

Capstone Project -2

Bike Sharing Demand Prediction

(Supervised Machine Learning regression)

TEAM MEMBERS

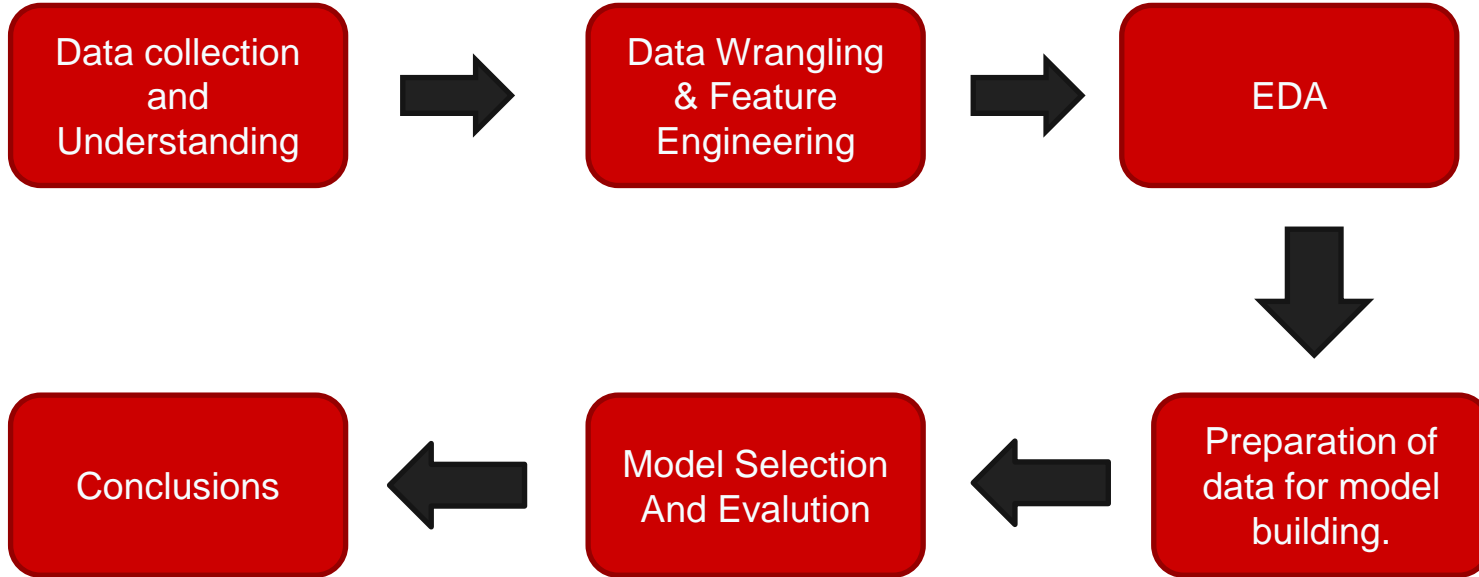
Vimal Kumar

Vishal Kumar Yadav

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.



- 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)

As we know we had 8760 observations and 14 features.

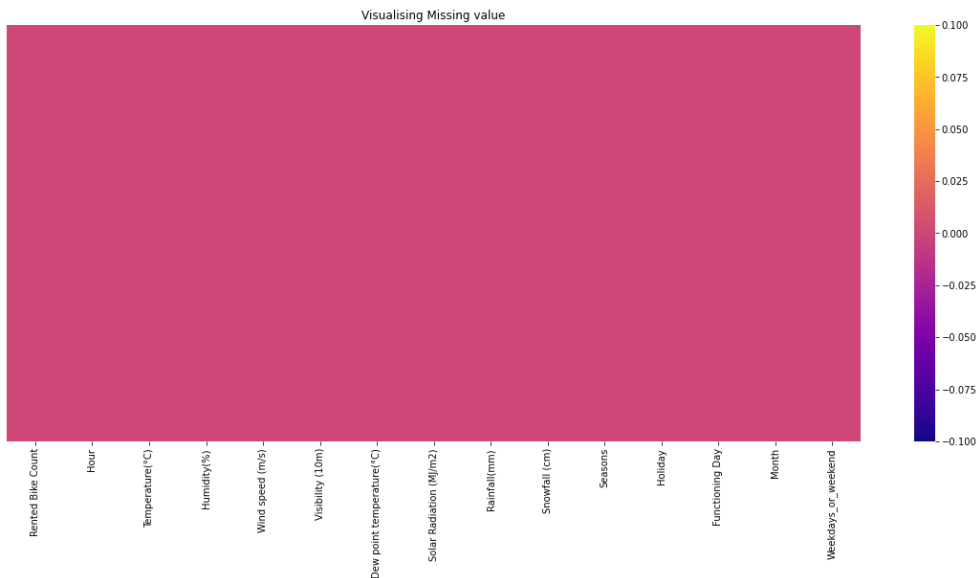
➤ **Categorical Features:** Seasons, Holiday and Functioning day.

➤ **Numerical Columns:**

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

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.



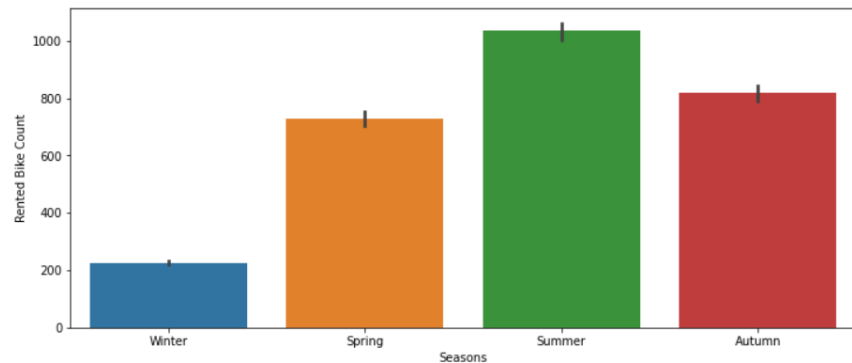
```
[ ] #change the date datatype to extract "Month","day","year".
bike_df['Date']=bike_df['Date'].astype('datetime64[ns]')
```

```
[ ] #checking duplicates value
duplicates=bike_df.duplicated().sum()
print(f"we have {duplicates} rows in our bike data.")
```

we have 0 rows in our bike data.

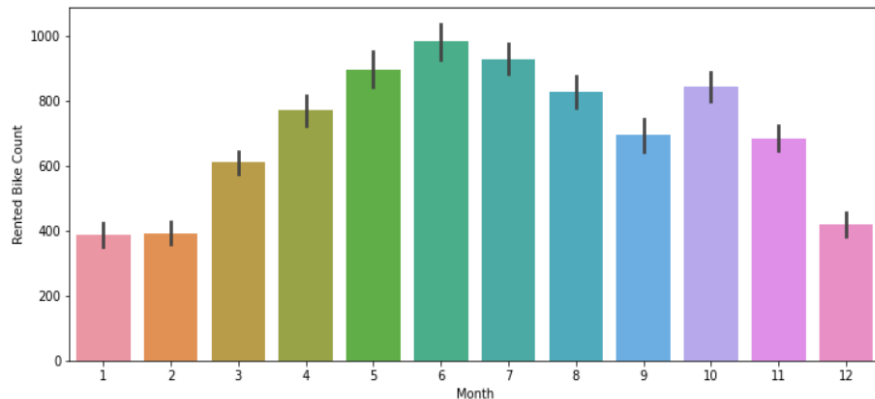
```
[ ] #creating new columns"Month","year","Day".
bike_df['Month']=bike_df['Date'].dt.month

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

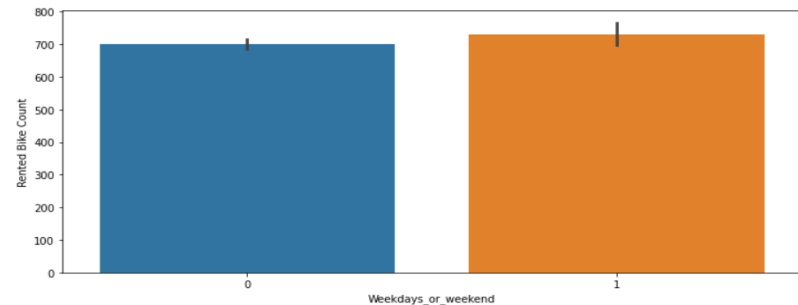
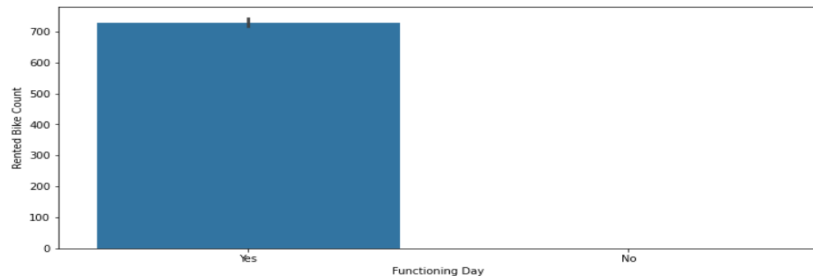
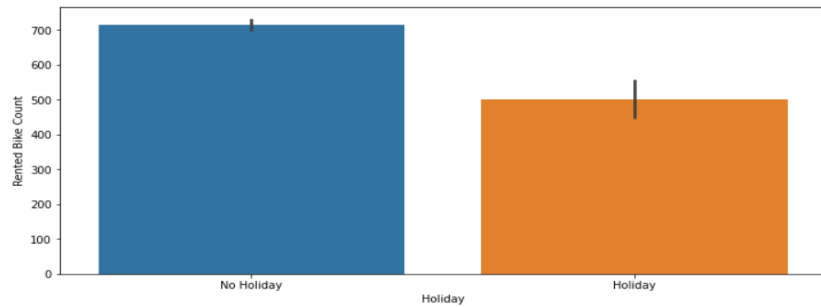


- Relation of rented bike count with categorical features:

Summer season had 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.



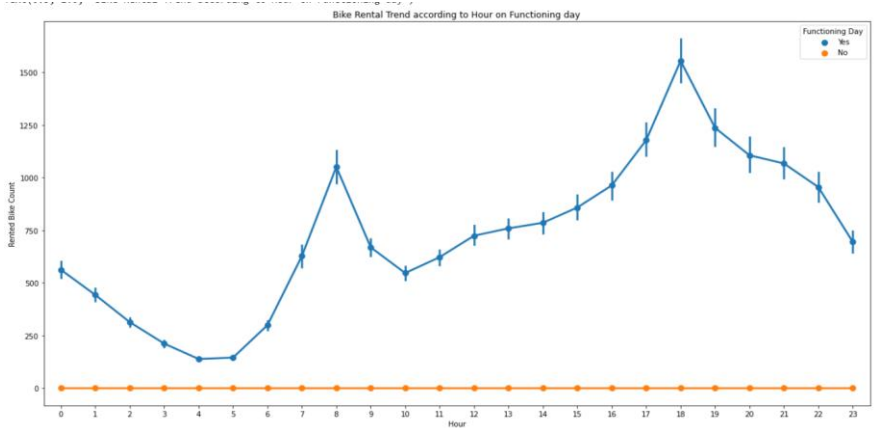
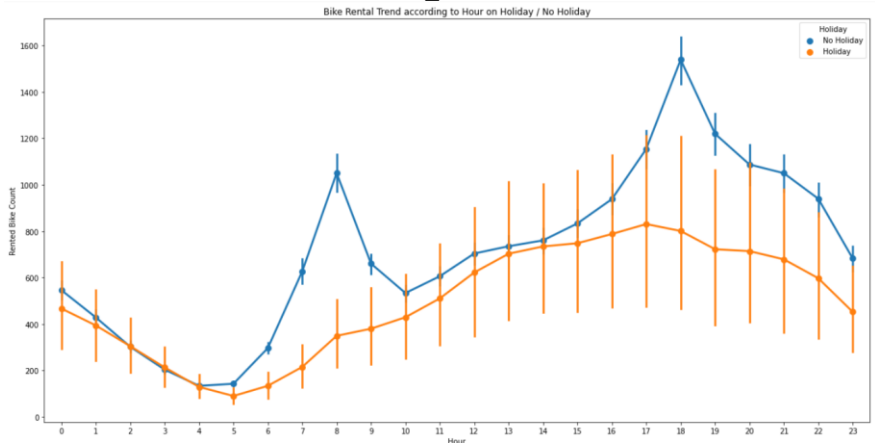
- From March Bike Rent Count started increasing and it was highest in June



Conclusion:

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

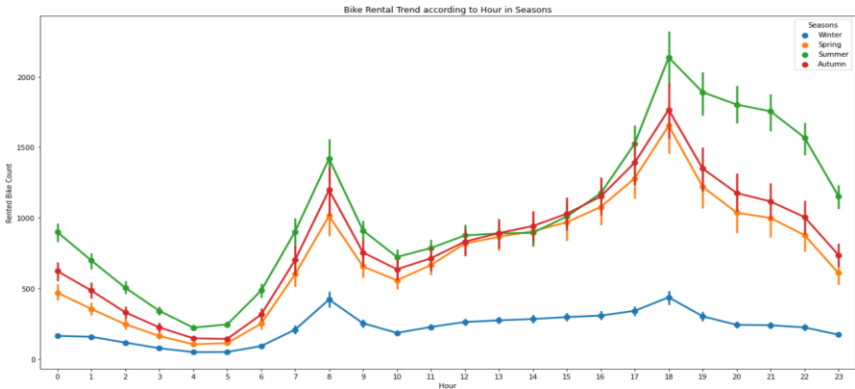
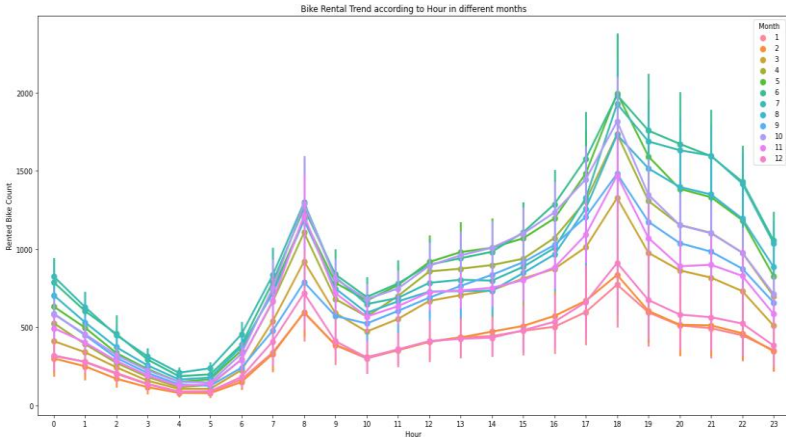
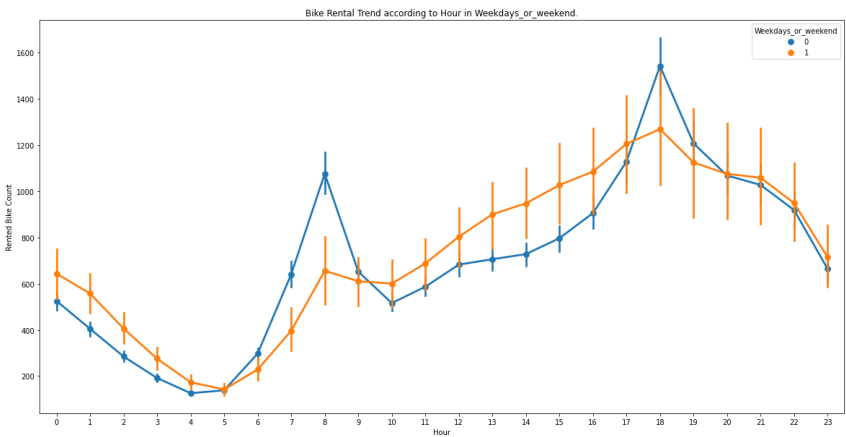
Bike Rent Trend According to Hour in different scenarios.



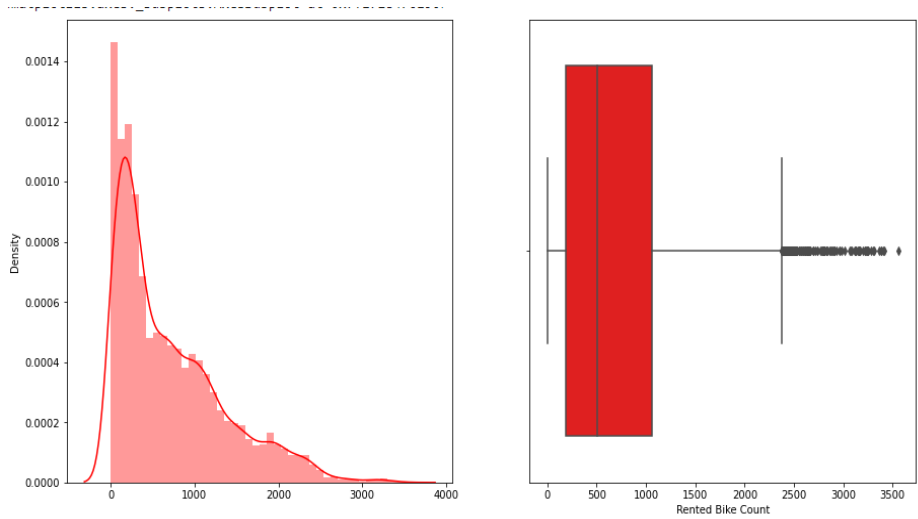
Observations:

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.

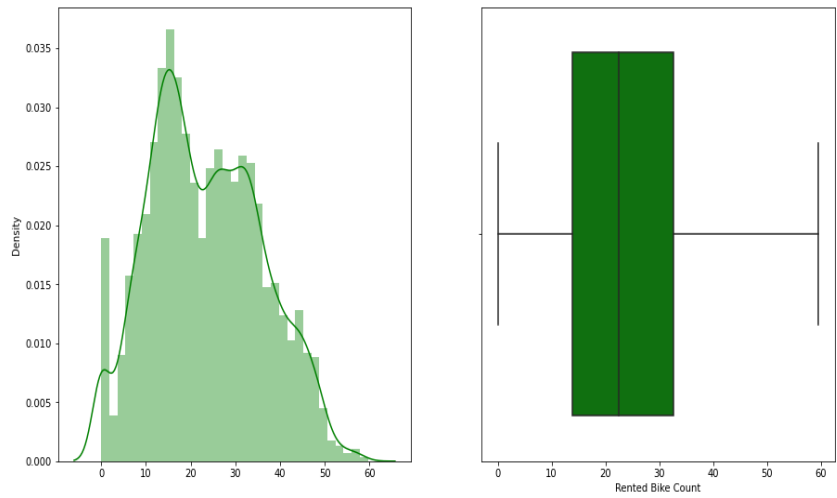
Bike Rent Trend according to hour in different scenarios.



Distribution of target variable- Bike Rent Count



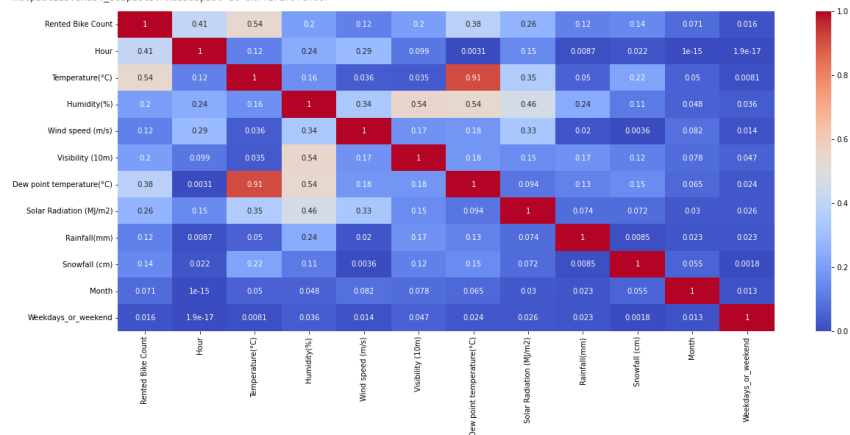
Distribution is rightly skewed and some outliers are observed.



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

Preparation of data for model building :

`<matplotlib.axes._subplots.AxesSubplot at 0x7f172797cfd0>`



➤ 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 :

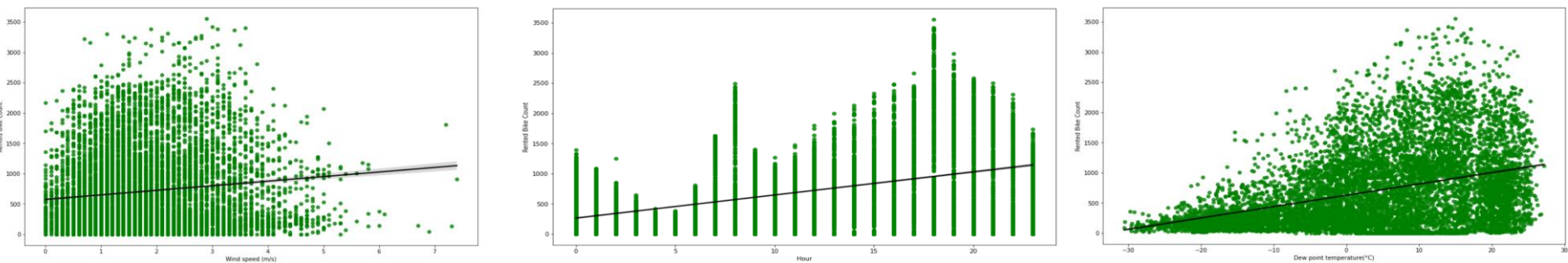
As this is the regression problem we are trying to predict continuous value. For this we used following regression models.

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

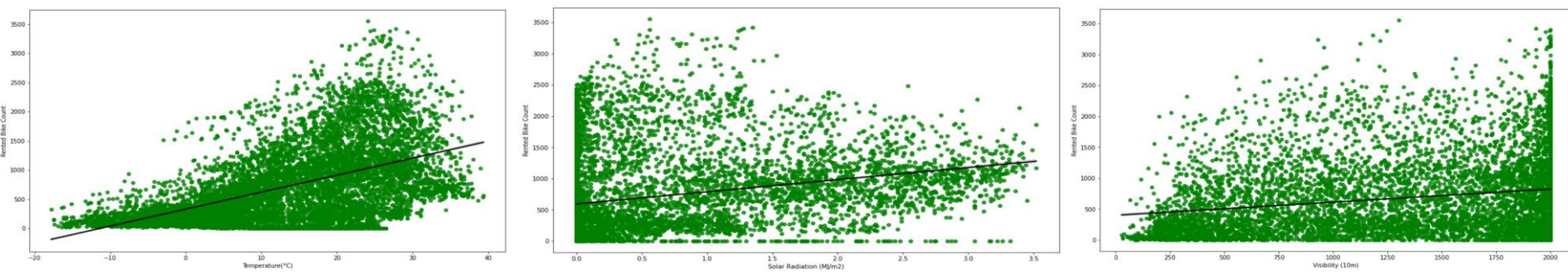
Assumptions of regression line:

- 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.

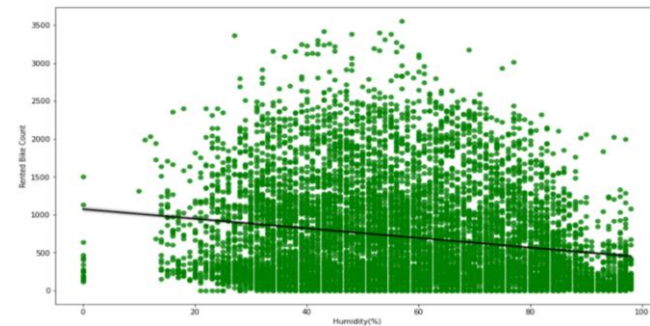
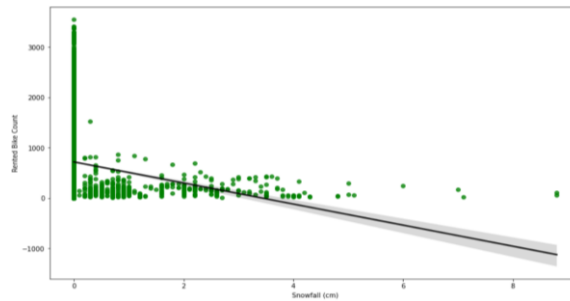
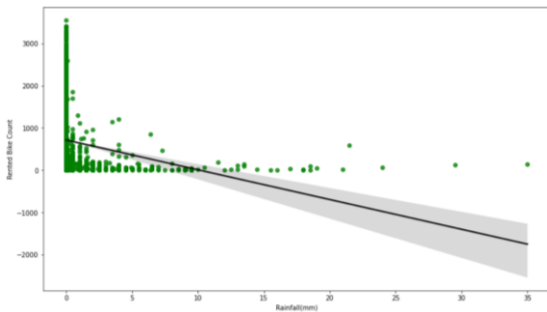
- Before and after applying these models we checked our regression assumptions by distribution of residuals, scatter plot of actual and predicted values, removing multi-collinearity among independent variables.



➤ 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 :



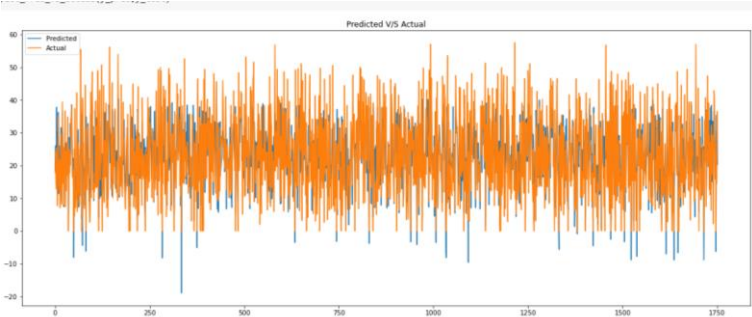
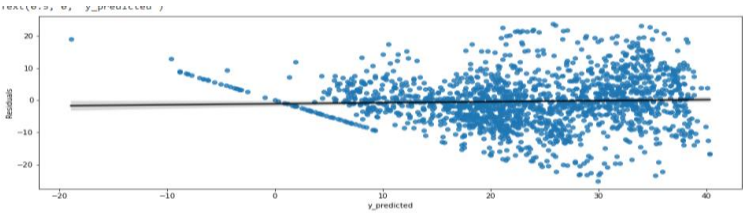
- 'Rainfall', 'Snowfall', 'Humidity' these features are negatively related with the target variable which means the rented bike count decreases when these features increase.

Linear regression, Lasso and Ridge Regression:

➤ Linear Regression

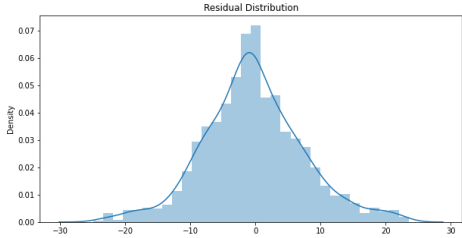
Scores on Train set

The Mean Absolute Error (MAE) is 5.839648887883752.
The Mean Squared Error(MSE) is 59.78437886830986.
The Root Mean Squared Error(RMSE) is 7.732035881209415.
The R2 Score is 0.6126276858762716.

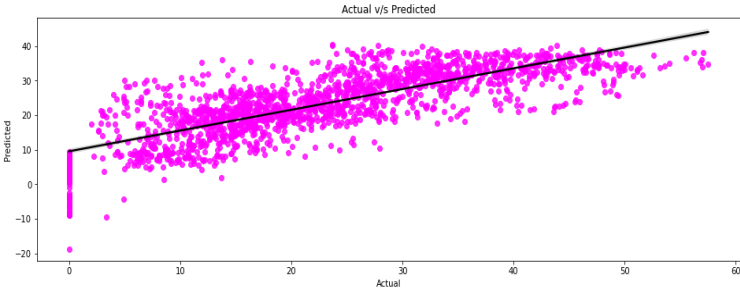


Scores on Test set

The Mean Absolute Error (MAE) is 5.910694345961074.
The Mean Squared Error(MSE) is 60.09070395561458.
The Root Mean Squared Error(RMSE) is 7.7518193964781315.
The R2 Score is 0.6184384273739414.



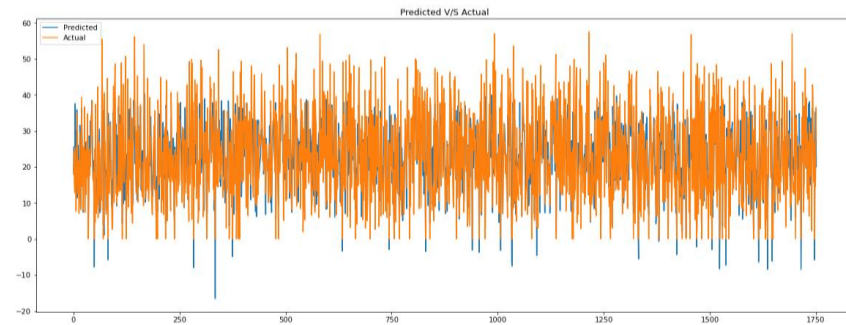
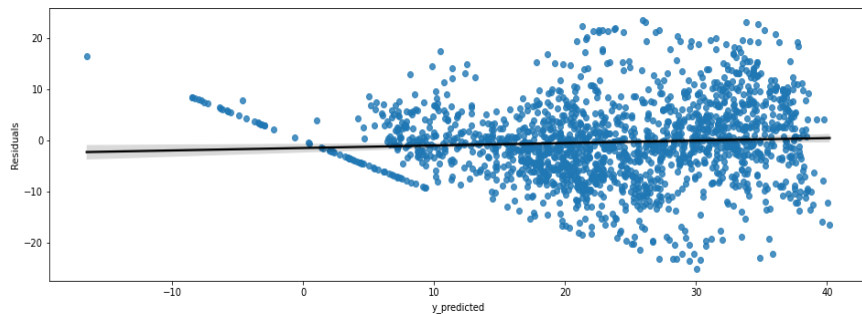
Means of residuals should be zero or close to 0 as much as Possible.It is done to check whether our line is actually the line of "best fit"



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

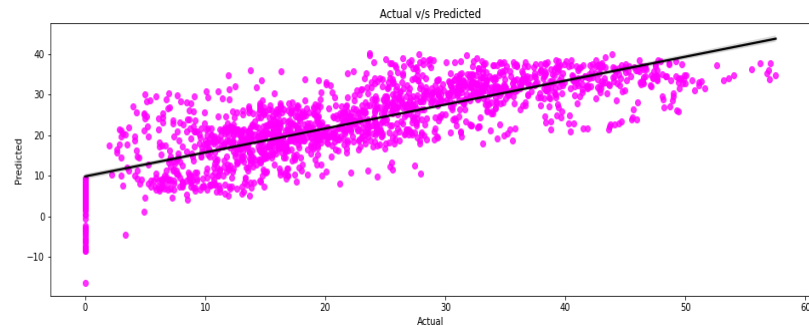
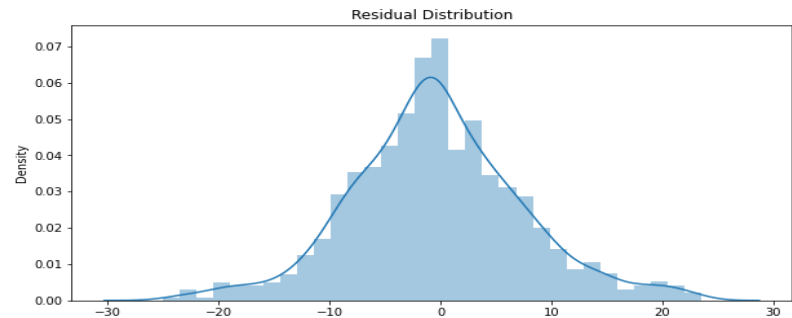
Scores on Train Set

The Mean Absolute Error (MAE) is 5.854107283160169.
The Mean Squared Error (MSE) is 59.944611212141645.
The Root Mean Squared Error (RMSE) is 7.742390536012869.
The R2 Score is 0.61158946192877.



Scores on Test Set

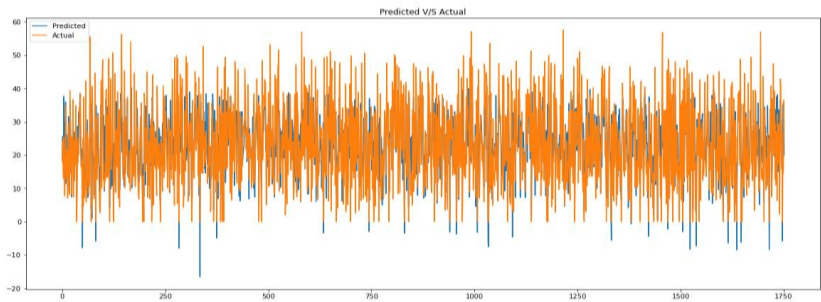
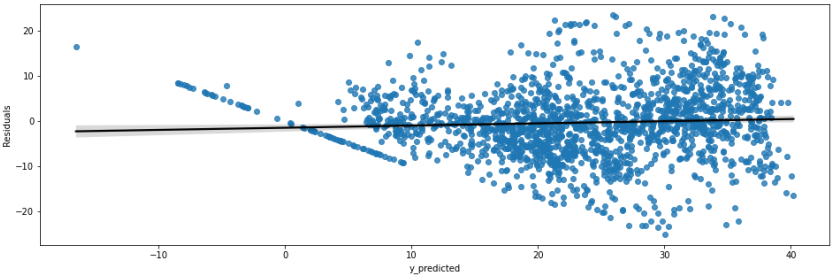
The Mean Absolute Error (MAE) is 5.930553036623993.
The Mean Squared Error (MSE) is 60.47089480694955.
The Root Mean Squared Error (RMSE) is 7.776303415309202.
The R2 Score is 0.6160243065601703.



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

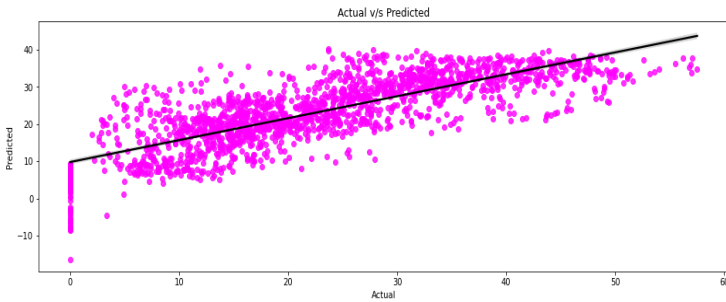
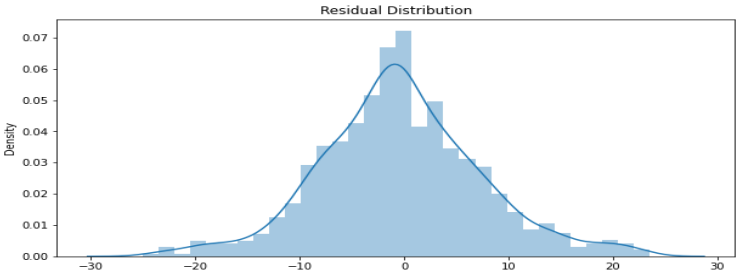
Scores on Train set

The Mean Absolute Error (MAE) is 5.854107283160169.
The Mean Squared Error(MSE) is 59.944611212141645.
The Root Mean Squared Error(RMSE) is 7.742390536012869.
The R2 Score is 0.61158946192877.



Scores on Test set

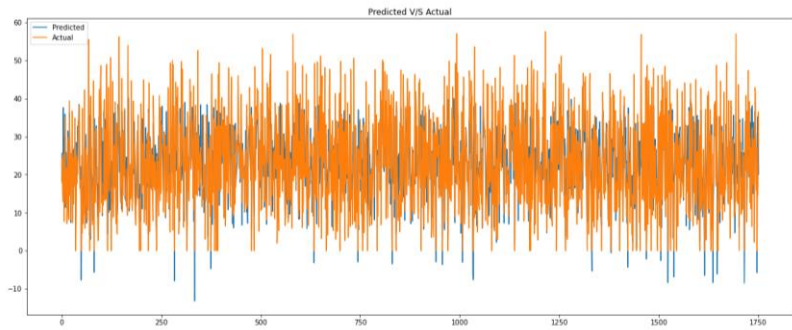
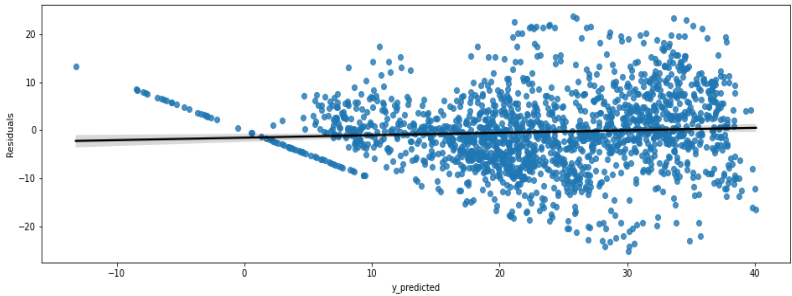
The Mean Absolute Error (MAE) is 5.930553036623993.
The Mean Squared Error(MSE) is 60.47089480694955.
The Root Mean Squared Error(RMSE) is 7.776303415309202.
The R2 Score is 0.6160243065601703.



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

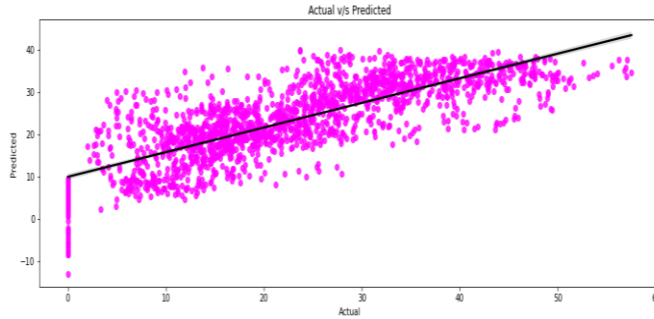
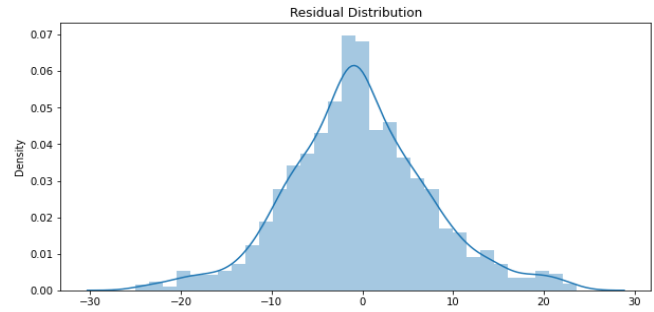
Scores on Train set

The Mean Absolute Error (MAE) is 5.875650253860865.
The Mean Squared Error(MSE) is 60.35480865020068.
The Root Mean Squared Error(RMSE) is 7.768835733248624.
The R2 Score is 0.6089315915312441.



Scores on Test set

The Mean Absolute Error (MAE) is 5.957471344152484.
The Mean Squared Error(MSE) is 61.1969615583184.
The Root Mean Squared Error(RMSE) is 7.82284868563354.
The R2 Score is 0.6114139566516649.

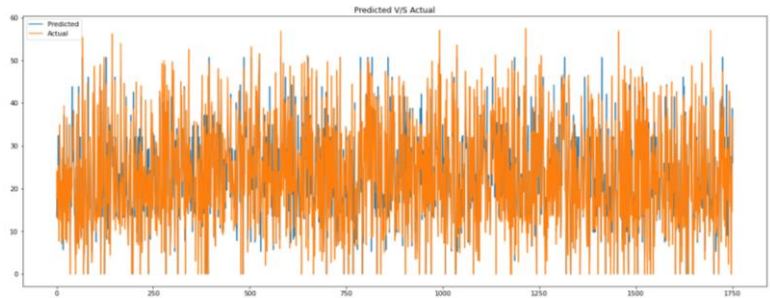
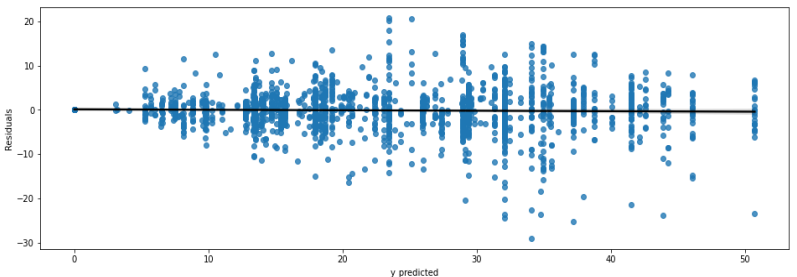


➤ Decision Tree regression(Hyper-parameter tuned- max_depth=9,max_features='auto')

The number of features to consider when looking for the best fit.

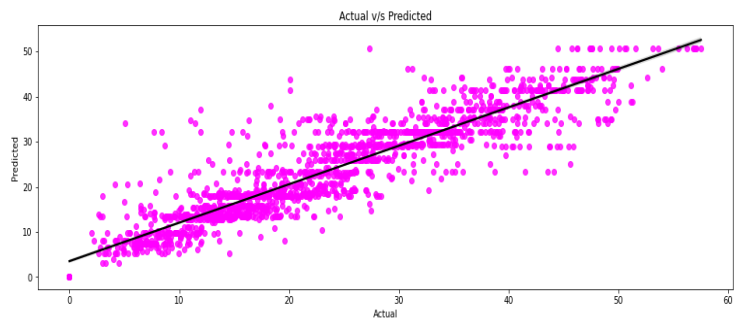
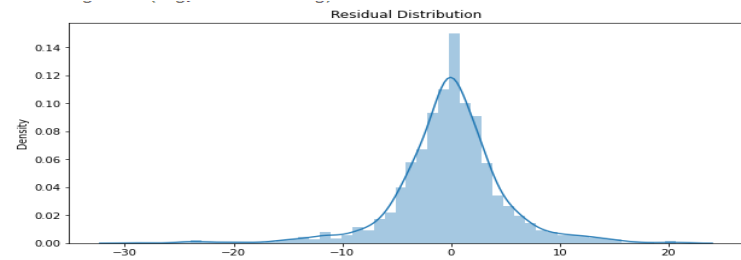
Scores on Train set

The Mean Absolute Error (MAE) is 3.037059512840259.
The Mean Squared Error(MSE) is 20.16583849207843.
The Root Mean Squared Error(RMSE) is 4.4906389848303805.
The R2 Score is 0.8693356413365273.



Scores on Test set

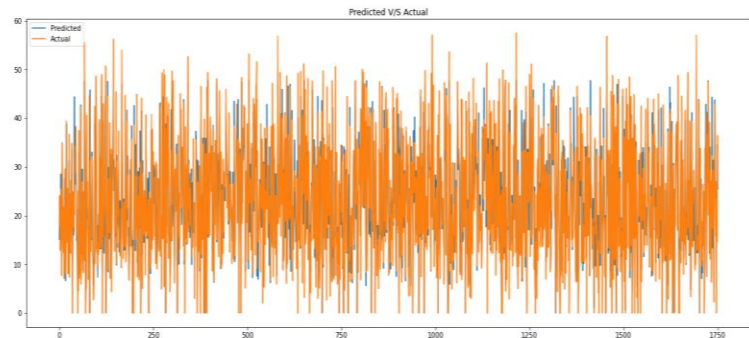
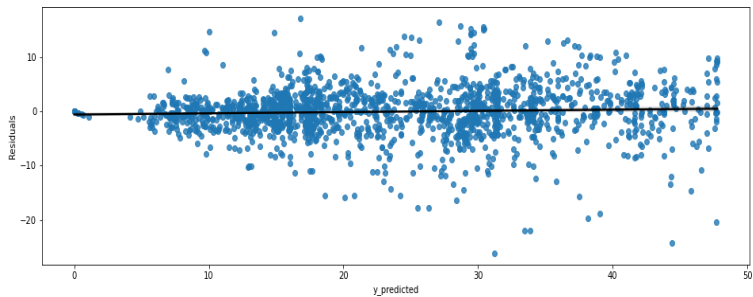
The Mean Absolute Error (MAE) is 3.3574782925331563.
The Mean Squared Error(MSE) is 24.714386718966907.
The Root Mean Squared Error(RMSE) is 4.971356627618553.
The R2 Score is 0.8430695658026737.



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

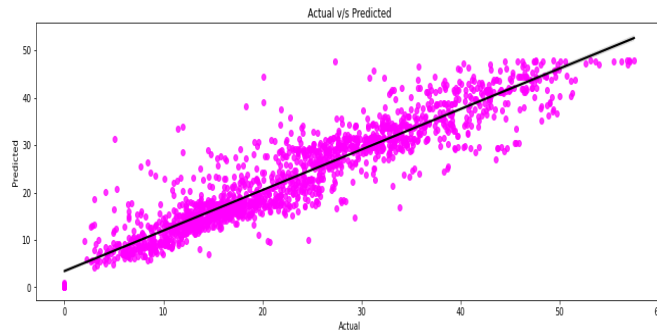
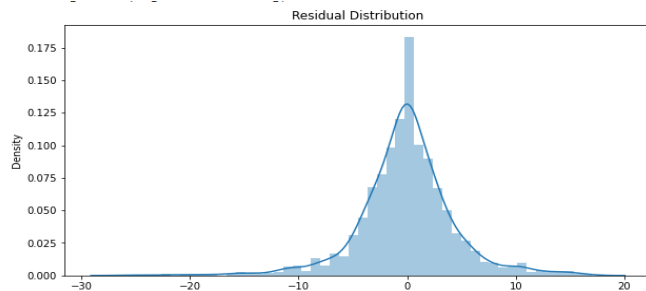
Scores on Train set

The Mean Absolute Error (MAE) is 2.748258687457823.
The Mean Squared Error(MSE) is 16.37518639531495.
The Root Mean Squared Error(RMSE) is 4.046626545076152.
The R2 Score is 0.8938971355354675.



Scores on Test set

The Mean Absolute Error (MAE) is 3.044886711051888.
The Mean Squared Error(MSE) is 19.97092794428468.
The Root Mean Squared Error(RMSE) is 4.46884418317918.
The R2 Score is 0.873189392508174.

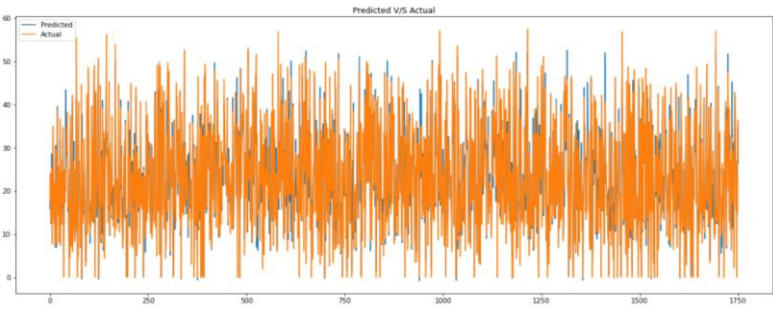
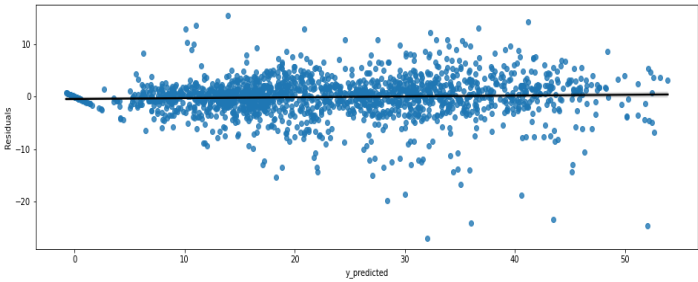




Gradient boosting regression(Hyper-parameter tuned- 'learning_rate': 0.04, 'max_depth': 8, 'n_estimators': 150, 'subsample': 0.9)

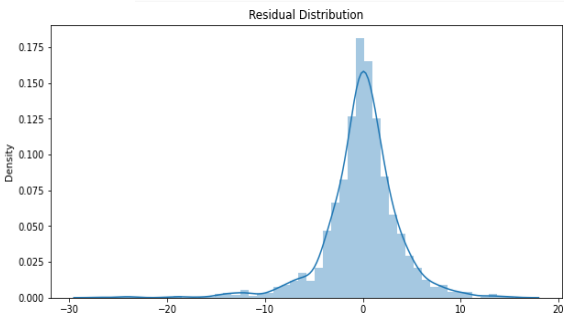
Scores on Train set

The Mean Absolute Error (MAE) is 1.7402446058264194.
The Mean Squared Error(MSE) is 6.5101720716965055.
The Root Mean Squared Error(RMSE) is 2.5515038843193056.
The R2 Score is 0.9578174019953979.

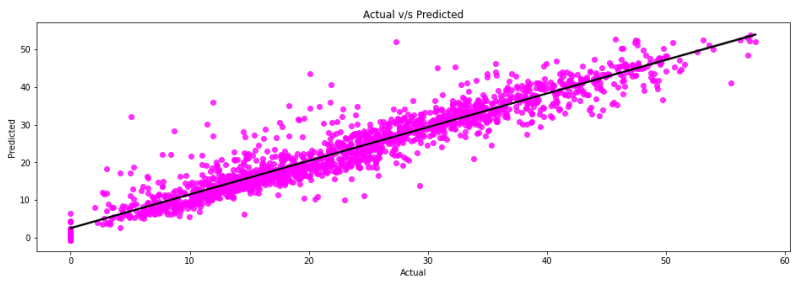


Scores on Test set

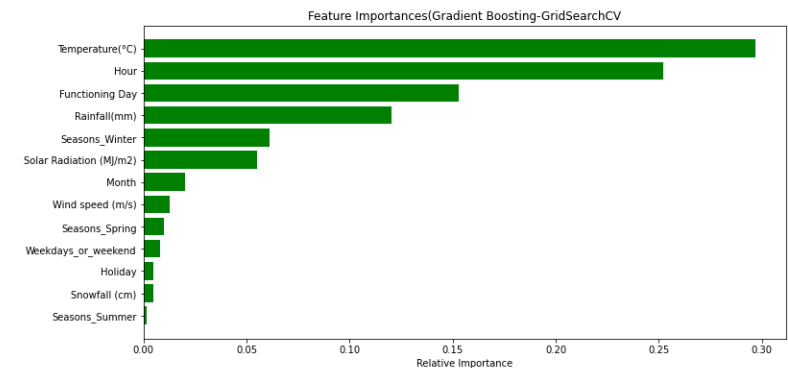
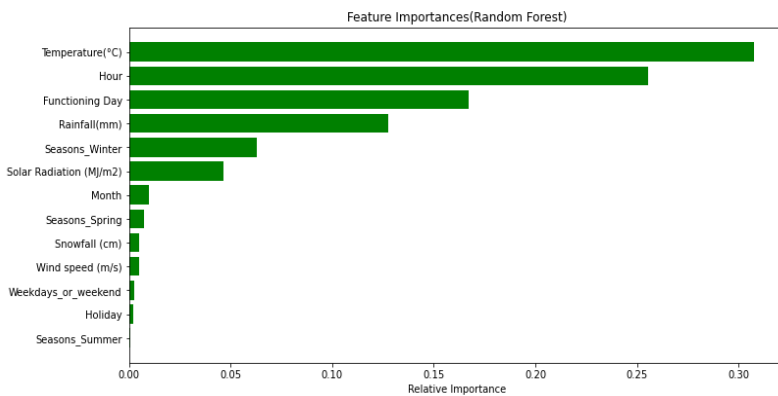
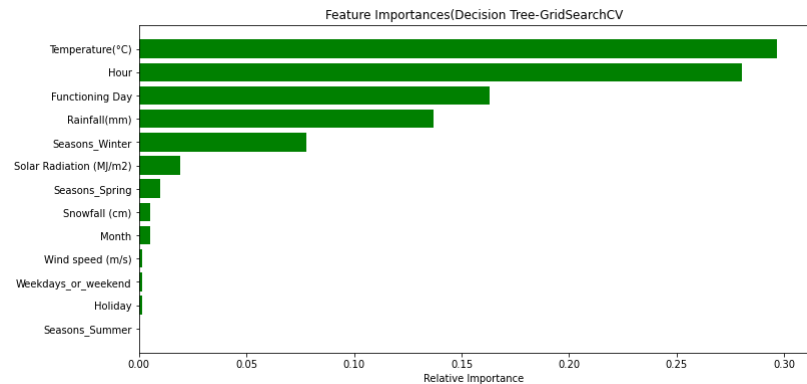
The Mean Absolute Error (MAE) is 2.55614385530248.
The Mean Squared Error(MSE) is 14.756768227597828.
The Root Mean Squared Error(RMSE) is 3.8414539210561705.
The R2 Score is 0.9062980574173423.



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



Features importance's :



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	5.8396	59.7844	7.7320	0.6126
	1	Lasso	5.8541	59.9446	7.7424	0.6116
	2	Ridge GridSearchCV	5.8541	59.9446	7.7424	0.6116
	3	ElasticNet(GridSearchCV-Tunned)	5.8757	60.3548	7.7688	0.6089
	4	Decision Tree Regressor-GridSearchCV	3.0371	20.1658	4.4906	0.8693
	5	Random Forest	1.0111	2.4352	1.5605	0.9842
	6	Random Forest-GridSearchCv	2.7483	16.3752	4.0466	0.8939
	7	Gardient boosting Regression	3.2338	21.1826	4.6025	0.8627
	8	Gradient Boosting Regression(GridSearchCV)	1.7402	6.5102	2.5515	0.9578
Test set	0	Linear Regression	5.9107	60.0907	7.7518	0.6184
	1	Lasso	5.9306	60.4709	7.7763	0.6160
	2	Ridge(GridsearchCv Tunned)	5.9306	60.4709	7.7763	0.6160
	3	ElasticNet(GridSearchCV-Tunned)	5.9575	61.1970	7.8228	0.6114
	4	Decision Tree Regressor(GridsearchCV)	3.3575	24.7144	4.9714	0.8431
	5	Radom forest	2.6370	15.7381	3.9671	0.9001
	6	Random Forest-GridSearchCv	3.0449	19.9709	4.4689	0.8732
	7	Gradient Boosting Regression	3.3368	22.9967	4.7955	0.8540
	8	Gradient Boosting Regression(GridSearchCV)	2.5561	14.7568	3.8415	0.9063

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.

THANK YOU