# Assignment-based Subjective Questions

## (Submitted by chirag pallan)

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

Sol. 1

- Count of bike rentals increased and became popular in year 2019 than 2018 (from 'year' variable)

- Fount of bike rentals is more during clear weather (from 'Weathersit' variable)

- Fall and summer are more favorable for bike rentals than spring (from 'Season' variable)

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

Sol. 2

- To avoid multicollinearity (if we don't drop, dummy variables will be correlated) and affects the model adversely

- To avoid redundant features

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

Sol. 3

- Count (target variable) has significantly high correlation with temperature (temp)

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

Sol. 4

- Residual errors follow normal distribution

- Maintains linear relation between dependent variable (test and predicted)

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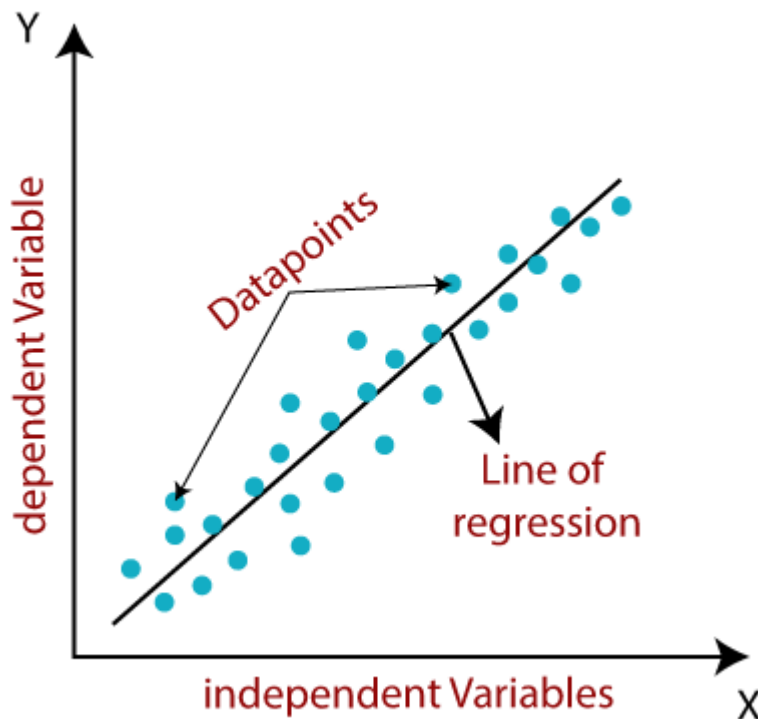
Sol. 5

-  TEMPERATURE (0.4354)

- WEATHER SITUATION – LIGHT AND SNOWY  (0.2837)

- YEAR  (0.2461)

# General Subjective Questions

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

Sol. 1

- Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables.

- Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (X) variables, hence called as linear regression.

- The linear regression model provides a sloped straight line representing the relationship between the variables. Consider the below image:



Mathematically, we can represent a linear regression as:

$y = a_0 + a_1 x + \varepsilon$

Where,

y= Dependent Variable (Target Variable)

X= Independent Variable (predictor Variable)

$a_0$= intercept of the line (Gives an additional degree of freedom)

$a_1$ = Linear regression coefficient (scale factor to each input value)

ε = random error

**Types of Linear Regression**

Linear regression can be further divided into two types of the algorithm:

> **- Simple Linear Regression**
>
> If a single independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Simple Linear Regression.
>
> **- Multiple Linear regression**
>
> If more than one independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Multiple Linear Regression.

**Finding the best fit line:**

- When working with linear regression, our main goal is to find the best fit line that means the error between predicted values and actual values should be minimized. The best fit line will have the least error.
- - The different values for weights or the coefficient of lines (a0, a1) gives a different line of regression, so we need to calculate the best values for a0 and a1 to find the best fit line, so to calculate this we use cost function.

**Cost function**

- - The different values for weights or coefficient of lines (a0, a1) gives the different line of regression, and the cost function is used to estimate the values of the coefficient for the best fit line.
- - Cost function optimizes the regression coefficients or weights. It measures how a linear regression model is performing.
- - We can use the cost function to find the accuracy of the mapping function, which maps the input variable to the output variable. This mapping function is also known as Hypothesis function.
- - For Linear Regression, we use the Mean Squared Error (MSE) cost function, which is the average of squared error occurred between the predicted values.

**Residuals**

- - The distance between the actual value and predicted values is called residual. If the observed points are far from the regression line, then the residual will be high, and so cost function will high. If the scatter points are close to the regression line, then the residual will be small and hence the cost function.

**Gradient Descent**

- Gradient descent is used to minimize the MSE by calculating the gradient of the cost function.
- A regression model uses gradient descent to update the coefficients of the line by reducing the cost function.
- It is done by a random selection of values of coefficient and then iteratively update the values to reach the minimum cost function.

**Model Performance**

- The Goodness of fit determines how the line of regression fits the set of observations. The process of finding the best model out of various models is called optimization. It can be achieved by below method

**R-squared method**

- R-squared is a statistical method that determines the goodness of fit.

- It measures the strength of the relationship between the dependent and independent variables on a scale of 0-100%.

- The high value of R-square determines the less difference between the predicted values and actual values and hence represents a good model.

- It is also called a coefficient of determination, or coefficient of multiple determination for multiple regression.

- It can be calculated from the below formula:

$$R\text{-squared} = \frac{\text{Explained variation}}{\text{Total Variation}}$$

**Assumptions of Linear Regression**

Below are some important assumptions of Linear Regression. These are some formal checks while building a Linear Regression model, which ensures to get the best possible result from the given dataset.

- **Linear relationship between the features and target**

  Linear regression assumes the linear relationship between the dependent and independent variables.

- **Small or no multicollinearity between the features**

  Multicollinearity means high-correlation between the independent variables. Due to multicollinearity, it may difficult to find the true relationship between the predictors and target variables. Or we can say, it is difficult to determine which predictor variable is affecting the target variable and which is not. So, the model assumes either little or no multicollinearity between the features or independent variables.

- **Homoscedasticity Assumption**

  Homoscedasticity is a situation when the error term is the same for all the values of independent variables. With homoscedasticity, there should be no clear pattern distribution of data in the scatter plot.

- **Normal distribution of error terms**

  Linear regression assumes that the error term should follow the normal distribution pattern. If error terms are not normally distributed, then confidence intervals will become either too wide or too narrow, which may cause difficulties in finding coefficients. It can be checked using the **q-q plot**. If the plot shows a straight line without any deviation, which means the error is normally distributed.

- **No autocorrelations**

    The linear regression model assumes no autocorrelation in error terms. If there will be any correlation in the error term, then it will drastically reduce the accuracy of the model. Autocorrelation usually occurs if there is a dependency between residual errors.

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

Sol. 2

- Anscombe's quartet comprises four data sets that have nearly identical simple descriptive statistics, yet have very different distributions and appear very different when graphed. Each dataset consists of eleven (x,y) points. They were constructed in 1973 by the statistician Francis
- Anscombe to demonstrate both the importance of graphing data when analyzing it, and the effect of outliers and other influential observations on statistical properties. He described the article as being intended to counter the impression among statisticians that "numerical calculations are exact, but graphs are rough.

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

Sol. 3

- Correlation coefficient formulas are used to find how strong a relationship is between data. The formulas return a value between -1 and 1, where, 1 indicates a strong positive relationship, -1 indicates a strong negative relationship, A result of zero indicates no relationship at all.

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

Sol. 4

- Scaling of data is method of making scale of the data in some particular range, there are various scaling methods like standard scalar, min-max scalar etc.
- Scaling is performed to avoid the effect of large values of coefficients.

| S.NO. | Normalization | Standardization |
|-------|---------------|-----------------|
| 1. | Minimum and maximum value of features are used for scaling | Mean and standard deviation is used for scaling. |
| 2. | It is used when features are of different scales. | It is used when we want to ensure zero mean and unit standard deviation. |

| S.NO. | Normalization | Standardization |
|-------|---------------|-----------------|
| 3. | Scales values between [0, 1] or [-1, 1]. | It is not bounded to a certain range. |
| 4. | It is really affected by outliers. | It is much less affected by outliers. |
| 5. | Scikit-Learn provides a transformer called MinMaxScaler for Normalization. | Scikit-Learn provides a transformer called StandardScaler for standardization. |
| 6. | This transformation squishes the n-dimensional data into an n-dimensional unit hypercube. | It translates the data to the mean vector of original data to the origin and squishes or expands. |
| 7. | It is useful when we don't know about the distribution | It is useful when the feature distribution is Normal or Gaussian. |
| 8. | It is a often called as Scaling Normalization | It is a often called as Z-Score Normalization. |

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

Sol. 5

- If there is perfect correlation, then VIF = infinity. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get R2 =1, which lead to 1/(1-R2) infinity. To solve this problem we need to drop one of the variables from the dataset which is causing this perfect multicollinearity.
- An infinite VIF value indicates that the corresponding variable may be expressed exactly by a linear combination of other variables (which show an infinite VIF as well).

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

Sol. 6

- Q-Q Plots (Quantile-Quantile plots) are plots of two quantiles against each other. A quantile is a fraction where certain values fall below that quantile. For example, the median is a quantile where 50% of the data fall below that point and 50% lie above it. The purpose of Q-Q plots is to find out if two sets of data come from the same distribution. A 45 degree angle is plotted on the Q-Q plot; if the two data sets come from a common distribution, the points will fall on that reference line.A Q Q plot showing the 45 degree reference line.
- If the two distributions being compared are similar, the points in the Q–Q plot will approximately lie on the line y = x. If the distributions are linearly related, the points in the Q–Q plot will approximately lie on a line, but not necessarily on the line y = x. Q–Q plots can also be used as a graphical means of estimating parameters in a location-scale family of distributions.
- A Q–Q plot is used to compare the shapes of distributions, providing a graphical view of how properties such as location, scale, and skewness are similar or different in the two distributions.