# ASSIGNMENT-BASED SUBJECTIVE QUESTIONS

BY: SACHIN KUMAR SINGH

**1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?**

Ans => **The demand of bike is less in the month of spring when compared with other seasons. The demand bike increased in the year 2019 when compared with year 2018.**

**2. Why is it important to use drop_first=True during dummy variable creation?**

Ans=> **If we do not use drop_first = True, then n dummy variables will be created, and these predictors(n dummy variables) are themselves correlated which is known as multicollinearity and it, in turn, leads to Dummy Variable Trap.**

**3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?**

Ans=>After checking atemp and temp both have same correlation with target variable of 0.63 which is the highest among all numerical variables.

**4. How did you validate the assumptions of Linear Regression after building the model on the training set?**

Ans=> According to this assumption there is linear relationship between the features and target. Linear regression captures only linear relationship. This can be validated by plotting a scatter plot between the features and the target.

- **5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?**

- Ans=> **Based on final model top three features contributing significantly towards explaining the demand are:**

  - **Temperature (0.73125)**

  - **weathersit : Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds (--0.27750)**

  - **year (0.24235)**

# GENERAL SUBJECTIVE QUESTIONS

**1. Explain the linear regression algorithm in detail.**

**Ans=> Linear regression is a supervised machine learning method that is used by the Train Using AutoML tool and finds a linear equation that best describes the correlation of the explanatory variables with the dependent variable. This is achieved by fitting a line to the data using least squares. The line tries to minimize the sum of the squares of the residuals. The residual is the distance between the line and the actual value of the explanatory variable. ng a line to the data using least squares.**

**2. Explain the Anscombe's quartet in detail.**

Ans=> Anscombe's quartet comprises a set of four dataset, having identical descriptive statistical properties in terms of means, variance, R-Squared, correlations, and linear regression lines but having different representations when we scatter plot on graph. Anscombe's quartet comprises a set of four dataset, having identical descriptive statistical properties in terms of means, variance, R-Squared, correlations, and linear regression lines but having different representations when we scatter plot on graph.

**3. What is Pearson's R?**

**Ans=>** The Pearson correlation coefficient (r) is the most common way of measuring a linear correlation. It is a number between –1 and 1 that measures the strength and direction of the relationship between two variables.

**4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?**

Ans=>*It is a step of data Pre-Processing which is applied to independent variables to normalize the data within a particular range. It also helps in speeding up the calculations in an algorithm.*

*Most of the times, collected data set contains features highly varying in magnitudes, units and range. If scaling is not done then algorithm only takes magnitude in account and not units hence incorrect modelling. To solve this issue, we have to do scaling to bring all the variables to the same level of magnitude.*
*It is important to note that scaling just affects the coefficients and none of the other parameters like t-statistic, F-statistic, p-values, R-squared, etc.*

*Normalization brings all of the data in the range of 0 and 1, sklearn .preprocessing .MinMaxScaler helps to implement normalization in python.*

*Standardization replaces the values by their Z scores. It brings all of the data into a standard normal distribution which has mean ($\mu$) zero and standard deviation one ($\sigma$).*

**5. You might have observed that sometimes the value of VIF is infinite. Why does this happen?**

**Ans=> If all the independent variables are orthogonal to each other, then VIF = 1.0. If there is perfect correlation, then VIF = infinity. A large value of VIF indicates that there is a correlation between the variables.**

**6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.**

**Ans=> Q-Q plots are also known as Quantile-Quantile plots. As the name suggests, they plot the quantiles of a sample distribution against quantiles of a theoretical distribution. Doing this helps us determine if a dataset follows any particular type of probability distribution like normal, uniform, exponential.**

# THANK YOU