

Capstone Project

Seoul Bike Sharing Demand Prediction

Problem Statement

Currently Rental bikes are introduced in many urban cities for the enhancement of mobility comfort. 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 crucial part is the prediction of bike count required at each hour for the stable supply of rental bikes.

Content

- ☐ Data Pipeline
- ☐ Data Description
- ☐ Exploratory Data Analysis
- ☐ Models performed
- ☐ Model Validation & Selection
- ☐ Evaluation Matrix of All the models
- ☐ Model Explainability - SHAP
- ☐ Challenges
- ☐ Conclusion

Data Pipeline

- Exploratory Data Analysis (EDA): In this part we have done some EDA on the features to see the trend.
- Data Processing: In this part we went through each attributes and encoded the categorical features.
- Model Creation: Finally in this part we created the various models. These various models are being analysed and we tried to study various models so as to get the best performing model for our project.

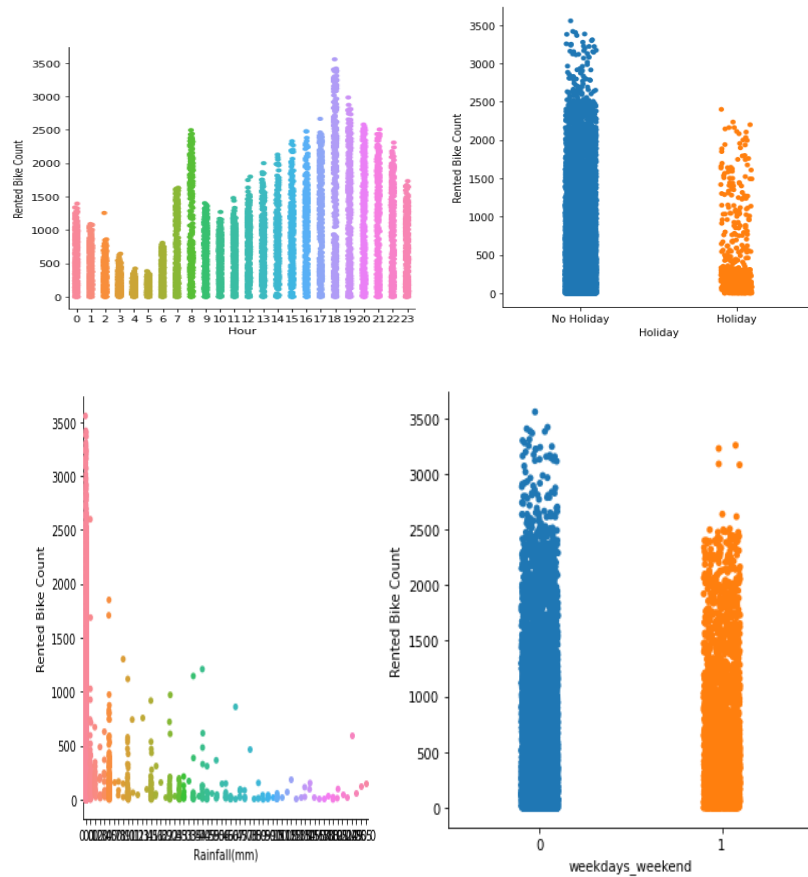
Data Description

Dependent variable:

- Rented Bike count - Count of bikes rented at each hour

Independent variables:

- Date : year-month-day
- Hour - Hour of the day
- Temperature-Temperature in Celsius
- Humidity - %
- Windspeed - m/s
- Visibility - 10 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)



Observation

From all these point plot we have observed a lot from every column like :

Season

In the season column, we are able to understand that the demand is low in the winter season.

Holiday

In the Holiday column, The demand is low during holidays, but in no holidays the demand is high, it may be because people use bikes to go to their work.

Functioning Day

In the Functioning Day column, If there is no Functioning Day then there is no demand

Days of week

In the Days of week column, We can observe from this column that the pattern of weekdays and weekends is different, in the weekend the demand becomes high in the afternoon. While the demand for office timings is high during weekdays, we can further change this column to weekdays and weekends.

month

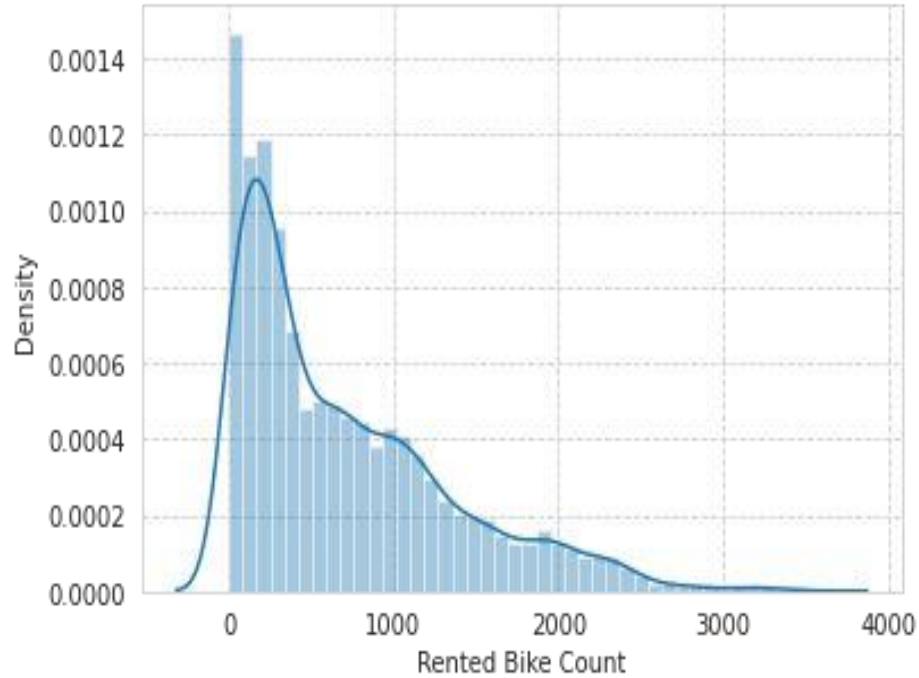
In the month column, We can clearly see that the demand is low in December January & February, It is cold in these months and we have already seen in season column that demand is less in winters.

year

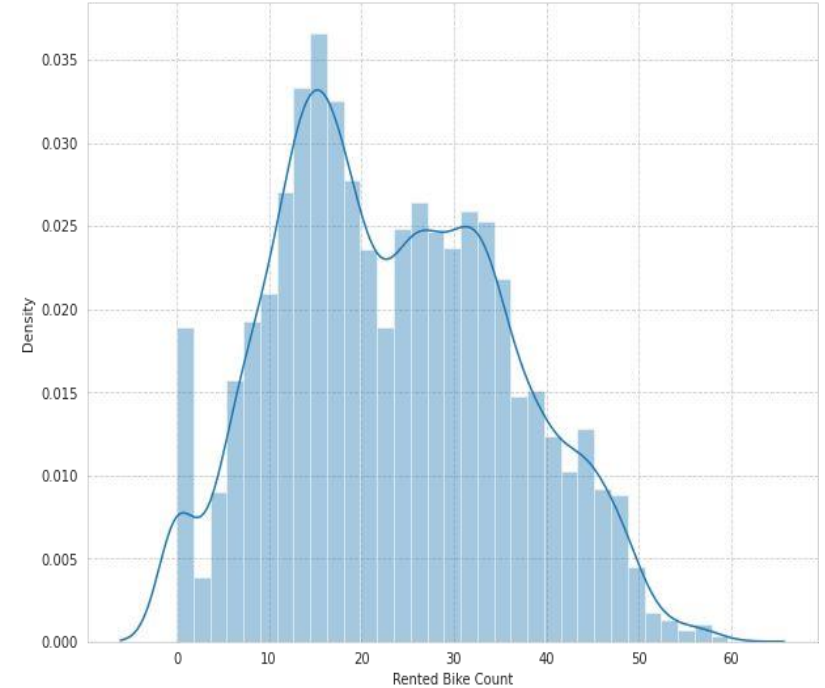
The demand was less in 2017 and higher in 2018, it may be because it was new in 2017 and people did not know much about it.

Correlation Graph

EDA (contd...)

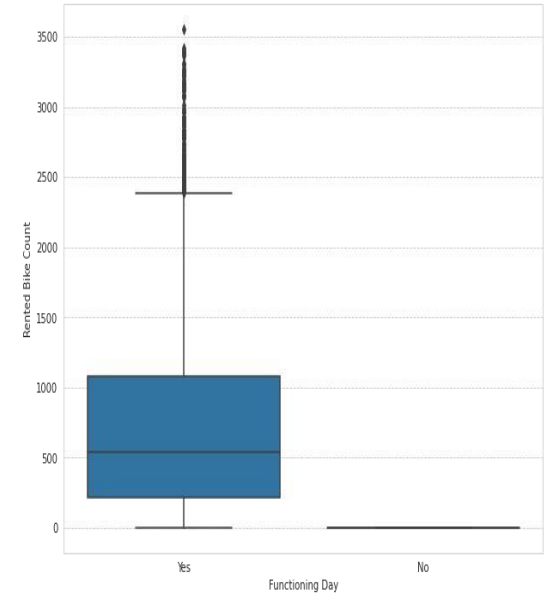
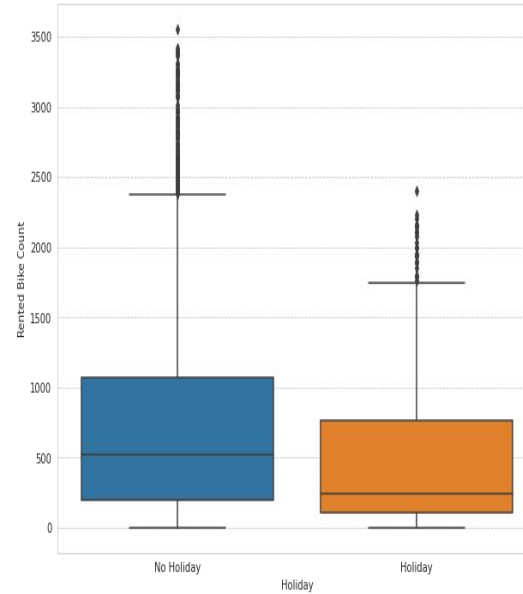
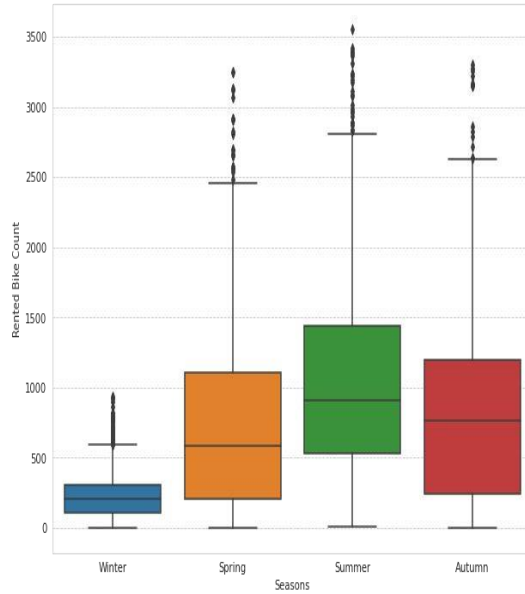


Distribution of rented bike count



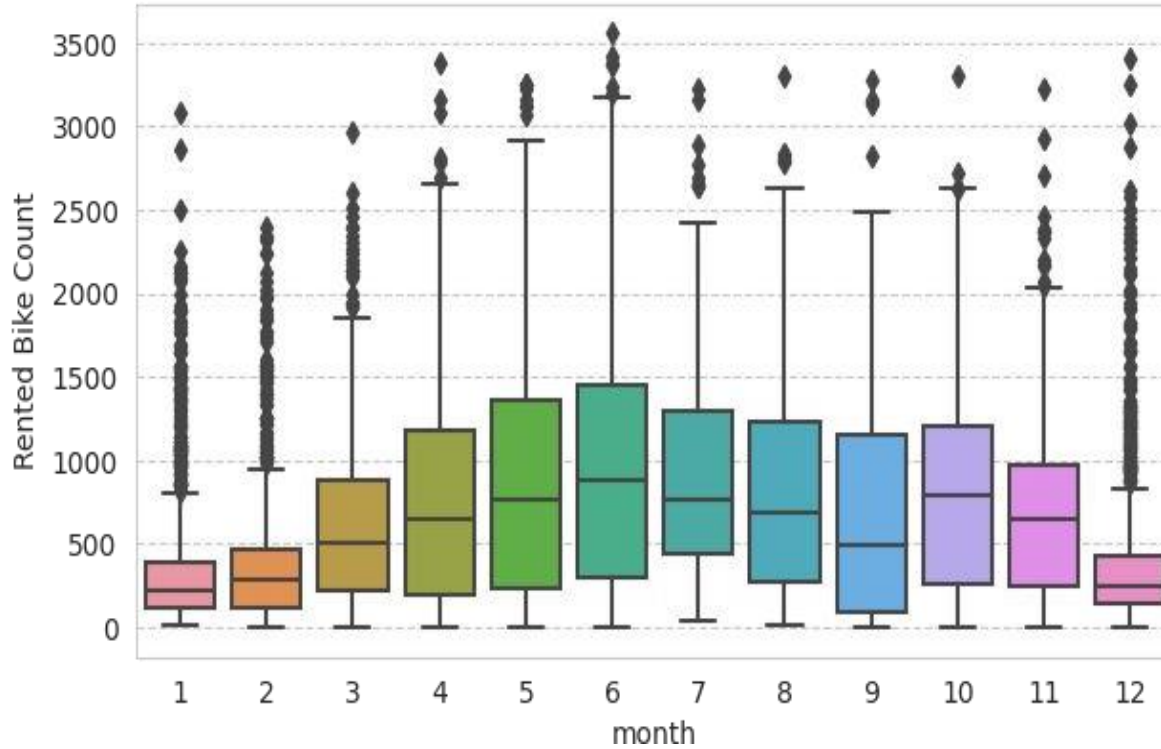
Square root transformation of rented bike count

EDA (contd...)



- Less demand on winter seasons
- Slightly Higher demand during Non holidays
- Almost no demand on non functioning day

EDA (contd...)



- We can see that there is less demand of Rented bike in the month of December, January, February i.e. during winter seasons
- Also demand of bike is maximum during May, June, July i.e. Summer seasons



Correlation Graph

From the correlation graph with Heat map we saw that dew point temp and temperature is highly correlated. Then we checked VIF and concluded that these two features are affecting VIF score also. so we decided to drop one of these feature and to do this we checked which feature is least correlated with Dependent variable and we identified it to be Dew point temperature and therefore we dropped the Dew point temperature.

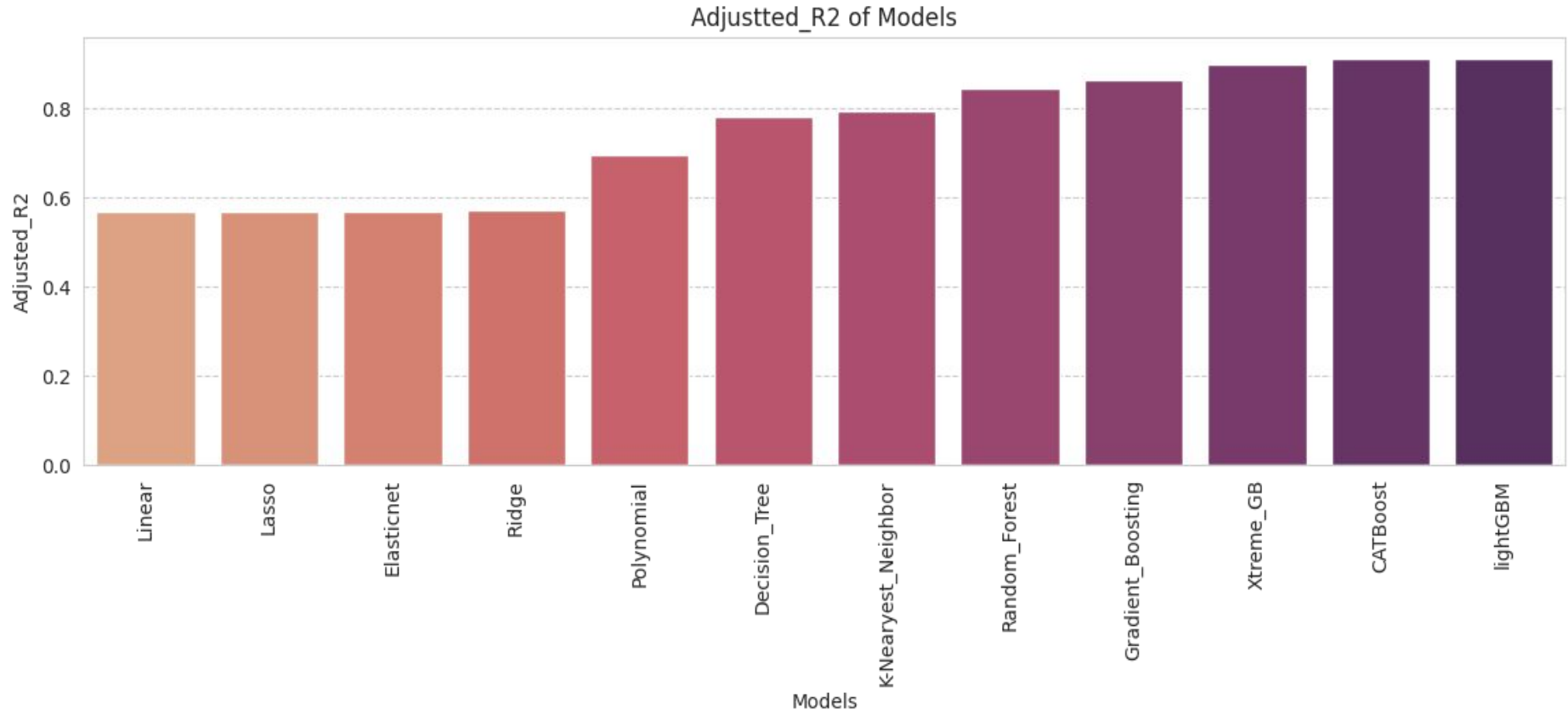
Model's Performed

- Linear Regression with regularizations
- Polynomial Regression
- K nearest neighbours
- Decision tree
- Random forest
- Gradient Boost
- eXtreme Gradient Boost
- lightGBM
- CatBoost

	Models	Mean_square_error	Root_Mean_square_error	R2	Adjusted_R2
0	Linear	175590.552873	419.035264	0.572911	0.569766
1	Lasso	175560.907118	418.999889	0.572983	0.569839
2	Ridge	175248.935066	418.627442	0.573742	0.570603
3	Elasticnet	175346.867499	418.744394	0.573504	0.570363
4	Polynomial	123952.860328	352.069397	0.698509	0.696289
5	K-Nearyest_Neighbor	83411.759209	288.810940	0.796159	0.794659
6	Decision_Tree	88506.087215	297.499726	0.783710	0.782117
7	Random_Forest	62790.180423	250.579689	0.846554	0.845424
8	Gradient_Boosting	55090.172685	234.712958	0.865371	0.864380
9	Xtreme_GB	40812.801816	202.021785	0.900262	0.899528
10	CATBoost	36339.421527	190.629015	0.911194	0.910540
11	lightGBM	35410.753754	188.177453	0.913464	0.912826

The Best model is Random Forest but it is over fitted, that's why We are using Hyperparameter tuning so that we can reduce the overfitting and increase the accuracy.

Adjusted R2 of Model's Performed

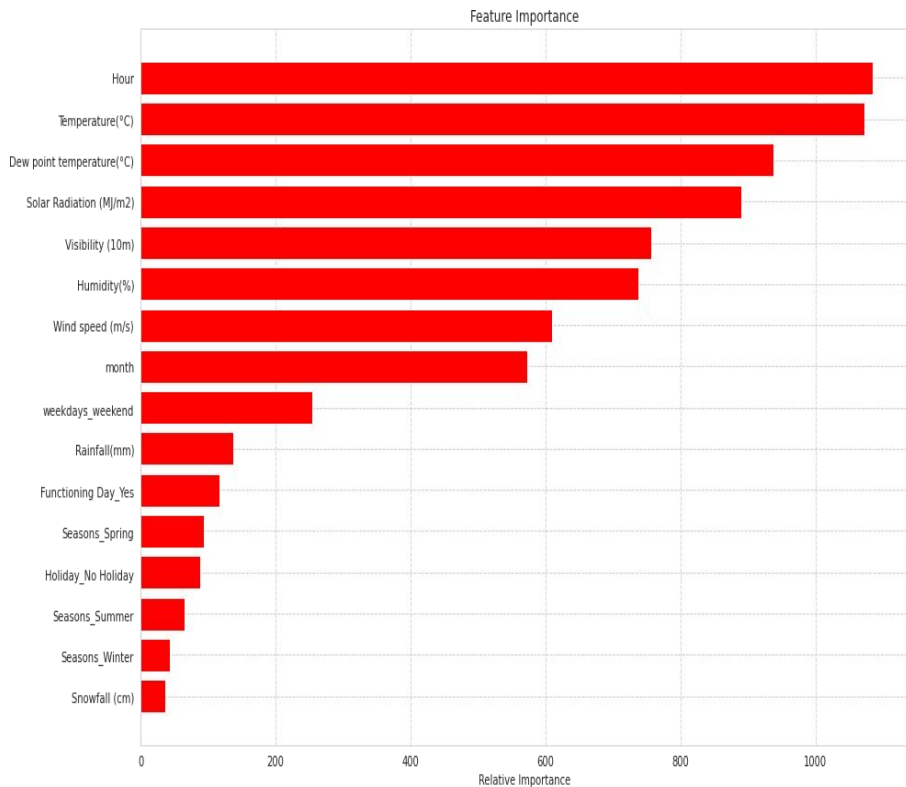


Model Validation & Selection(continued)

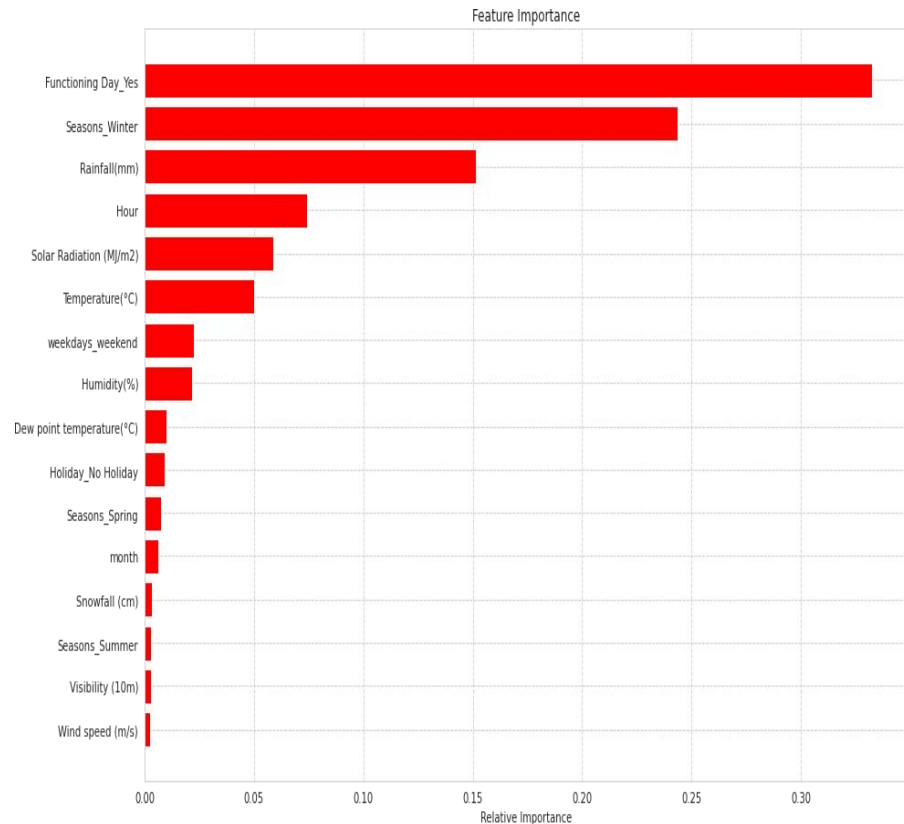
- **Observation 1:** As seen in the Model Evaluation Matrices table, Linear Regression, KNN is not giving great results.
- **Observation 2:** Random forest & GBR have performed equally good in terms of adjusted r^2 .
- **Observation 3:** We are getting the best results from lightGBM and CatBoost.



Feature Importance

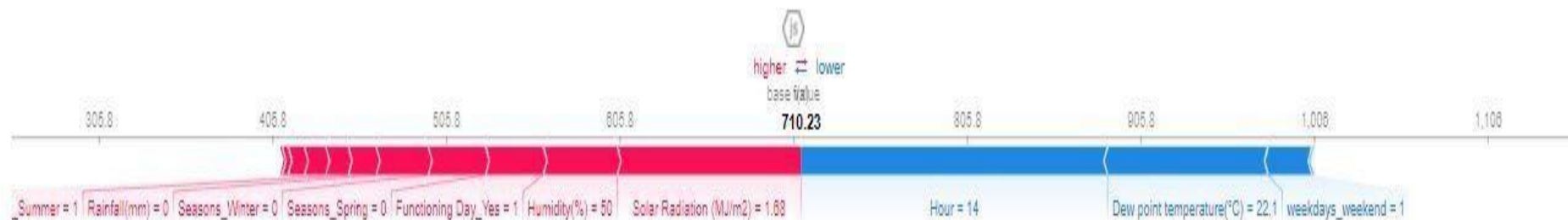


lightGBM



CatBoost

Model Explainability - SHAP



lightGBM



CatBoost

Challenges

- A huge amount of data needed to be dealt while doing the project which is quite an important task and also even small inferences need to be kept in mind.
- As dataset was quite big enough which led more computation time.



Conclusion

- It is quite evident from the results that lightGBM and Catboost is the best model that can be used for the Bike Sharing Demand Prediction since the performance metrics (mse,rmse) shows lower and (r2,adjusted_r2) show a higher value for the lightGBM and Catboost models.
- So, we can use either lightGBM or catboost model for the above problem



**THANK
YOU**