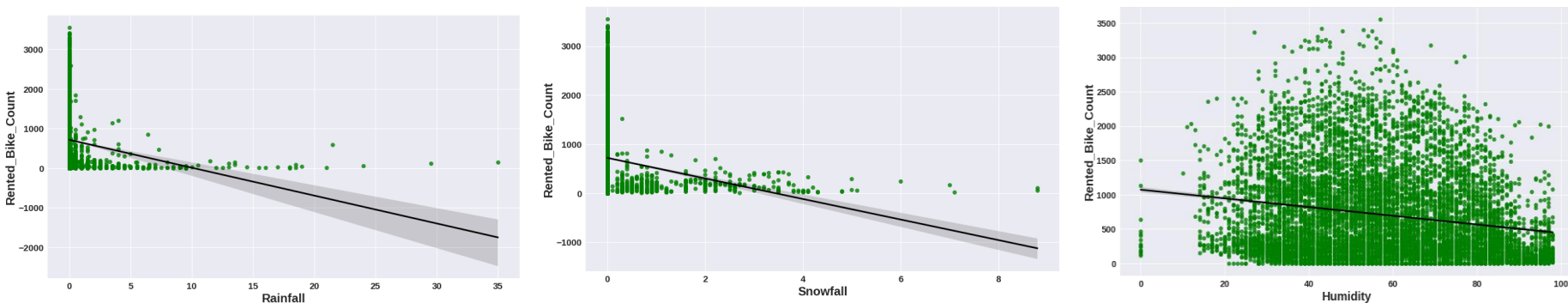
Model Selection and Evaluation :



From the above regression plot of all numerical features we see that the columns 'Temperature', 'Wind_speed', 'Visibility', 'Dew_point_temperature', 'Solar_Radiation' are positively relation to the target variable, which means the rented bike count increases with increase of these features.



❖ Model Selection and Evaluation :



- **'Rainfall', 'Snowfall', 'Humidity' these features are negatively related with the target variable which means the rented bike count decreases when these features increase.**

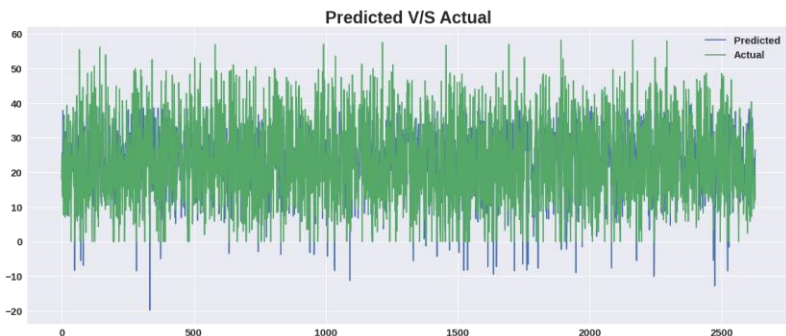
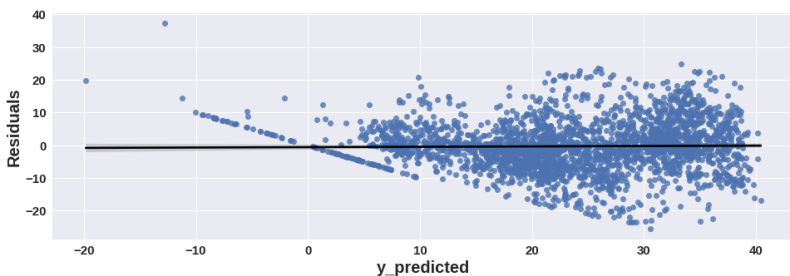
❖ Model Selection and Evaluation :

Linear regression, Lasso and Ridge Regression:

➤ Linear Regression

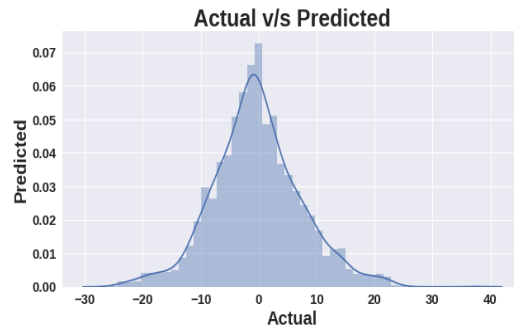
Scores on Train set

The Mean Absolute Error (MAE) is 5.855397241788345.
 The Mean Squared Error(MSE) is 60.29949292444555.
 The Root Mean Squared Error(RMSE) is 7.765274813195316.
 The R2 Score is 0.6123528085603556.

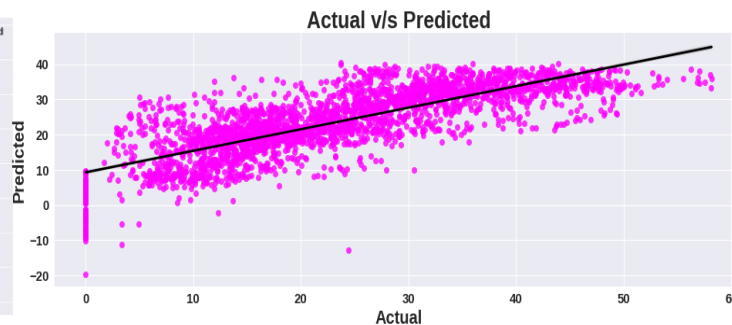


Scores on Test set

The Mean Absolute Error (MAE) is 5.834169822951748.
 The Mean Squared Error(MSE) is 58.624247223024895.
 The Root Mean Squared Error(RMSE) is 7.656647257319936.
 The R2 Score is 0.618326967365199.



Mean of residuals should be zero or close to 0 as much as possible. It is done to check whether our line is actually the line of “best fit”



❖ Model Selection and Evaluation :

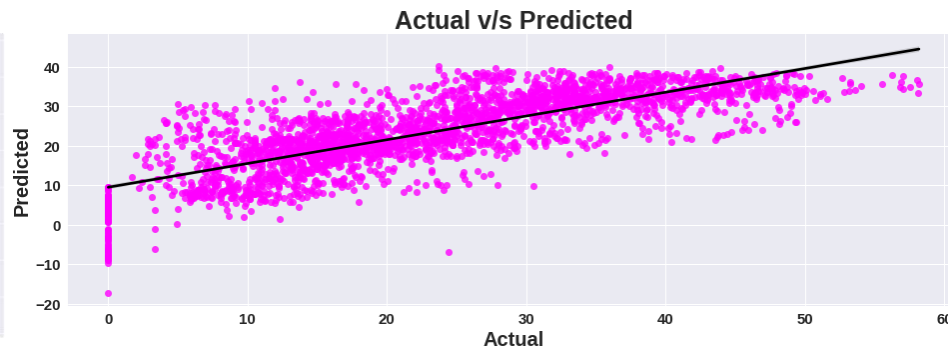
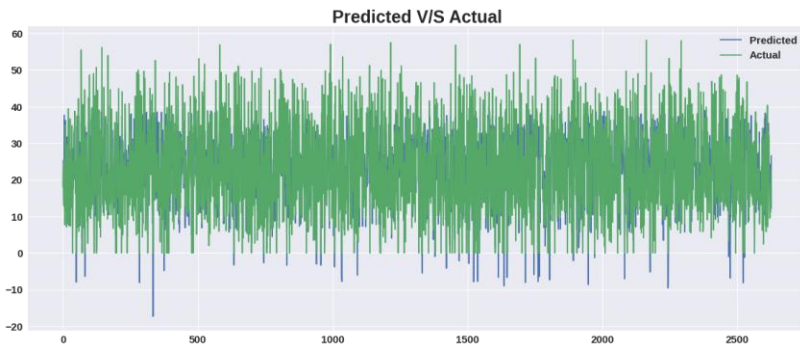
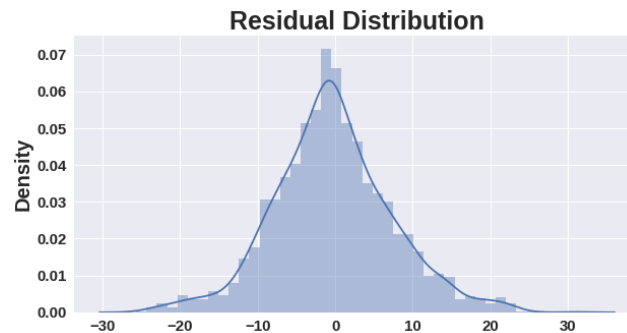
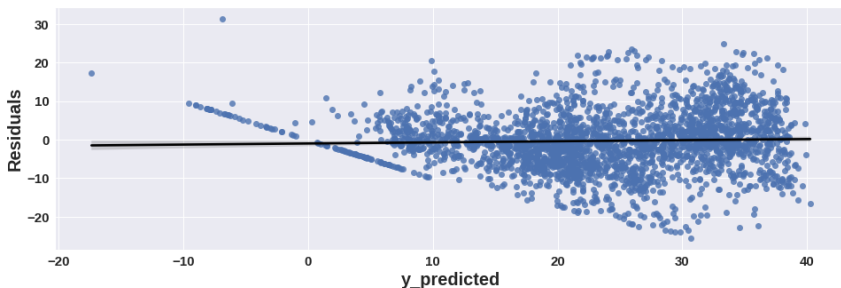
➤ Lasso (Hyper-parameter tuned- $\alpha=0.01$)

Scores on Train set

The Mean Absolute Error (MAE) is 5.869103531726283.
 The Mean Squared Error(MSE) is 60.46402436494349.
 The Root Mean Squared Error(RMSE) is 7.775861647749624.
 The R2 Score is 0.6112950857219155.

Scores on Test set

The Mean Absolute Error (MAE) is 5.850566426263689.
 The Mean Squared Error(MSE) is 58.792684087499225.
 The Root Mean Squared Error(RMSE) is 7.667638755673042.
 The R2 Score is 0.61723035952942.



❖ Model Selection and Evaluation :

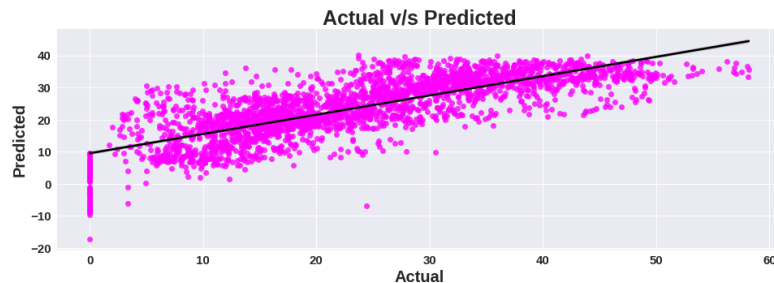
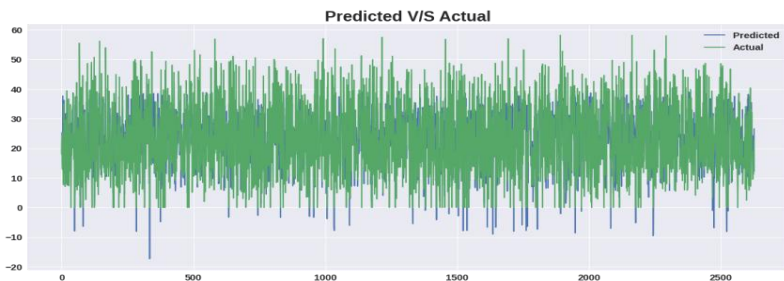
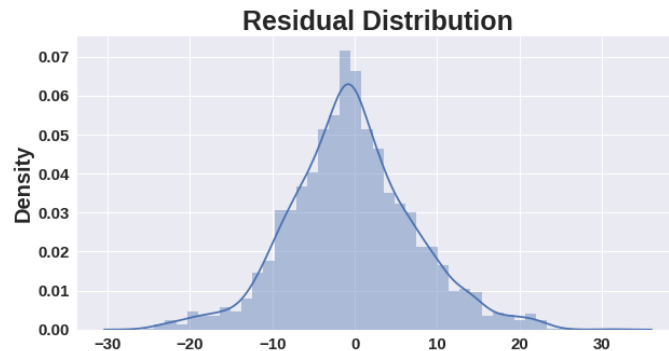
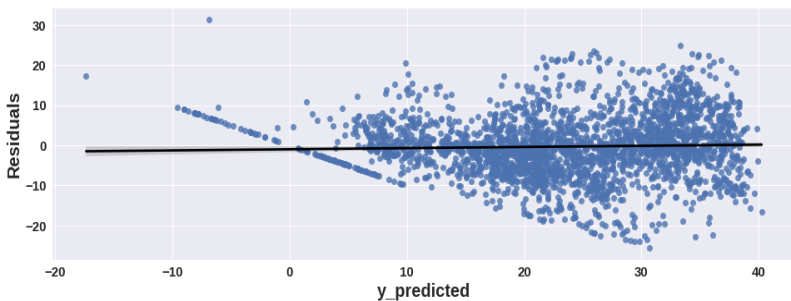
➤ Ridge (Hyper-parameter tuned- $\alpha=0.1$)

Scores on Train set

The Mean Absolute Error (MAE) is 5.869103531726283.
 The Mean Squared Error(MSE) is 60.46402436494349.
 The Root Mean Squared Error(RMSE) is 7.775861647749624.
 The R2 Score is 0.6112950857219155.

Scores on Test set

The Mean Absolute Error (MAE) is 5.850566426263689.
 The Mean Squared Error(MSE) is 58.792684087499225.
 The Root Mean Squared Error(RMSE) is 7.667638755673042.
 The R2 Score is 0.61723035952942.



❖ Model Selection and Evaluation :

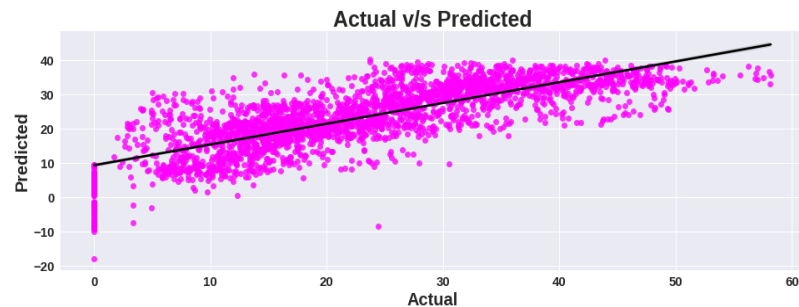
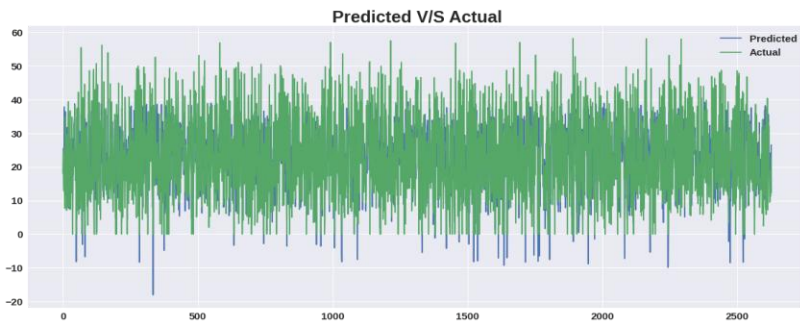
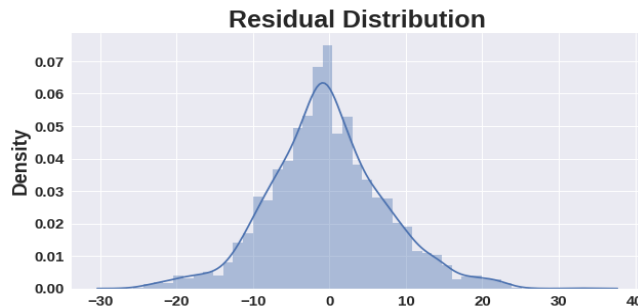
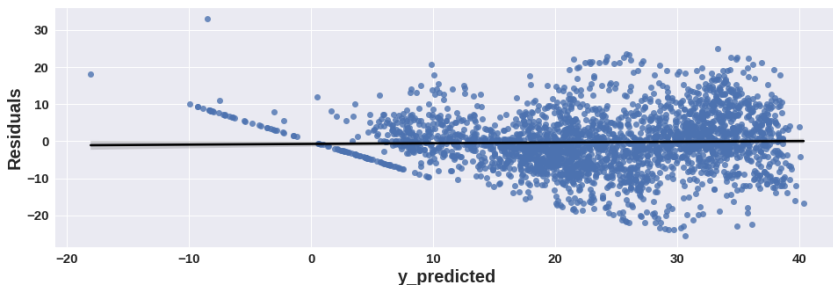
➤ Elastic Net (Hyper-parameter tuned- $\alpha=0.001, l1_ratio=0.5$)

Scores on Train set

The Mean Absolute Error (MAE) is 5.8627247571297865.
 The Mean Squared Error(MSE) is 60.35398893457537.
 The Root Mean Squared Error(RMSE) is 7.76878297641113.
 The R2 Score is 0.6120024702034815.

Scores on Test set

The Mean Absolute Error (MAE) is 5.841404955954685.
 The Mean Squared Error(MSE) is 58.702377424389866.
 The Root Mean Squared Error(RMSE) is 7.661747674283581.
 The R2 Score is 0.6178183008610153.



❖ Model Selection and Evaluation :

➤ Decision Tree regression(Hyper-parameter tuned- max_depth=9,max_features='auto')

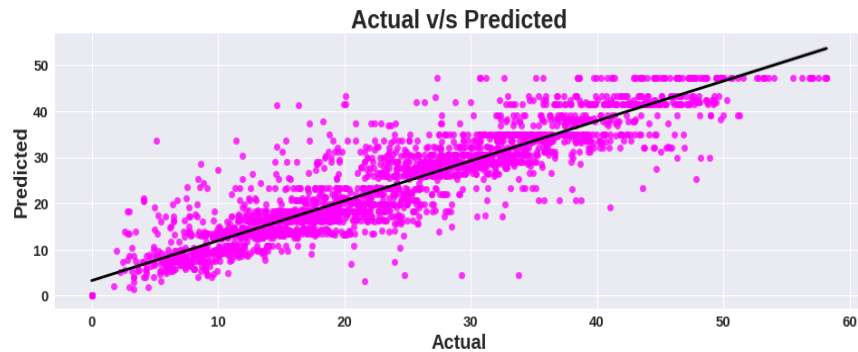
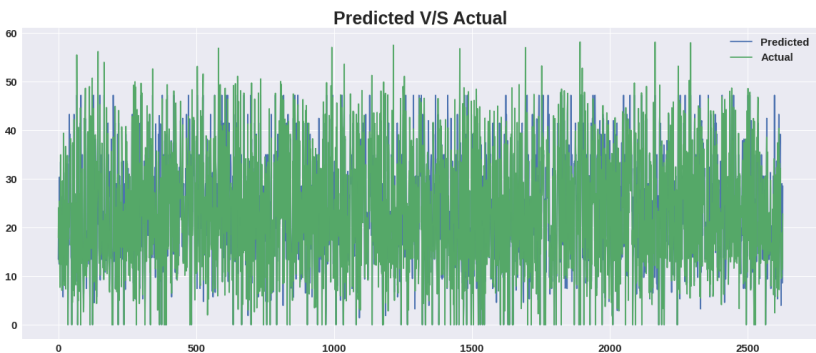
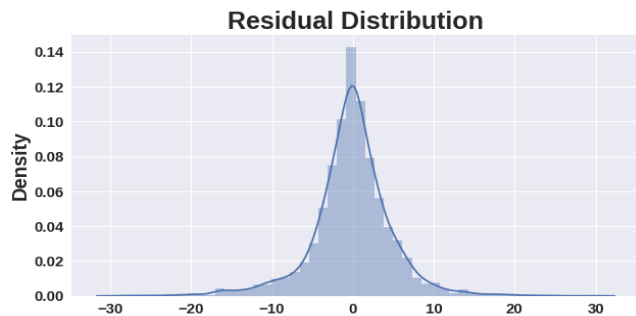
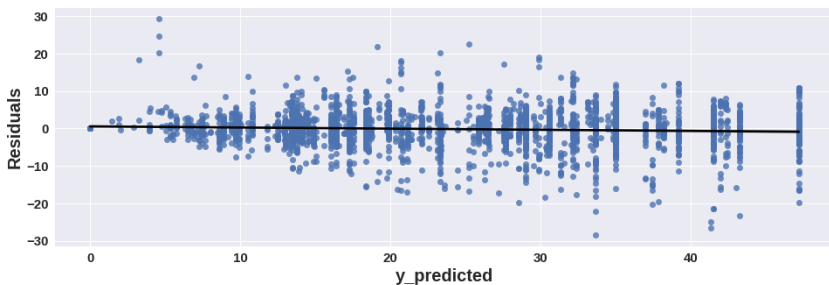
The number of features to consider when looking for the best split

Scores on Train set

The Mean Absolute Error (MAE) is 2.8855165215690715.
The Mean Squared Error(MSE) is 18.444625087726916.
The Root Mean Squared Error(RMSE) is 4.294720606480347.
The R2 Score is 0.8814250872495163.

Scores on Test set

The Mean Absolute Error (MAE) is 3.3992026410244094.
The Mean Squared Error(MSE) is 24.910895604820194.
The Root Mean Squared Error(RMSE) is 4.9910816067081285.
The R2 Score is 0.8378176689421706.



❖ Model Selection and Evaluation :

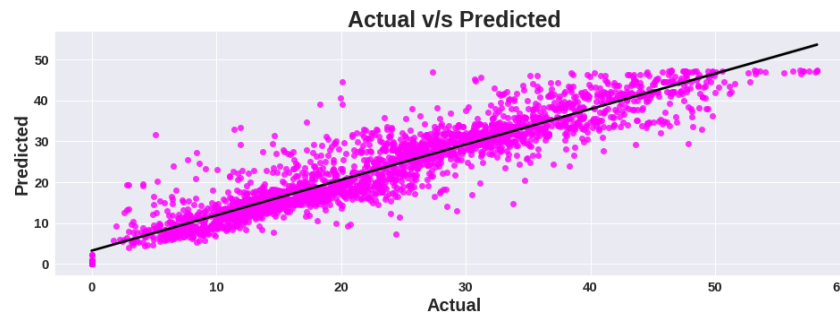
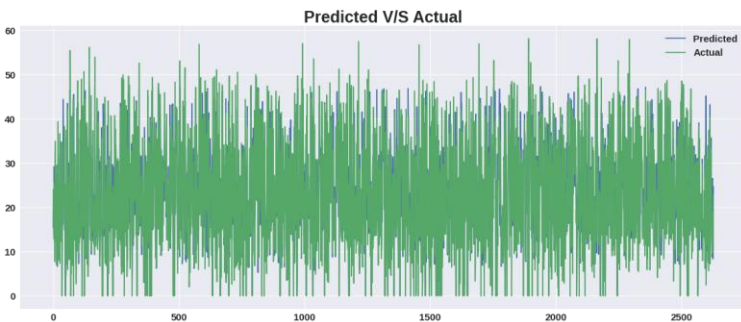
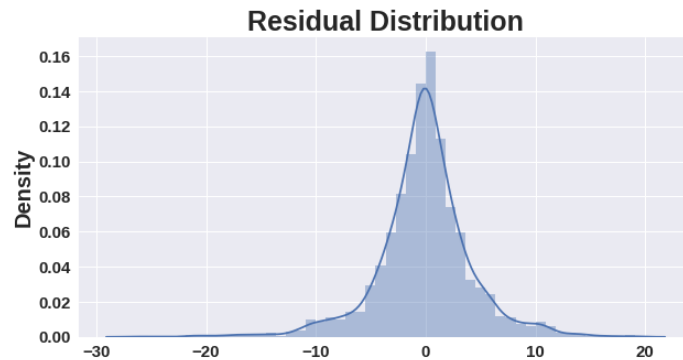
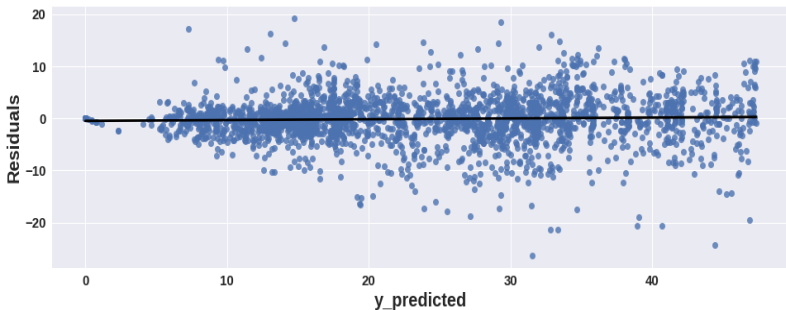
➤ Random forest regression(Hyper-parameter tuned- 'max_depth': 9, 'n_estimators': 100')

Scores on Train set

The Mean Absolute Error (MAE) is 2.6247545856850936.
 The Mean Squared Error(MSE) is 14.905429807049964.
 The Root Mean Squared Error(RMSE) is 3.8607550825000496.
 The R2 Score is 0.904177502634334.

Scores on Test set

The Mean Absolute Error (MAE) is 2.947737482949683.
 The Mean Squared Error(MSE) is 18.68756387544933.
 The Root Mean Squared Error(RMSE) is 4.32291150446656.
 The R2 Score is 0.8783346564815596.



❖ Model Selection and Evaluation :

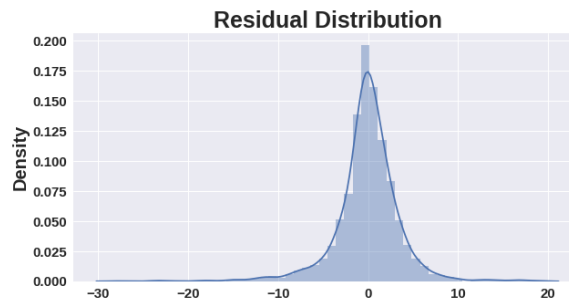
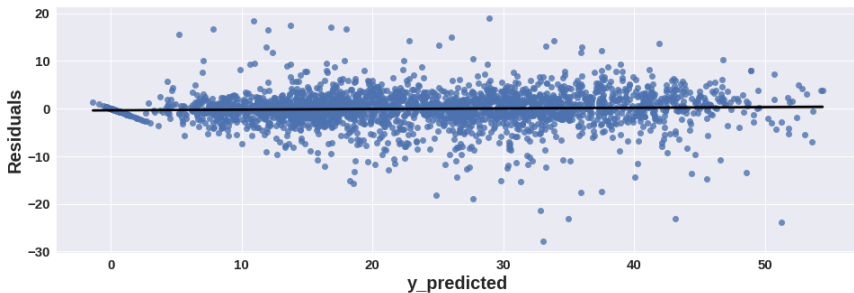
➤ Gradient boosting regression(Hyper-parameter tuned- 'learning_rate': 0.04, 'max_depth': 8, 'n_estimators': 150, 'subsample': 0.9)

Scores on Train set

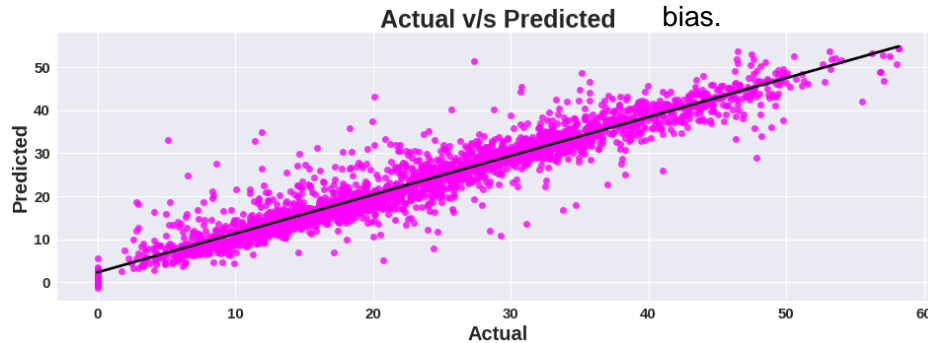
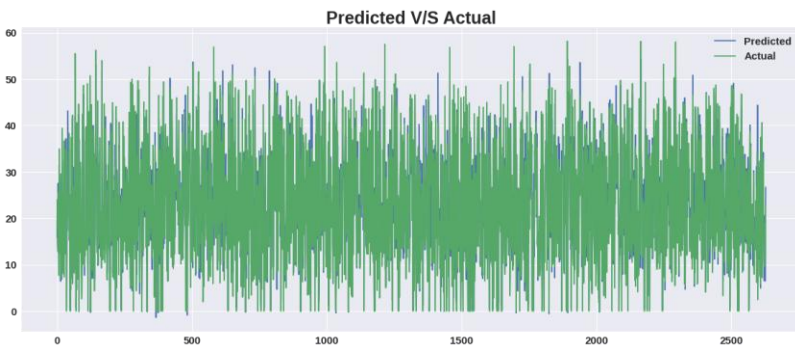
The Mean Absolute Error (MAE) is 1.5177256084968402.
 The Mean Squared Error(MSE) is 4.794733133724516.
 The Root Mean Squared Error(RMSE) is 2.1896879078363005.
 The R2 Score is 0.9691761117241932.

Scores on Test set

The Mean Absolute Error (MAE) is 2.3713827898940885.
 The Mean Squared Error(MSE) is 13.217709880555015.
 The Root Mean Squared Error(RMSE) is 3.635616850075791.
 The R2 Score is 0.9139461288874848.

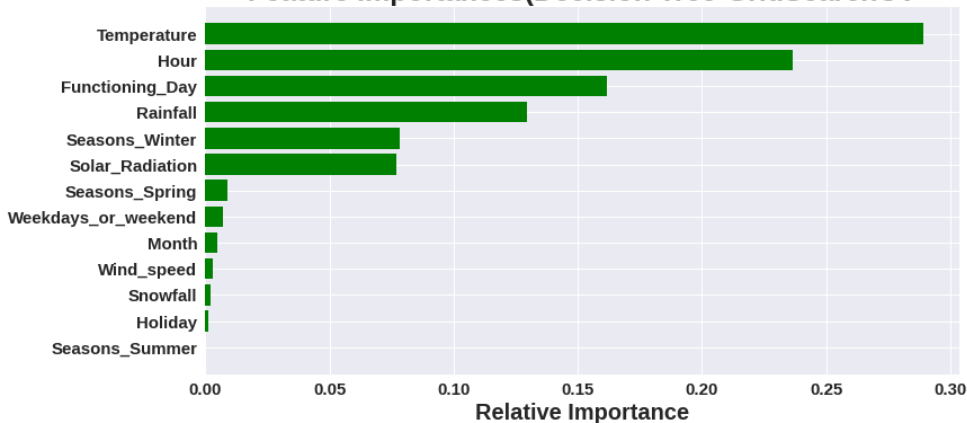


➤ Learning rate shrinks the contribution of each tree by learning_rate. There is a trade-off between learning_rate and n_estimators.
 ➤ Choosing subsample < 1.0 leads to a reduction of variance and an increase in bias.

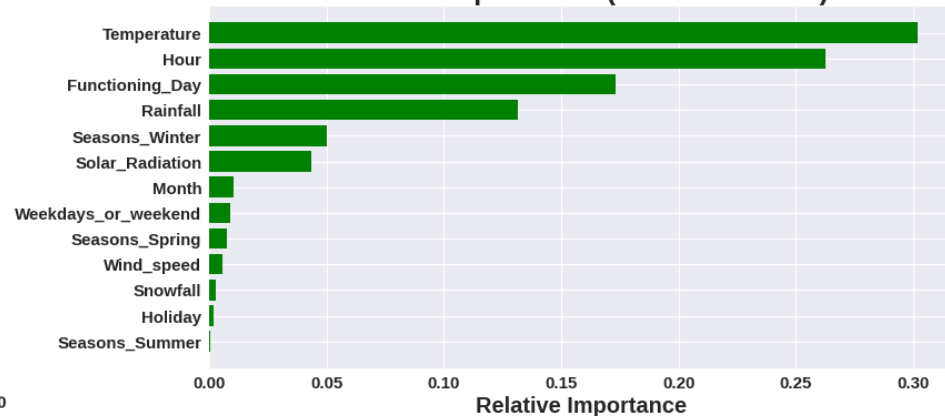


❖ Feature importance's :

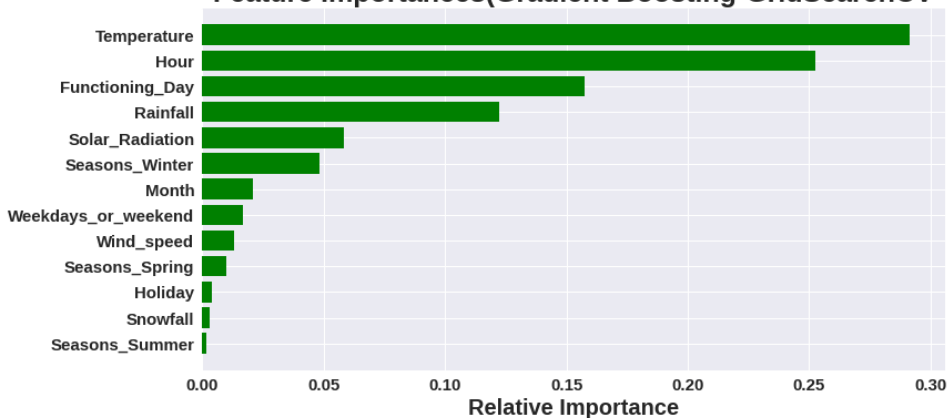
Feature Importances(Decision Tree-GridSearchCV)



Feature Importances(Random Forest)



Feature Importances(Gradient Boosting-GridSearchCV)



From all 3 models we can say that temperature, hour, functioning day are the top three important features.

❖ Conclusion:

	Model	MAE	MSE	RMSE	R2_score
Training set	0 Linear Regression	5.8555	60.2995	7.7653	0.6124
	1 Lasoo	5.8691	60.4640	7.7759	0.6113
	2 Ridge GridSearchCV	5.8691	60.4640	7.7759	0.6113
	3 ElasticNet(GridSearchCV-Tunned)	5.8932	60.9027	7.8040	0.6085
	4 Decision Tree Regressor-GridSearchCV	2.8855	18.4446	4.2947	0.8814
	5 Random Forest	0.9448	2.1723	1.4739	0.9860
	6 Random Forest-GridSearchCv	2.6205	14.7855	3.8452	0.9049
	7 Gardient boosting Regression	3.1772	20.5277	4.5308	0.8680
	8 Gradient Boosting Regression(GridSearchCV)	1.5432	5.0540	2.2481	0.9675
Test set	0 Linear Regression	5.8342	58.6242	7.6566	0.6183
	1 Lasso	5.8506	58.7927	7.6676	0.6172
	2 Ridge(GridsearchCv Tunned)	5.8506	58.7927	7.6676	0.6172
	3 ElasticNet(GridSearchCV-Tunned)	5.8711	59.2909	7.7001	0.6140
	4 Decision Tree Regressor(GridsearchCV)	3.3973	24.8655	4.9865	0.8381
	5 Radom forest	2.4778	14.2254	3.7717	0.9074
	6 Random Forest-GridSearchCv	2.9422	18.5797	4.3104	0.8790
	7 Gradient Boosting Regression	3.2843	21.6821	4.6564	0.8588
	8 Gradient Boosting Regression(GridSearchCV)	2.3978	13.4016	3.6608	0.9127

As we have calculated MAE,MSE,RMSE and R2 score for each model. Based on r2 score will decide our model performance.

Our assumption: if the difference of R2 score between Train data and Test is more than 5 % we will consider it as over fitting.

Linear, Lasso, Ridge and Elastic Net:

Linear, Lasso, Ridge and Elastic regression models have almost similar R2 scores(61%) on both training and test data.(Even after using GridserachCV we have got similar results as of base models).

Decision Tree Regression:

On Decision tree regressor model, without hyper-parameter tuning, we got r2 score as 100% on training data and on test data it was very less. Thus our model memorized the data. So it was a over fitted model.

After hyper-parameter tuning we got r2 score as 88% on training data and 83% on test data which is quite good for us.

Random Forest:

On Random Forest regressor model, without hyper-parameter tuning we got r2 score as 98% on training data and 90% on test data. Thus our model memorized the data. So it was a over fitted model, as per our assumption

After hyper-parameter tuning we got r2 score as 90% on training data and 87% on test data which is very good for us.

Gradient Boosting Regression(Gradient Boosting Machine):

On Random Forest regressor model, without hyper-parameter tuning we got r2 score as 86% on training data and 85% on test data. Our model performed well without hyper-parameter tuning. After hyper-parameter tuning we got r2 score as 96% on training data and 91% on test data, thus we improved the model performance by hyper-parameter tuning.

❖ Conclusion:

Thus Gradient Boosting Regression(GridSearchCV) and Random forest(GridSearchCv) gives good r^2 scores. We can deploy this models.

Signing off...

THANK YOU