

BoomBikes: Linear Regression Analysis

By Prem Kumar Subudhi



Problem Statement

Boom Bikes going loss because of Covid and finding difficult to sustain, need to come up with a business plan to recover from the current situation and book profit in future.

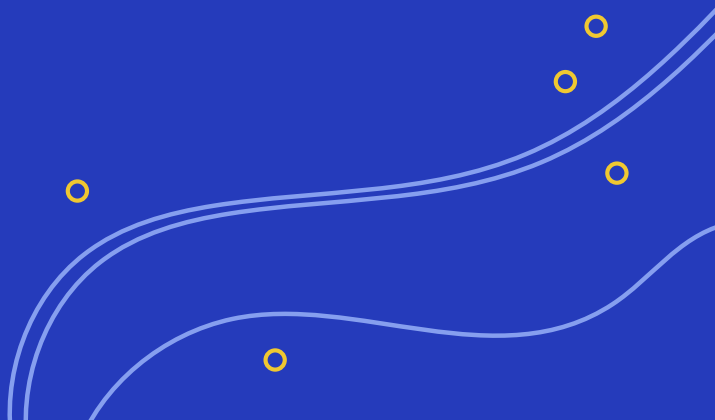


Task

Find the factors on which the demand for these shared bikes depends. Specifically, they want to understand the factors affecting the demand for these shared bikes in the American market. The company wants to know the variables those are significant in predicting the demand for shared bikes and how well those variables describe the bike demands.

The Goal of problem statement

1. Data Understanding and loading, drop un-necessary features that does not contribute for the model.
2. Create a linear model with all the features that may contribute or effect the price (i.e CNT feature)
3. The model should be interpretable and easily understandable.

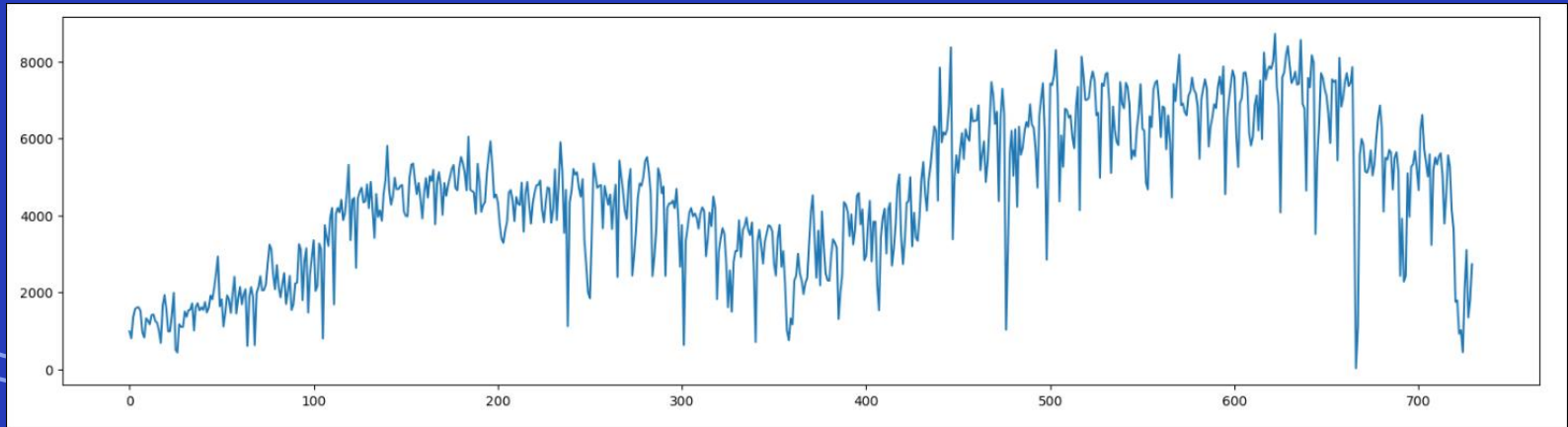




Linear regression model

Simple / multiple linear regression

The feature C_n has
total rental bikes including both casual and registered, which is continuous
in nature as shown in below diagram.



Linear regression model evaluation



The measure R-squared

Mercury is the closest planet to the Sun and **the smallest one in the entire Solar System**



Adjusted R-squared

Venus has a beautiful name and is **the second planet from the Sun**. It's hot and its atmosphere is poisonous



Residual analysis

Earth is the third planet from the Sun and **the only one that harbors life** in the Solar System



Assumptions of linear regression

Linearity & normality

The change in independent variable increase proportionately with independent variables. If the residuals are not normally distributed, it might affect the accuracy of confidence intervals and hypothesis tests

Independence

All the features used in the model are assumed to be independent with each other, otherwise the model will introduce errors.

Homoscedasticity

When this assumption is met, it implies that the model's predictions are equally accurate across all levels of the independent variable(s).

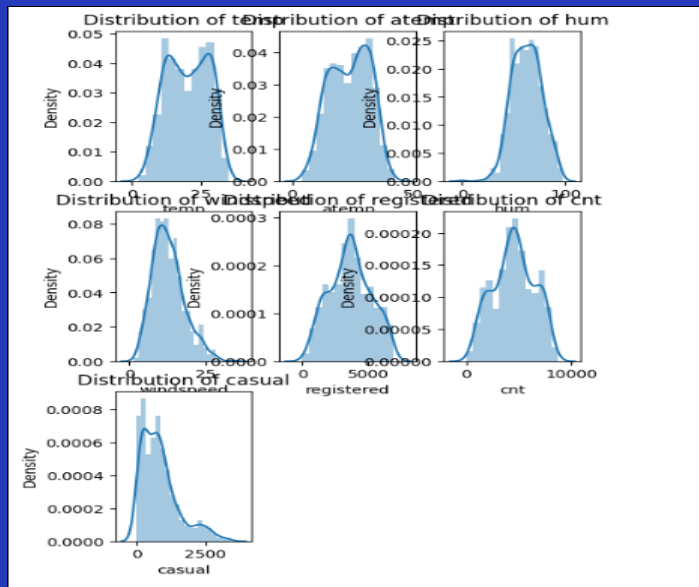
No multicollinearity

Above independence nature led to non multicollinearity, the model will become weak when it has multicollinearity introduced.



Distribution Plot

Visualizations



Analysis

Despite being red, Mars is **actually a cold place**. It is full of iron oxide dust

Assessment

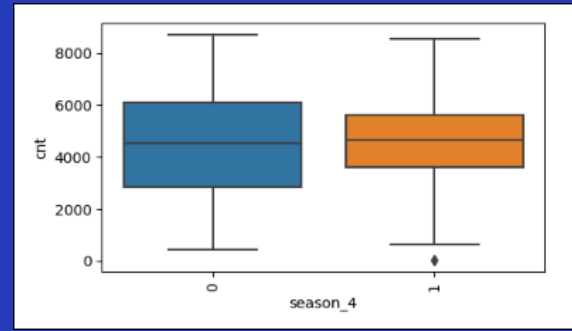
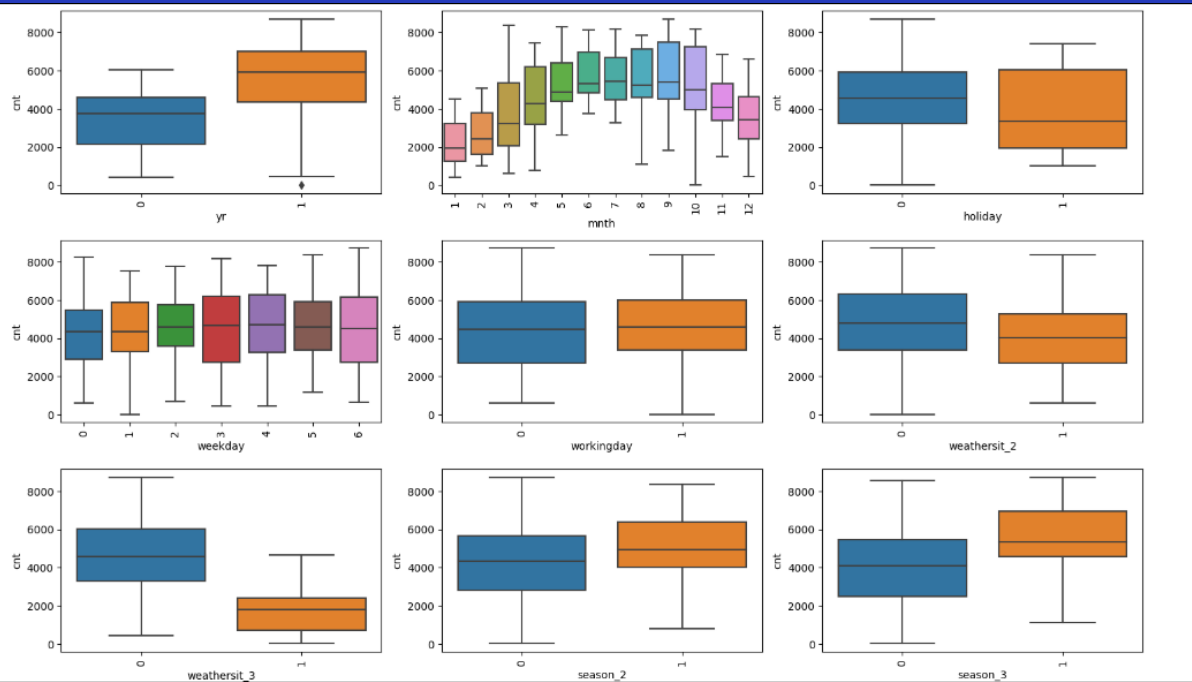
Jupiter is a gas giant and **the biggest planet** in the entire Solar System

Patterns

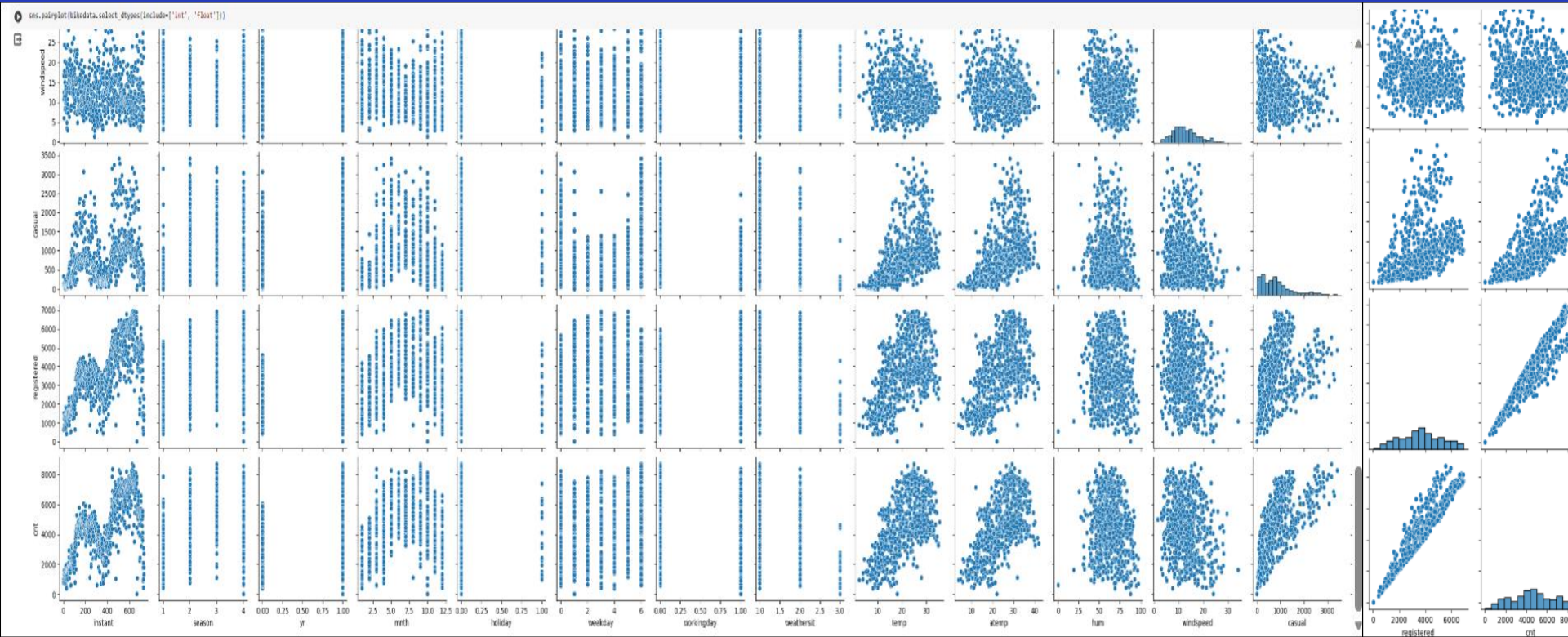
Earth is the third planet from the Sun and **the only one that harbors life**

Data splits

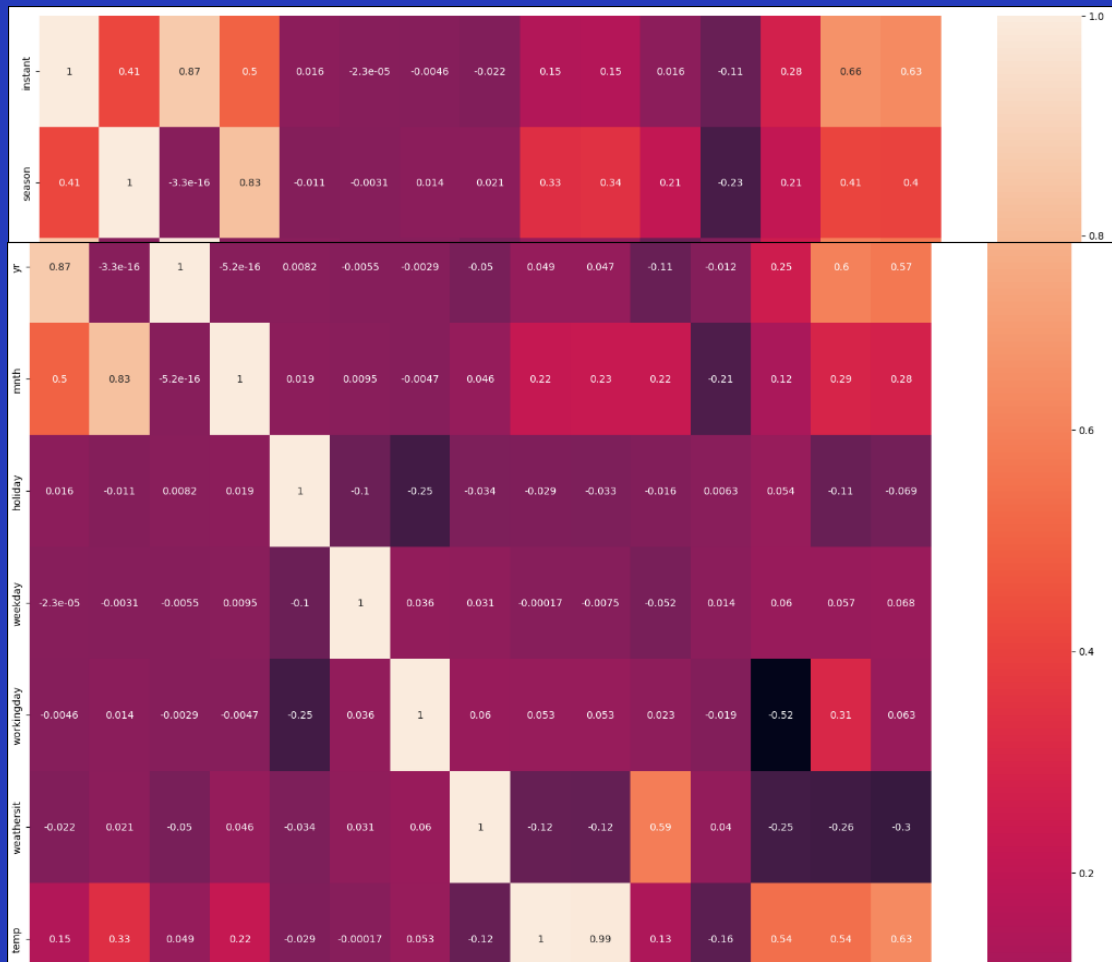
Created dummy variables for Categorical variables



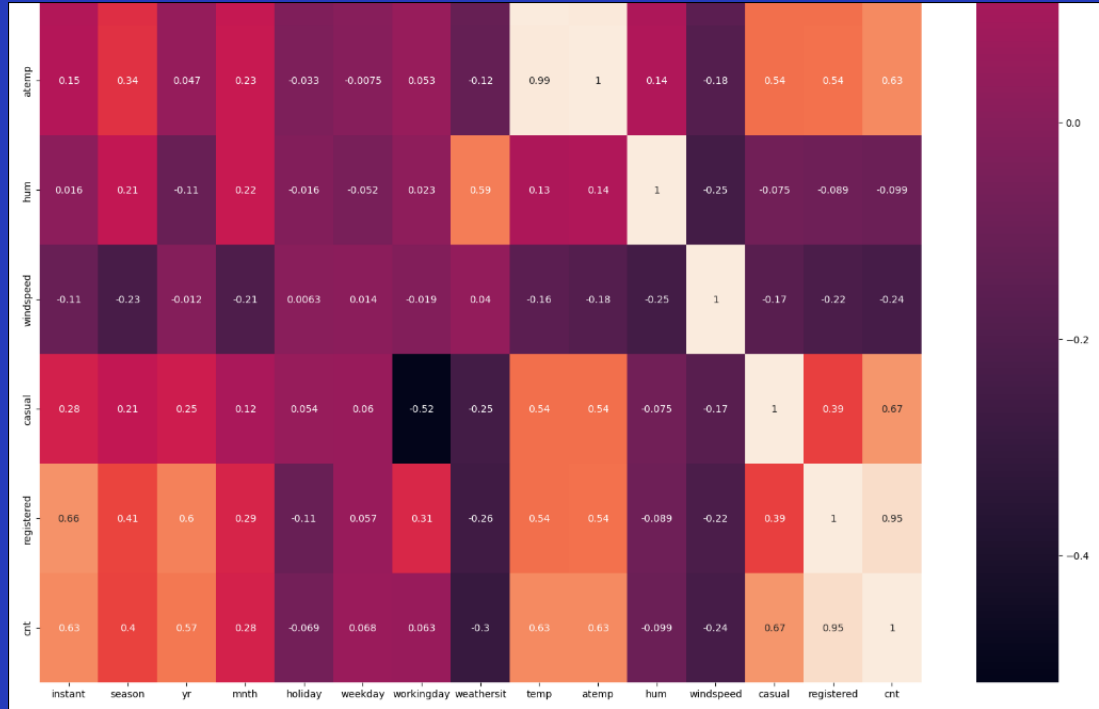
Seaborn Pairplot



Heat Map



Heat Map



Linear regression analysis: key steps

Data preparation

Convert the Categorical variables to dummy variables and remove the first column

Step 01

Step 02

Step 03

Step 04

Model evaluation

R²_Score of the model expected to be high

Model building

Regression model
$$Y = \text{constant} + \sum (\text{Beta}_0 * \text{Independent variables})$$

Results

Predicted result from the model

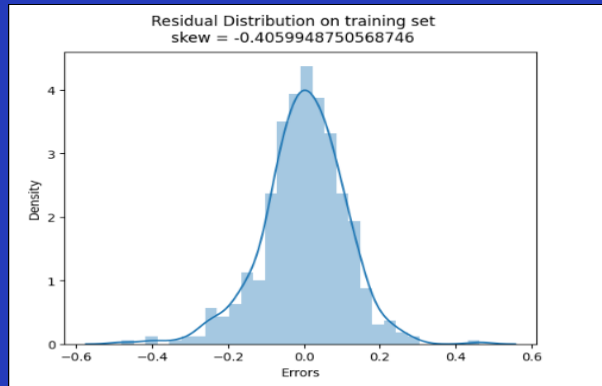
Linear regression key information

| Aspect of linear regression | Description/Information |
|---------------------------------|--|
| Purpose | Predictive modeling, understanding relationships, hypothesis testing, etc. |
| Assumptions | Linearity, independence, homoscedasticity, normality, no multicollinearity |
| Model fit metrics | R-squared, adjusted R-squared, F-statistic, p-values |
| Data preparation | Collect, clean, and preprocess data. Handle missing values and outliers |
| Exploratory Data Analysis (EDA) | Visualize data through scatter plots, histograms, correlation matrices |

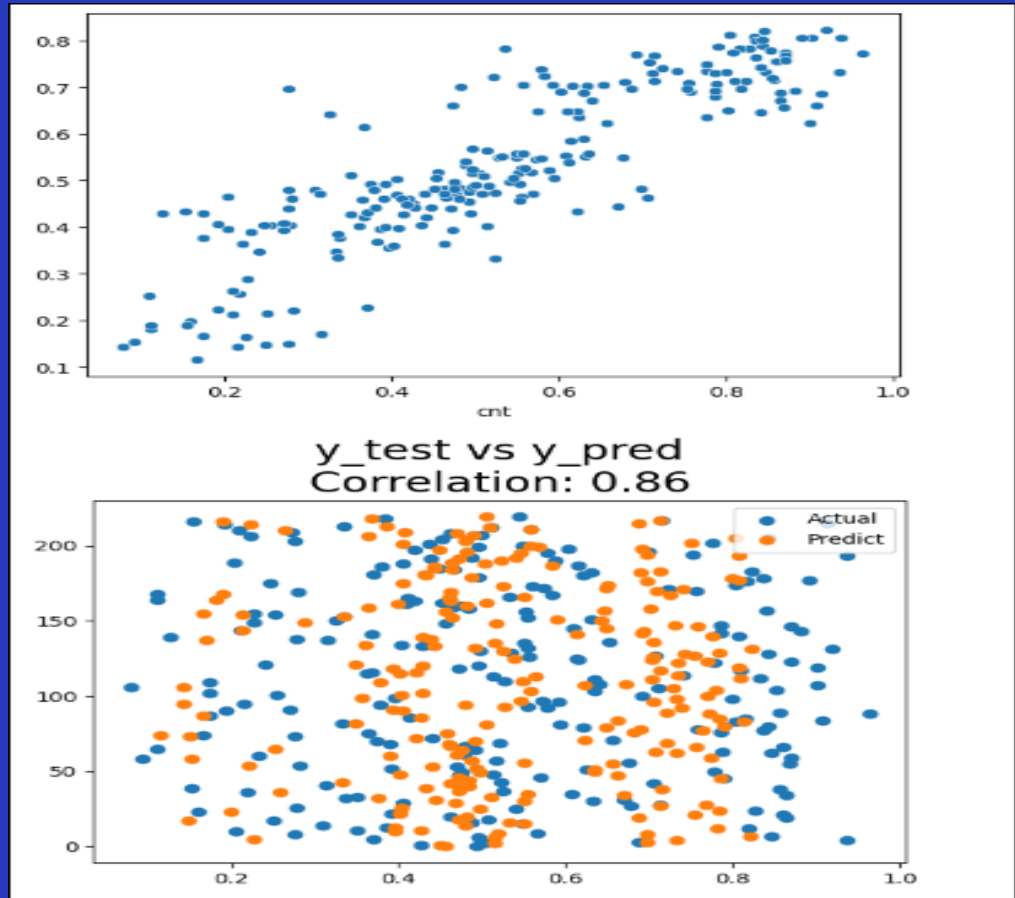
Diagnostic tools for model validation

Residual analysis

Residual plot



scatter diagram
for the test set
vs predicted.

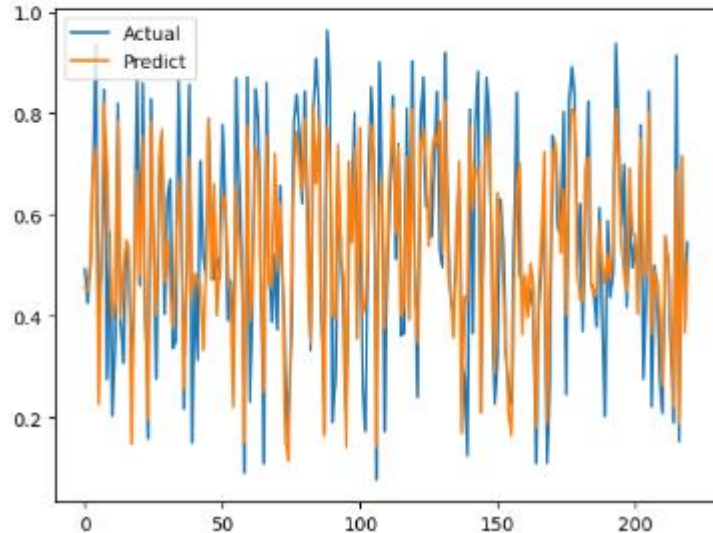


Predicted values from the model

Plot the CNT test-actual vs test-predicted

```
plt.plot(list(y2_test),label='Actual')  
plt.plot(list(y2_test_pred),label='Predict')  
plt.legend()
```

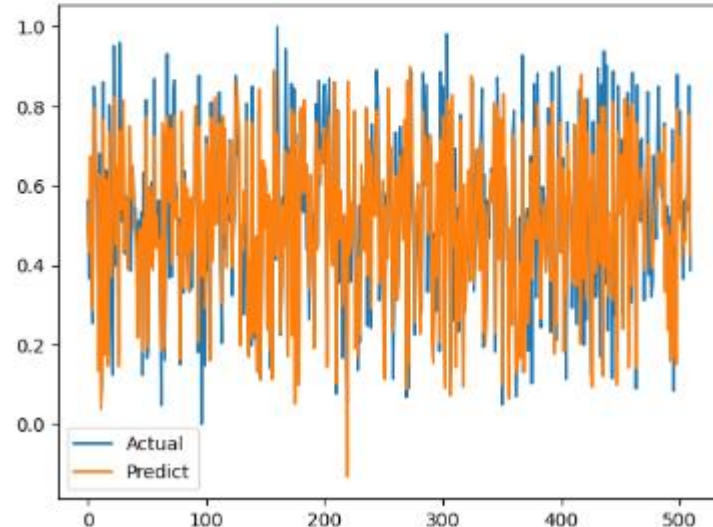
<matplotlib.legend.Legend at 0x78b504173610>



Plot the CNT trained-actual vs trained-predicted

```
plt.plot(list(y1_train),label='Actual')  
plt.plot(y1_train_pred, label='Predict')  
plt.legend()
```

<matplotlib.legend.Legend at 0x78b510b90a00>



The model created :

$$\text{CNT} = 0.26 + .24 * \text{yr} + 0.04 * \text{weekday} - 0.15 * \text{windspeed} + 0.27 * \\ \text{season_2} + 0.31 * \text{season_3} + 0.23 * \text{season_4} - 0.07 * \\ \text{weathersit_2} + -0.28 * \text{weathersit_3}$$

(As holiday has high P-value it will be excluded.)

OLS vs XGBoost

Train R2 Value ::> 0.919340607472633

Test R2 Value: 0.8516

Feature Importances:

| | |
|----------------|----------|
| yr | 0.366905 |
| mnth | 0.112416 |
| holiday | 0.012894 |
| weekday | 0.010726 |
| workingday | 0.019876 |
| atemp | 0.208465 |
| hum | 0.039786 |
| windspeed | 0.021618 |
| season_2 | 0.011544 |
| season_3 | 0.004960 |
| season_4 | 0.080601 |
| weathersit_2 | 0.031936 |
| weathersit_3 | 0.078272 |
| dtype: float32 | |

The XG boost method give the R2 value better than compared to OLS result.

The background is a solid blue color. It features several white, wavy, horizontal lines that flow from the left and right edges towards the center. Scattered throughout the background are small, yellow-outlined circles. The word "End" is centered in a large, white, sans-serif font.

End