

# Analysis of Categorical Variables

- **Question:** From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?
- **Seasonal Impact:** The highest number of bike bookings occurs during the Fall season, suggesting that cooler weather conditions are more favorable for biking.
- **Monthly Trends:** September consistently sees the highest number of bookings, followed by October and August. This indicates a peak in biking activities towards the end of summer and beginning of fall.
- **Trend Analysis:** From March to October, there is a noticeable upward trend in bookings, highlighting increased outdoor activity during these months.
- **Weather Influence:** Clear weather conditions significantly boost bike bookings, whereas rainy weather leads to a marked decrease in bookings.
- **Annual Growth:** The data shows a substantial increase in bookings in 2019 compared to the previous year, reflecting positive growth and popularity of the bike-sharing service.



Photo by Freddie Geng on Unsplash

# Importance of `drop\_first=True` in Dummy Variable Creation

- **Question:** Why is it important to use `drop_first=True` during dummy variable creation?
- **Preventing Multicollinearity:** Using `drop_first=True` helps to avoid multicollinearity by removing one category level, thereby preventing perfect collinearity among dummy variables.
- **Reducing Redundancy:** By dropping the first category, the total number of dummy variables is reduced, which simplifies the model and improves efficiency.
- **Model Interpretability:** This practice ensures that the model remains interpretable and free from redundant variables, making it easier to understand and analyze the effects of other predictors.

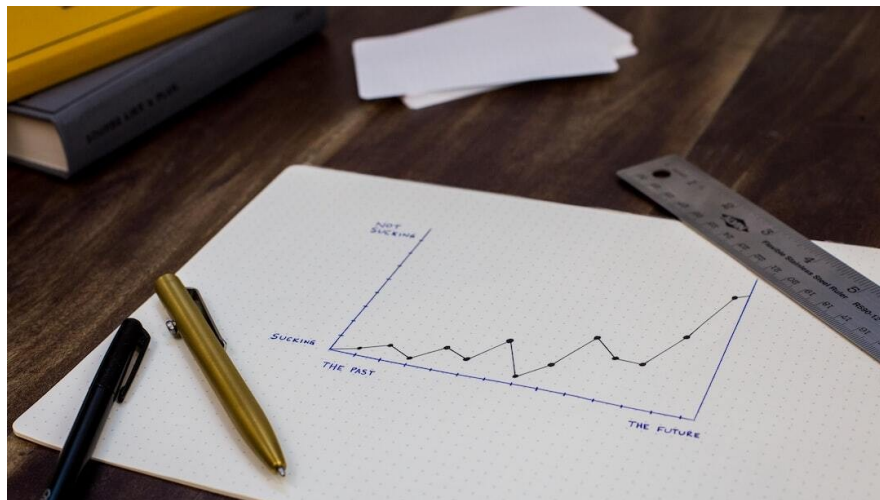


Photo by Isaac Smith on Unsplash

# Pair-Plot Analysis of Numerical Variables

- **Question:** Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?
- **Highest Correlation:** In the pair-plot analysis, the temperature variable shows the highest correlation with the target variable, with a correlation coefficient of 0.63. This strong positive correlation indicates that higher temperatures are associated with an increase in bike bookings.

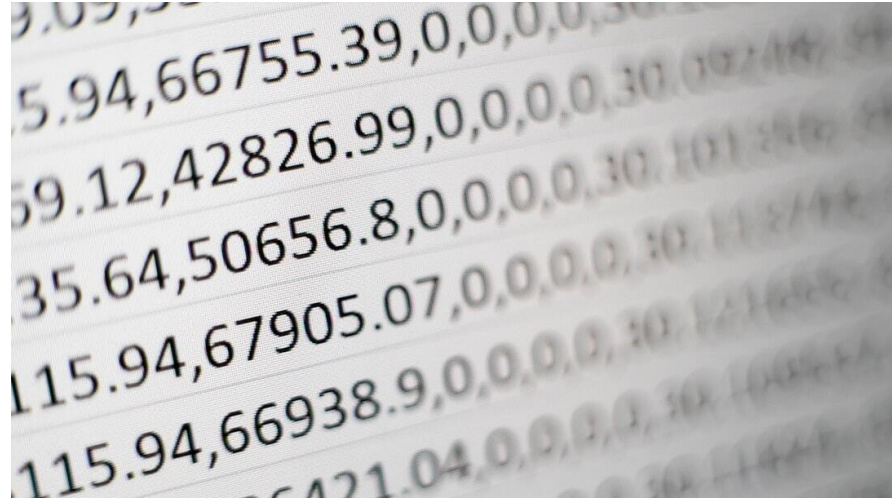


Photo by Mika Baumeister on Unsplash

# Validation of Linear Regression Assumptions

- **Question:** How did you validate the assumptions of Linear Regression after building the model on the training set?
- **Residual Analysis:** Plotting the residuals using a histogram or Q-Q plot to check for normality. If the residuals follow a normal distribution, the assumption is met.
- **Multicollinearity Check:** Calculating the Variance Inflation Factor (VIF) to ensure no high multicollinearity among the predictors.
- **Homoscedasticity Check:** Plotting residuals against the predicted values or each predictor to ensure there is no pattern, indicating constant variance of residuals (homoscedasticity).

# Top Features Contributing to Bike Demand

- **Question:** Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?
- **Temperature:** Higher temperatures correlate with more bike rentals.
- **Year:** The year variable indicates an increasing trend in bike usage over time.
- **Weather Conditions (Light Rain):** Light rain conditions have a measurable impact, suggesting users are more willing to rent bikes even in mild adverse weather.



Photo by Viktor Keri on Unsplash

# Explaining Linear Regression Algorithm

- **Question:** Explain the linear regression algorithm in detail.
- **Model Definition:** Linear regression models the relationship between a dependent variable (Y) and one or more independent variables (X) using a linear equation.
- **Least Squares Method:** The goal is to find the values of the coefficients that minimize the sum of squared errors (SSE) between observed and predicted values.
- **Algorithm Steps:** Includes data preparation, model building, model training, evaluation, and tuning.

# Anscombe's Quartet

- **Question:** Explain the Anscombe's quartet in detail.
- **Purpose:** Demonstrates the importance of data visualization in statistical analysis and the potential for outliers to distort statistical results.
- **Components:** Each dataset in the quartet has the same statistical properties but different distributions when graphed.
- **Conclusion:** Highlights that summary statistics alone can be misleading without considering the underlying data distribution.



# Pearson's R

- **Question:** What is Pearson's R?
- **Definition:** Pearson's correlation coefficient measures the strength and direction of a linear relationship between two continuous variables.
- **Range:** Ranges from -1 to 1, where 1 indicates a perfect positive linear relationship, 0 indicates no linear relationship, and -1 indicates a perfect negative linear relationship.
- **Calculation:** Calculated using the covariance of the variables divided by the product of their standard deviations.



Photo by SOULSANA on Unsplash



# Scaling in Machine Learning

- **Question:** What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?
- **Purpose of Scaling:** Scaling transforms numerical features to a common scale, which is crucial for algorithms sensitive to the scale of input features.
- **Normalized Scaling:** Also known as min-max scaling, it rescales features to a range of  $[0, 1]$  or  $[-1, 1]$ . Useful when the data is not normally distributed.
- **Standardized Scaling:** Transforms features to have a mean of 0 and a standard deviation of 1. Useful when the data needs to be normalized for statistical analysis and is less sensitive to outliers.



Photo by Emanuele Diviso on Unsplash

# Variance Inflation Factor (VIF)

- **Question:** You might have observed that sometimes the value of VIF is infinite. Why does this happen?
- **Perfect Multicollinearity:** Occurs when one variable is a perfect linear combination of others, leading to an undefined variance and infinite VIF.
- **Near Perfect Multicollinearity:** High correlations between independent variables can significantly inflate the VIF value, indicating multicollinearity issues.
- **Impact on Model:** High VIF values can lead to unstable coefficient estimates and inflated standard errors, compromising model reliability.

# Q-Q Plot in Linear Regression

- **Question:** What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.
- **Purpose:** Q-Q plot compares quantiles of observed data against a theoretical distribution, typically the normal distribution.
- **Normality Assumption:** Used to check if residuals follow a normal distribution, a key assumption in linear regression.
- **Outlier Detection:** Identifies outliers or deviations from normality, guiding model refinement and validation.

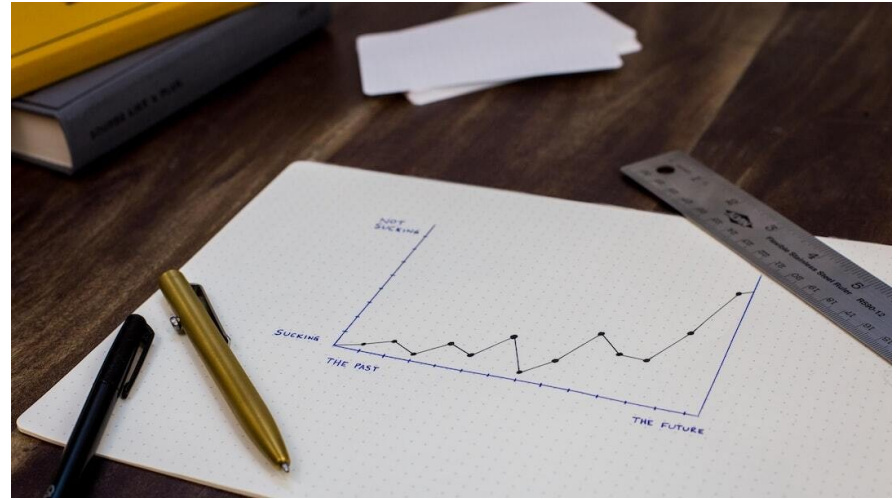


Photo by Isaac Smith on Unsplash

# Conclusion

- **Key Learnings:** Summarized the analysis of categorical variables, importance of dummy variable creation, pair-plot analysis, validation of linear regression assumptions, and identified top features contributing to bike demand.
- **Statistical Insights:** Explored Pearson's correlation coefficient, scaling techniques, and the importance of VIF and Q-Q plots in linear regression.
- **Practical Applications:** Demonstrated the importance of data visualization and thorough statistical analysis in understanding and predicting bike-sharing demand.