

Predicting Bike Rentals Using Multiple Linear Regression

A Case Study of BoomBikes

Sathish Madhiyalagan

Objective

The primary objective of this project is to develop a predictive model using multiple linear regression to accurately forecast the daily bike rental counts for BoomBikes. By analyzing various influencing factors such as weather conditions, day of the week, and seasonal trends, the model aims to provide actionable insights that can help BoomBikes optimize their operations, enhance customer satisfaction, and increase revenue.

Software and Libraries

- ❖ **Python:** Programming language for implementing the machine learning model.
- ❑ Pandas For data manipulation and analysis.
- ❑ Numpy For numerical computations.
- ❑ Matplotlib and seaborn For data visualization.
- ❑ Scikit-learn For building and evaluating the regression model.

Step 1 : Data Import and Preliminary Exploration

Importing and understanding the dataset is the initial step in our analysis. We begin by importing the data into our Python environment and gaining an overview of its structure and contents. Let's explore the dataset's features and their descriptions in data dictionary.

Upon loading and examining the dataset, it becomes evident that our primary objective is to predict the total bike rentals, as represented by the '**cnt**' column. This understanding guides our initial exploration and subsequent analysis.

Step 2 : EDA

- Print information about the dataset, such as the data types of each column and the number of non-null values.
- Calculate and display summary statistics (e.g., mean, standard deviation, min, max) for numerical columns.
- Identify and print the number of missing values in each column to assess data completeness.

- Determine and print the number of duplicate rows in the dataset, if any.
- Plot a histogram to visualize the distribution of the target variable 'cnt', which represents rental counts.

These steps are EDA first steps, helping to understand its characteristics and identify any potential issues or patterns.

After thoroughly reviewing my dataset, I'm pleased to confirm that everything appears to be in good order, with no missing values, duplicates.

Then we go with **Univariate analysis, Bivariate Analysis and Multivariate Analysis.**

- Univariate Analysis (categorical variables or numerical variables)

Univariate analysis examines individual variables to understand their distribution and properties, using techniques like histograms and summary statistics.

- Bivariate Analysis (Categorical vs. Categorical, Categorical vs. Numerical and Numerical vs. Numerical)

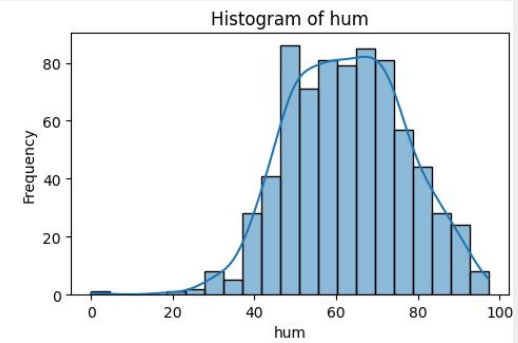
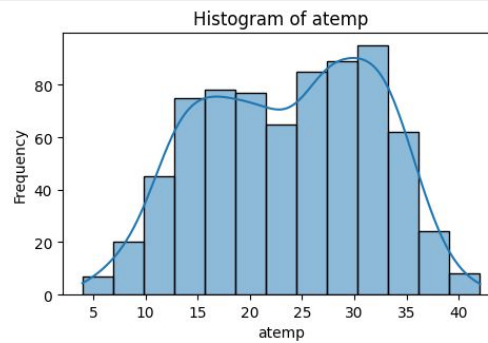
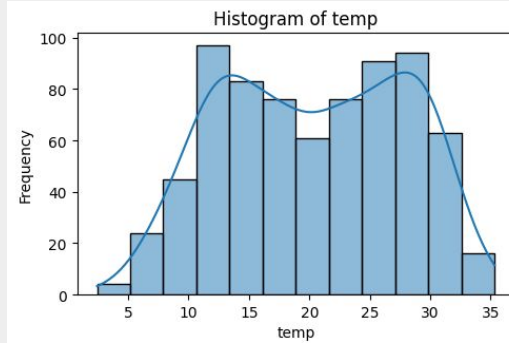
Bivariate analysis explores relationships between pairs of variables, revealing patterns and correlations through scatter plots and cross-tabulations.

- Multivariate Analysis (Its combined all)

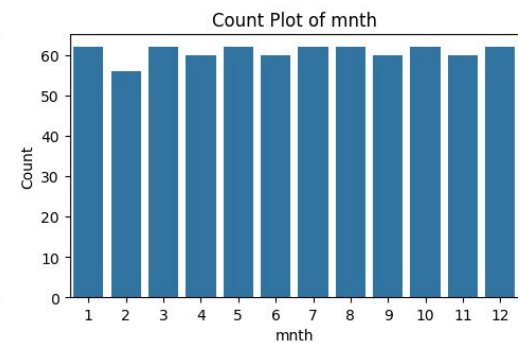
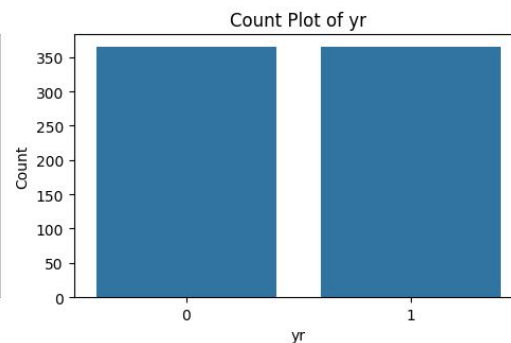
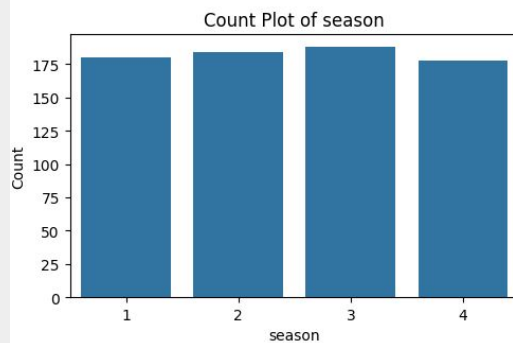
Multivariate analysis delves into interactions among multiple variables, uncovering complex patterns and relationships using techniques like PCA, factor analysis, and multiple regression.

Univariate analysis

Histograms for numerical variables

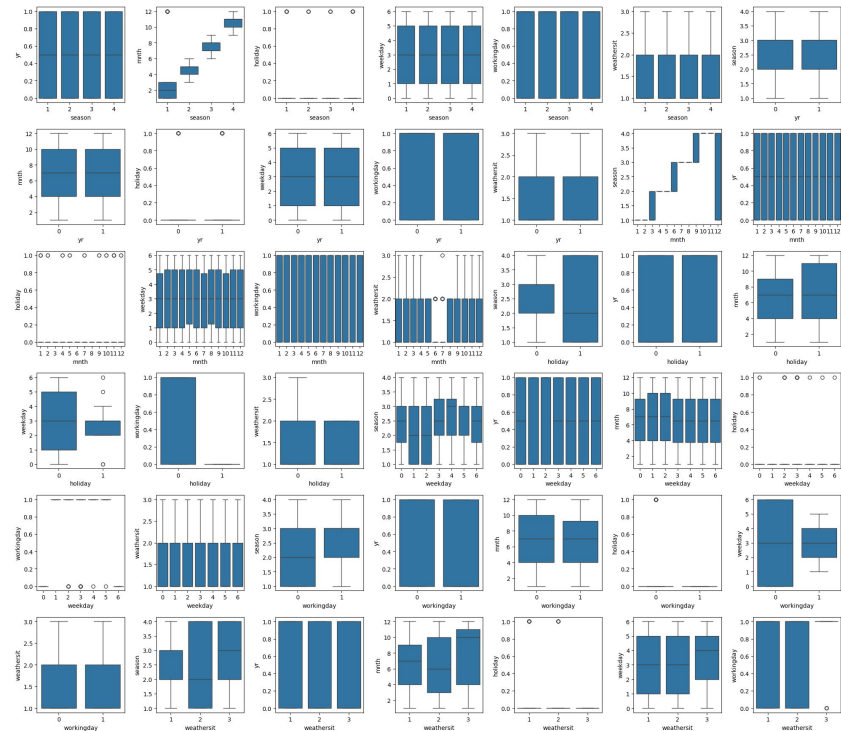
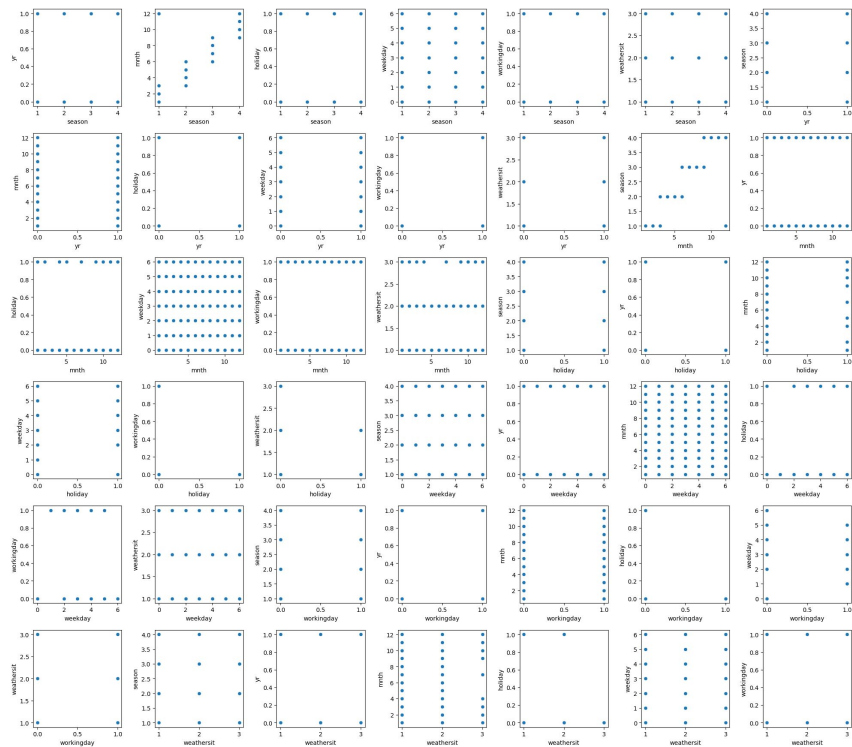


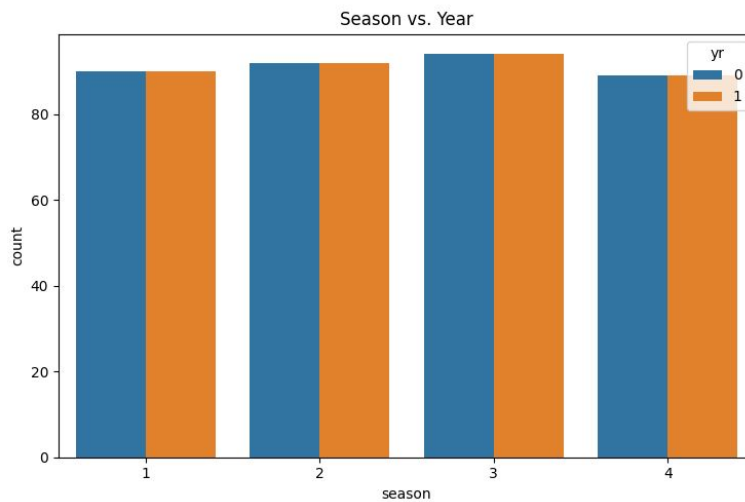
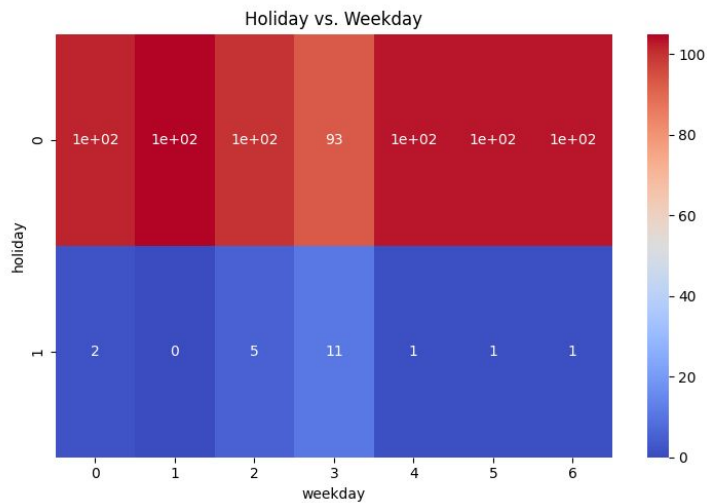
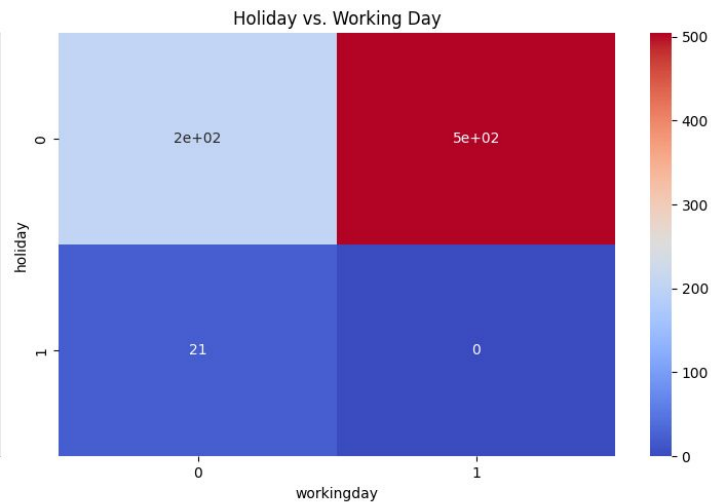
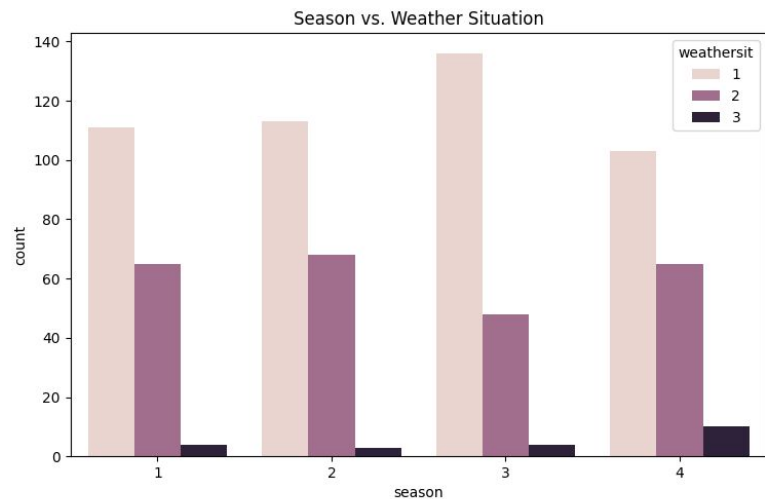
Count plots for categorical variables



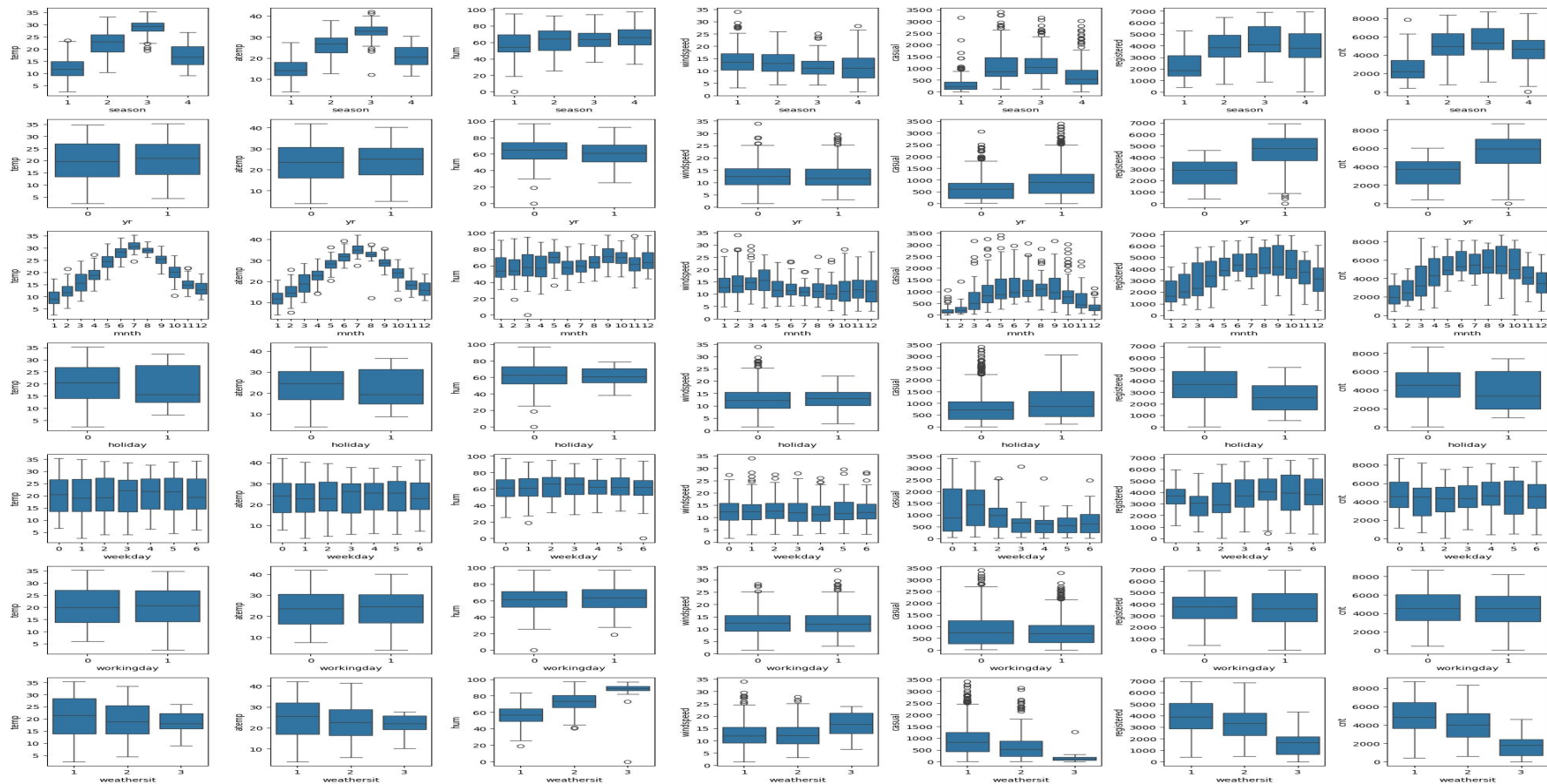
Bivariate Analysis

Categorical vs. Categorical

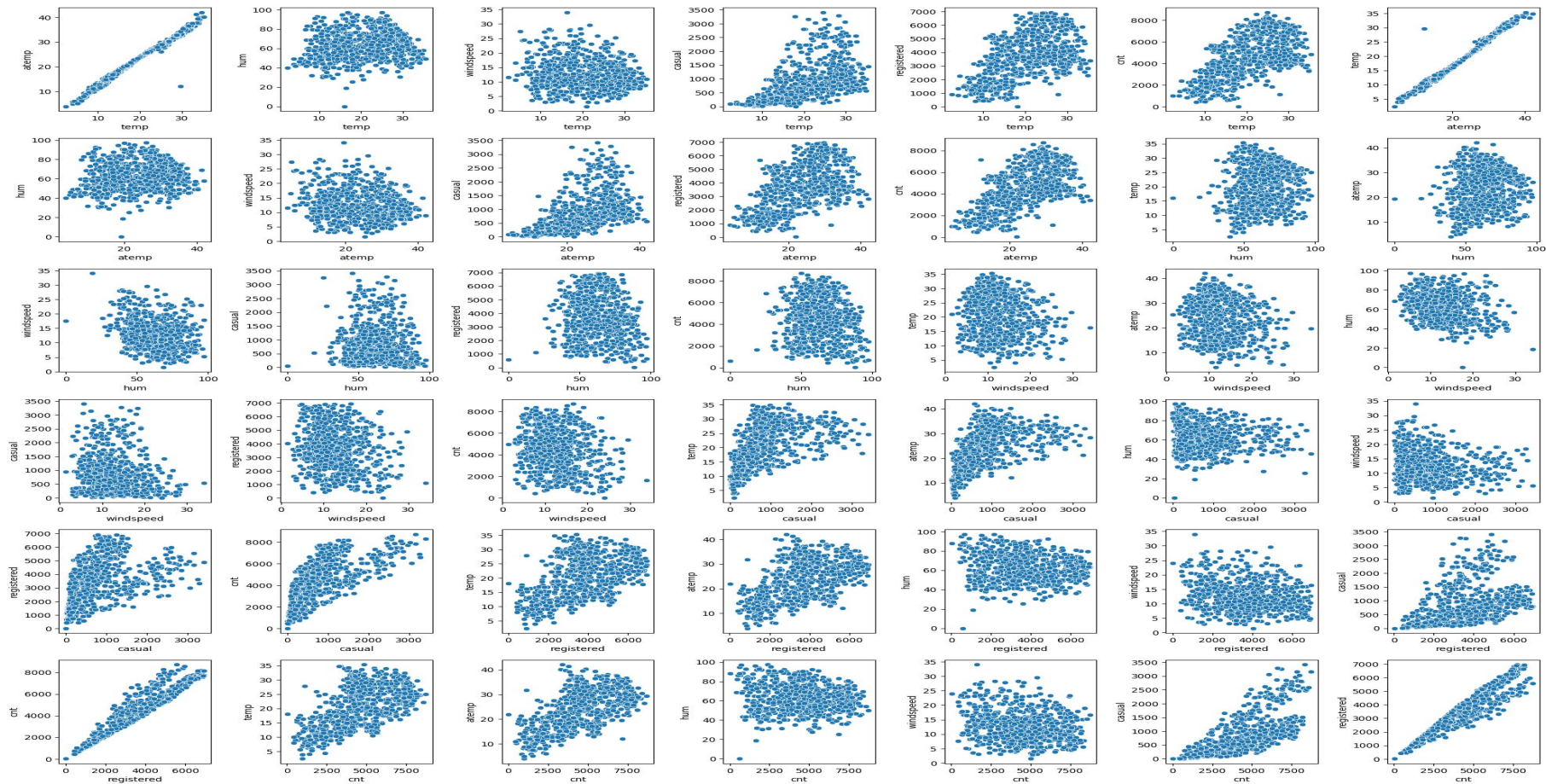




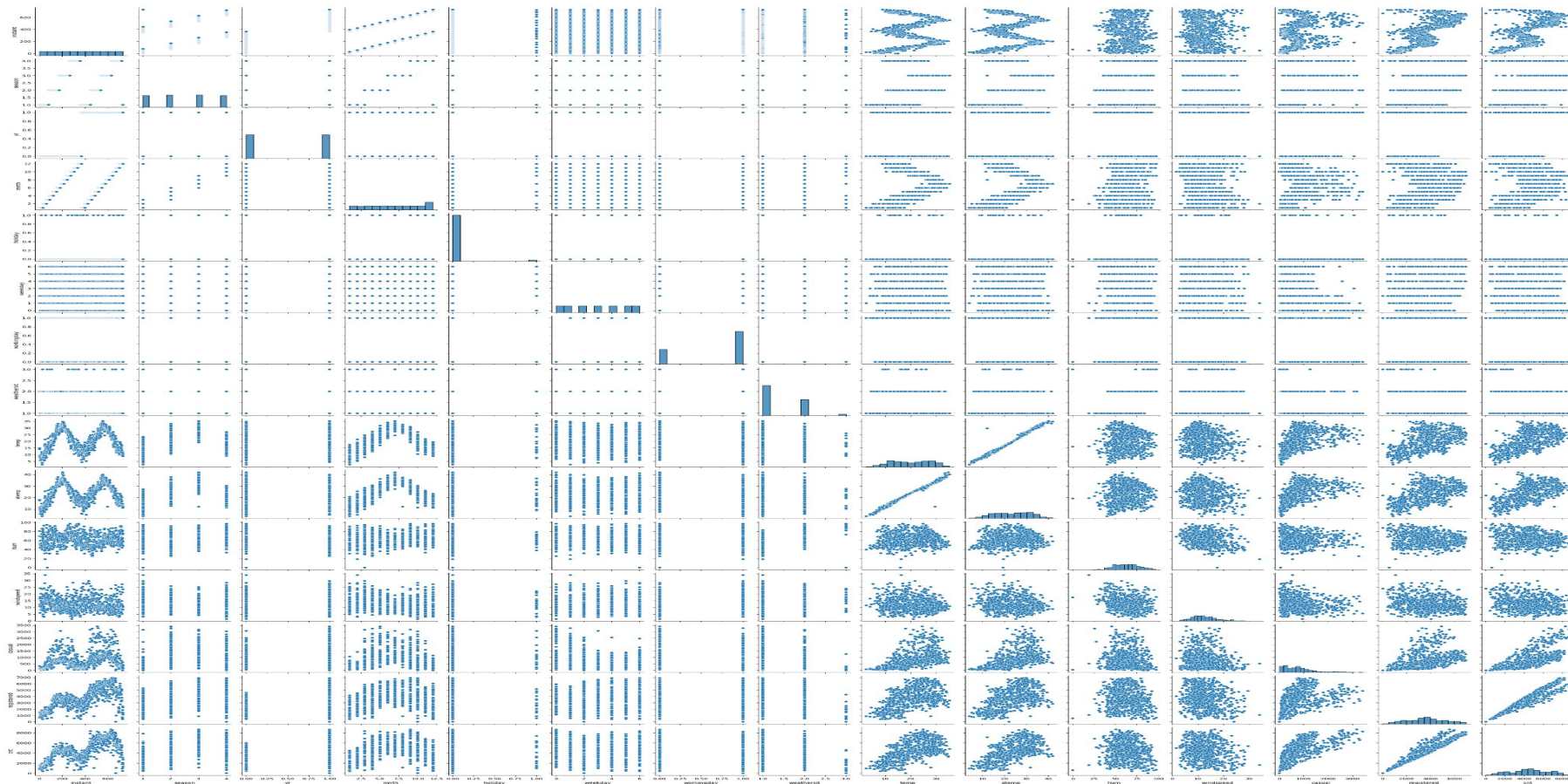
Categorical vs. Numerical



Numerical vs. Numerical



Multivariate Analysis



Step 3 : Pre Processing

In the preprocessing phase, it's crucial to identify and remove any columns that are not relevant to the prediction task or those that do not contribute meaningful information to the model. Based on my statement, it seems four columns that are not correct or necessary for this feature.

That four columns are (instant, dteday, casual, registered)

One more time check ,Missing Values and null values ,Categorical to encode numerical format

Step 4 : Dummy Variables Creation

Convert the selected columns (season, mnth, weekday, and weathersit) to categorical data type. This step is necessary for the next step, where we will create dummy variables.

Generate dummy variables for the categorical columns specified in `columns_to_encode`. The `drop_first=True` parameter is used to drop the first category in each column to avoid multicollinearity. This process creates new binary columns representing the presence of each category in the original columns.

Concatenate the original DataFrame `boomBikesAfterDrop` with the newly created dummy variables (`status`). This step appends the dummy variable columns to the original DataFrame.

Drop the original categorical columns (`season`, `mnth`, `weekday`, and `weathersit`) from the combined DataFrame (`bick`). This step ensures that only the dummy variables remain in the dataset, **preventing redundancy and multicollinearity**.

Step 5 : Splitting the Data into Training and Testing Sets

To effectively train and evaluate a machine learning model, it is crucial to divide the dataset into training and testing sets. This allows the model to learn from the training set and be evaluated on unseen data in the testing set, helping to ensure that the model generalizes well to new data.

- Import the `train_test_split` function from the `sklearn.model_selection` module. This function is used to split the dataset into training and testing sets.
- Set the random seed using `np.random.seed(0)`. This ensures that the splitting of data is reproducible, meaning that every time you run the code, the training and testing sets will have the same rows.
- Use the `train_test_split` function to split the `bickDf` DataFrame into training and testing sets. The `train_size=0.7` parameter specifies that 70% of the data should be used for training, while the `test_size=0.3` parameter specifies that 30% should be used for testing. The `random_state=100` parameter ensures that the split is consistent every time the code is run.

Rescaling the Features :

- Specify the list of numerical columns that need to be scaled. Import the `MinMaxScaler` class from the `sklearn.preprocessing` module to perform scaling.
- Create an instance of the `MinMaxScaler`. Fit the `MinMaxScaler` on the training set's numerical columns and transform the data, scaling the values to a range of `[0, 1]`.
- Display summary statistics for the scaled numerical columns to understand their distribution after scaling.
- Create a heatmap to visualize the correlations between the scaled numerical features in the training set.
- This visualization helps identify strongly correlated features, which can inform feature selection and engineering decisions in the modeling process.

So, we pick temp as the first variable and we'll try to **fit a regression line** to that.

yr	1	0.015	-0.003	0.11	0.1	-0.083	0.001	0.59	0.014	0.044	-0.023	-0.020	0.005	0.054	-0.026	0.01	0.018	0.055	0.002	0.060	0.19	-0.04	-0.034	0.026	0.011	0.033	-0.034	0.018	0.015	-0.061		
holiday	-0.015	1	-0.25	-0.066	0.71	-0.029	0.018	-0.096	0.063	-0.04	0.051	0.057	-0.053	0.009	0.049	0.046	0.047	-0.053	0.044	0.005	0.43	-0.004	0.069	0.035	0.19	-0.066	-0.032	0.025	0.038	-0.028		
workingday	-0.003	-0.25	1	0.076	0.003	0.021	0.002	0.029	0.007	0.03	-0.033	0.036	0.003	0.008	0.085	0.022	0.032	-0.012	0.013	0.035	0.026	0.026	0.025	0.27	0.23	0.16	0.26	0.25	-0.61	-0.069	0.0083	
temp	-0.11	-0.066	0.007	1	0.99	0.16	-0.19	0.64	0.13	0.7	-0.23	-0.3	-0.18	0.051	0.16	0.3	0.41	0.39	0.21	-0.019	-0.19	-0.27	0.038	0.002	0.046	0.075	-0.02	0.025	-0.09	-0.036		
atemp	0.1	-0.073	0.003	0.99	1	0.17	-0.22	0.65	0.14	0.67	-0.21	-0.3	-0.18	-0.04	0.16	0.3	0.41	0.36	0.21	-0.005	0.19	-0.26	-0.037	0.002	0.043	0.079	-0.019	0.012	-0.086	0.044		
hum	-0.085	0.029	0.021	0.16	0.17	1	-0.27	-0.06	-0.013	0.04	0.17	-0.13	-0.11	-0.13	0.15	-0.095	0.054	0.053	0.15	0.16	0.004	0.058	-0.03	0.034	0.026	-0.057	0.052	0.002	0.48	0.25		
windspeed	-0.001	0.018	0.002	-0.19	-0.22	-0.27	1	-0.25	0.11	-0.19	-0.091	0.13	0.14	0.18	-0.034	0.035	0.085	-0.11	-0.12	-0.019	0.01	-0.061	0.035	0.036	-0.058	0.022	0.002	0.016	-0.03	0.087		
cnt	0.59	-0.090	0.008	0.64	0.65	-0.06	-0.25	1	0.13	0.37	0.033	-0.27	-0.14	0.023	0.11	0.19	0.16	0.23	0.2	0.066	-0.048	-0.14	-0.071	-0.019	-0.02	0.11	-0.041	0.012	-0.18	-0.23		
season_2	-0.014	-0.063	0.03	0.13	0.14	-0.013	0.11	0.13	1	-0.34	-0.33	-0.15	0.088	0.51	0.53	0.25	-0.16	-0.19	-0.17	-0.18	-0.18	-0.17	-0.017	0.001	0.12	-0.066	-0.018	0.002	0.039	-0.045		
season_3	-0.044	-0.04	-0.033	0.7	0.67	0.04	-0.19	0.37	0.34	1	-0.34	-0.16	-0.2	-0.17	-0.18	0.068	0.48	0.55	0.35	-0.18	-0.18	-0.18	0.009	0.031	-0.041	0.036	0.008	0.022	-0.075	-0.025		
season_4	-0.023	0.051	0.036	-0.23	-0.21	0.17	-0.091	0.033	-0.33	-0.34	1	-0.15	-0.19	-0.17	-0.17	-0.16	-0.17	-0.19	-0.02	0.53	0.53	0.31	0.017	0.045	0.051	-0.054	0.008	0.019	0.023	0.11		
mnth_2	-0.02	0.057	0.0035	-0.3	-0.3	-0.13	0.13	-0.27	-0.15	-0.16	-0.15	1	-0.088	0.078	0.081	0.076	0.077	-0.087	0.079	0.082	-0.082	-0.081	0.028	0.017	0.009	0.025	2e-17	0.033	0.013	1e-17		
mnth_3	-0.005	0.053	0.006	-0.18	-0.18	-0.11	0.14	-0.14	0.088	-0.2	-0.19	-0.088	1	-0.096	-0.1	-0.094	0.095	-0.11	0.097	-0.1	-0.1	-0.1	-0.032	-0.012	0.034	0.016	0.031	0.01	-0.002	0.021		
mnth_4	-0.005	0.009	0.005	-0.051	-0.04	-0.13	0.18	0.023	0.51	-0.17	-0.17	-0.078	0.096	1	-0.089	0.033	0.084	0.095	0.086	-0.09	-0.09	-0.089	0.016	0.018	0.008	0.027	-0.039	0.003	0.011	0.0076		
mnth_5	-0.026	0.049	0.022	0.16	0.16	0.15	-0.034	0.11	0.53	-0.18	-0.17	-0.081	-0.1	-0.089	1	-0.086	0.087	-0.099	-0.09	-0.093	0.093	0.092	0.052	0.11	0.024	0.057	0.006	0.009	0.02	-0.053		
mnth_6	-0.01	-0.046	0.032	0.3	0.3	-0.095	0.035	0.19	0.25	0.068	-0.16	-0.076	0.094	0.083	0.086	1	-0.082	-0.093	0.084	0.087	0.087	0.086	0.023	0.003	0.039	0.047	0.009	0.051	-0.02	-0.095	-0.049	
mnth_7	-0.018	-0.047	0.012	0.41	0.41	-0.054	0.085	0.16	-0.16	0.48	-0.17	-0.077	0.095	0.084	0.087	0.082	1	-0.094	0.085	0.088	0.088	0.087	0.02	0.042	-0.049	0.03	-0.078	0.001	-0.11	0.0064		
mnth_8	-0.005	-0.053	0.013	0.39	0.36	0.053	-0.11	0.23	-0.19	0.55	-0.19	-0.087	-0.11	-0.095	0.099	0.093	0.094	1	-0.096	-0.1	-0.1	-0.099	0.075	0.063	0.072	0.019	0.003	0.033	0.017	-0.057		
mnth_9	-0.002	0.044	-0.035	0.21	0.21	0.15	-0.12	0.2	-0.17	0.35	-0.02	-0.079	0.097	0.086	-0.09	-0.084	0.085	-0.096	1	-0.091	-0.091	-0.09	-0.027	0.005	0.04	0.12	0.0027	0.02	0.015	0.045	0.034	
mnth_10	-0.013	0.005	-0.026	-0.013	0.005	0.16	-0.019	0.066	-0.18	-0.18	0.53	0.082	-0.1	-0.09	-0.093	0.087	0.088	-0.1	-0.099	1	-0.094	0.093	0.023	0.025	0.062	-0.026	0.049	0.014	0.013	0.15		
mnth_11	-0.019	0.13	0.026	-0.19	-0.19	0.049	0.01	-0.048	-0.18	-0.18	0.53	0.082	-0.1	-0.09	-0.093	0.087	0.088	-0.1	-0.091	0.094	1	-0.093	0.036	0.023	0.009	0.059	0.03	0.006	0.016	0.029		
mnth_12	-0.04	0.043	0.025	-0.27	-0.26	0.058	-0.061	-0.14	-0.17	-0.18	0.31	-0.081	-0.1	-0.089	0.092	0.086	0.087	-0.099	-0.09	-0.093	0.093	1	-0.026	-0.01	0.003	0.084	0.033	-0.012	0.019	-0.011		
weekday_1	-0.034	0.069	0.27	-0.038	0.037	-0.03	0.035	-0.071	-0.017	0.009	0.031	-0.028	0.032	0.016	-0.052	0.023	0.02	0.007	0.05	0.027	0.023	0.003	0.036	0.026	1	-0.18	-0.17	-0.17	-0.18	-0.17	-0.035	0.022
weekday_2	-0.026	0.035	0.23	-0.002	0.002	0.034	0.036	-0.019	0.001	0.01	0.031	0.017	-0.012	0.018	-0.01	0.003	0.042	-0.065	0.005	0.025	0.025	-0.011	-0.18	1	-0.17	-0.17	-0.18	-0.16	0.008	0.042		
weekday_3	-0.011	0.19	0.16	-0.046	0.043	0.026	-0.058	-0.02	-0.012	0.041	0.051	0.009	0.02	0.034	0.008	0.024	-0.047	0.049	0.072	0.012	0.062	0.009	0.038	-0.17	-0.17	1	-0.16	-0.16	-0.15	0.003	0.01	
weekday_4	-0.033	-0.066	0.26	0.075	0.079	-0.057	0.022	0.11	0.066	0.036	-0.054	0.025	0.016	0.027	0.057	0.009	0.04	0.03	0.019	0.027	0.02	0.009	0.084	-0.17	-0.17	-0.16	1	-0.17	-0.16	-0.06	-0.071	
weekday_5	-0.034	0.032	0.25	-0.02	-0.019	0.052	0.002	0.041	-0.018	0.008	0.008	0.02	-0.031	-0.039	0.006	-0.051	-0.078	0.039	0.02	-0.049	0.03	0.033	-0.18	-0.18	-0.16	-0.17	1	-0.16	0.003	0.092		
weekday_6	-0.018	0.025	-0.61	0.025	0.012	0.002	0.016	0.012	0.002	0.022	-0.019	0.033	0.01	-0.003	0.009	-0.02	-0.001	0.033	0.015	-0.014	0.006	0.012	-0.17	-0.16	-0.15	-0.16	-0.16	1	0.029	0.037		
weathersit_2	-0.015	0.038	-0.069	-0.09	-0.086	0.48	-0.03	-0.18	0.039	-0.075	0.023	-0.011	0.022	0.011	0.078	-0.095	-0.11	0.017	0.045	0.013	-0.016	0.019	0.039	0.088	0.009	-0.060	0.003	0.029	1	-0.13		
weathersit_3	-0.061	0.028	0.083	0.036	0.04	0.25	0.087	0.23	0.045	0.025	0.113	1e-17	0.021	0.007	0.053	0.049	0.006	-0.057	0.034	0.15	0.029	0.011	0.022	-0.042	0.001	-0.071	0.092	0.037	-0.13	1		
yr	1	0.015	-0.003	0.11	0.1	-0.083	0.001	0.59	0.014	0.044	-0.023	-0.020	0.005	0.054	-0.026	0.01	0.018	0.055	0.002	0.060	0.19	-0.04	-0.034	0.026	0.011	0.033	-0.034	0.018	0.015	-0.061		
holiday	-0.015	1	-0.25	-0.066	0.71	-0.029	0.018	-0.096	0.063	-0.04	0.051	0.057	-0.053	0.009	0.049	0.046	0.047	-0.053	0.044	0.005	0.43	-0.004	0.069	0.035	0.19	-0.066	-0.032	0.025	0.038	-0.028		
workingday	-0.003	-0.25	1	0.076	0.003	0.021	0.002	0.029	0.007	0.03	-0.033	0.036	0.003	0.008	0.085	0.022	0.032	-0.012	0.013	0.035	0.026	0.026	0.025	0.27	0.23	0.16	0.26	0.25	-0.61	-0.069	0.0083	
temp	-0.11	-0.066	0.007	1	0.99	0.16	-0.19	0.64	0.13	0.7	-0.23	-0.3	-0.18	0.051	0.16	0.3	0.41	0.39	0.21	-0.019	-0.19	-0.27	0.038	0.002	0.046	0.075	-0.02	0.025	-0.09	-0.036		
atemp	0.1	-0.073	0.003	0.99	1	0.17	-0.22	0.65	0.14	0.67	-0.21	-0.3	-0.18	-0.04	0.16	0.3	0.41	0.36	0.21	-0.005	0.19	-0.26	-0.037	0.002	0.043	0.079	-0.019	0.012	-0.086	0.044		
hum	-0.085	0.029	0.021	0.16	0.17	1	-0.27	-0.06	-0.013	0.04	0.17	-0.13	-0.11	-0.13	0.15	-0.095	0.054	0.053	0.15	0.16	0.004	0.058	-0.03	0.034	0.026	-0.057	0.052	0.002	0.48	0.25		
windspeed	-0.001	0.018	0.002	-0.19	-0.22	-0.27	1	-0.25	0.11	-0.19	-0.091	0.13	0.14	0.18	-0.034	0.035	0.085	-0.11	-0.12	-0.019	0.01	-0.061	0.035	0.036	-0.058	0.022	0.002	0.016	-0.03	0.087		
cnt	0.59	-0.090	0.008	0.64	0.65	-0.06	-0.25	1	0.13	0.37	0.033	-0.27	-0.14	0.023	0.11	0.19	0.16	0.23	0.2	0.066	-0.048	-0.14	-0.071	-0.019	-0.02	0.11	-0.041	0.012	-0.18	-0.23		
season_2	-0.014	-0.063	0.03	0.13	0.14	-0.013	0.11	0.13	1	-0.34	-0.33	-0.15	0.088	0.51	0.53	0.25	-0.16	-0.19	-0.17	-0.18	-0.18	-0.17	-0.017	0.001	0.12	-0.066	-0.018	0.002	0.039	-0.045		
season_3	-0.044	-0.04	-0.033	0.7	0.67	0.04	-0.19	0.37	0.34	1	-0.34	-0.16	-0.2	-0.17	-0.18	0.068	0.48	0.55	0.35	-0.18	-0.18	-0.18	0.009	0.031	-0.041	0.036	0.008	0.022	-0.075	-0.025		
season_4	-0.023	0.051	0.036	-0.23	-0.21	0.17	-0.091	0.033	-0.33	-0.34	1	-0.15	-0.19	-0.17	-0.17	-0.16	-0.17	-0.19	-0.02	0.53	0.53	0.31	0.017	0.045	0.051	-0.054	0.008	0.019	0.023	0.11		
mnth_2	-0.02	0.057	0.0035	-0.3	-0.3	-0.13	0.13	-0.27	-0.15	-0.16	-0.15	1	-0																			

Step 6 : Adding Variable and Features selection

→ Add Numerical Columns One by One:

- Added the 'temp' column as the first numerical feature.
- Evaluated the model's performance with the added feature using Ordinary Least Squares (OLS) regression.
- Obtained an R-squared value of 0.414, indicating that the 'temp' feature alone explains 41.4% of the variance in the target variable.
- Repeated the process by adding all numerical columns to the feature set.
- Evaluated the final model's performance, which achieved an impressive **R-squared value of 0.851**.
- The final model includes multiple numerical features along with categorical features like season, month, weekday, holiday, and weather conditions.

Step 7: Variance Inflation Factor (VIF) Checking

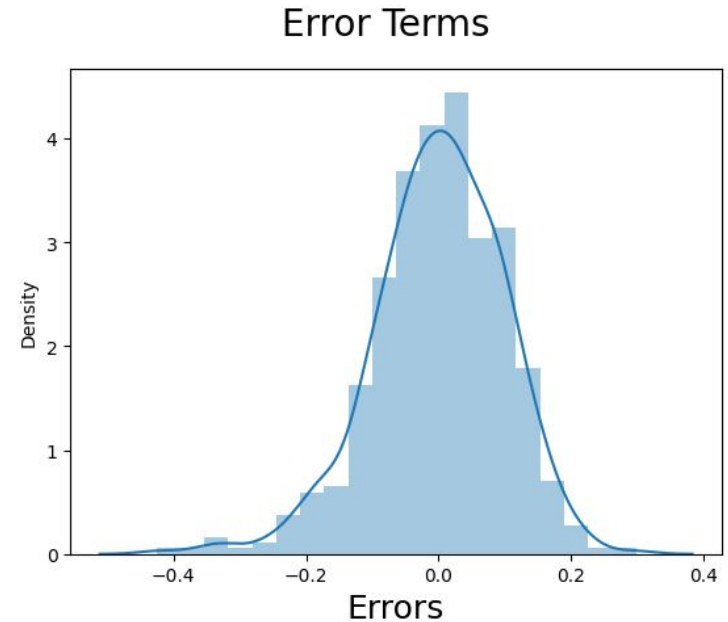
- The code you provided calculates the Variance Inflation Factor (VIF) for each feature variable in your dataset. The VIF helps assess multicollinearity among predictor variables in a regression analysis.
- **Variables with high VIF values, typically above 5**, indicate strong multicollinearity, suggesting that those variables are highly correlated with other predictor variables in the model.
- Based on the results of the VIF calculation, you should identify variables with high VIF values and consider dropping them from your model to address multicollinearity issues
- Look for variables with high VIF values, typically above 5, indicating multicollinearity.

Step 8: Dropping the variable and updating the model

- Look for variables with high VIF values, typically above 5, indicating multicollinearity.
- `['weekday_1', 'weekday_2', 'weekday_5', 'weekday_6', 'weekday_4', 'weekday_3', 'mnth_3', 'mnth_4', 'mnth_5', 'mnth_6', 'mnth_8', 'mnth_11', 'mnth_12', 'mnth_7', 'atemp', 'season_3', 'season_4', 'workingday']`
- Variables with p-values less than the chosen significance level (usually 0.05) are considered statistically significant.

Stop 9: Residual Analysis of the train data

- Scatter plot of residuals against predicted values.
- Histogram or density plot of residuals to check for normality.
- Q-Q plot to compare the distribution of residuals to a **normal distribution**.



Step 10: Making Predictions Using the Final Model

My final model exhibits a good R-squared value, indicating strong explanatory power, and features low p-values and VIF, suggesting that the selected variables are statistically significant and exhibit low multicollinearity, respectively.

Step 11: Model Evaluation

In the final step of model evaluation, my predictions align closely with the training sample data, indicating that the model performs well and is suitable for further feature development.

