

Assignment-based Subjective Questions

1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (3 marks)

Ans – Following is the effect of the categorical variables on the dependent variable:-

1. There are more users in Summer and Fall season which can be seen from month data as well
2. Bikes are used more when the weather is clear and less in light Rain/snow
3. Users are almost same whether its weekday or weekend,
4. There is continuous increase in ridership from 2018 to 2019
5. With increase in temperature, there is linear increase in ridership
6. When humidity is optimal i.e. 50-80, more people tend to rent bike
7. Lower wind speed is directly related to more ridership

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

Ans- When we create dummy variable, it will create columns equal to the number of the categorical variables. i.e. if there are n variables, it will create n columns with value set as 1 if applicable else 0. If nth variable is equivalent to other n-1 variable set to 0.

Hence if we drop 1 variable, we will not lose any data and we can reduce the number of columns for our analysis.

3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (1 mark)

Ans- From the pair plot, highest correlation of target variable is with Temp

4. How did you validate the assumptions of Linear Regression after building the model on the training set? (3 marks)

Ans- Relationship is validated using following method-

- a. Multicollinearity- All the VIFs are below 5 which means there is relationship between variables are insignificant.
- b. Number of error terms- We can see from graph Error terms are normally distributed
- c. Linear relationship Validation- From the graphs we can see linear relationship between variables.
- d. Residual analysis(homoscedasticity)- There is no relationship between residuals

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (2 marks)

Ans- Final model equation is-

```
cnt = 0.275 + 0.248 * yr + 0.56 * workingday - 0.189 * windspeed - 0.087 * Pleasant - 0.302 * Light + 0.219 * summer + 0.258 * fall + 0.184 * winter - 0.103 * jan + 0.072 * sep + 0.065 * sat
```

3 most significant variables are-

- a. workingday – working day has most demand of the bikes
- b. fall – In the fall season, demand is more

- c. Light – It has most -ve impact on the demand. When the weather is Light Rain, light snow, thunderstorms, demand is decreased.

General Subjective Questions

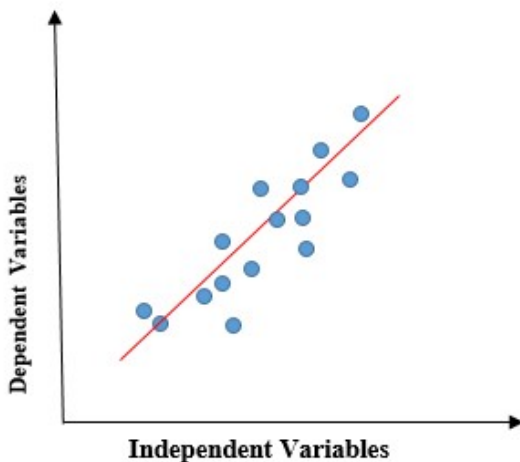
1. Explain the linear regression algorithm in detail. (4 marks)

Ans- Regression is a supervised learning technique that supports finding the correlation among variables. A regression problem is when the output variable is a real or continuous value.

In Regression, we plot a graph between the variables which best fit the given data points. The machine learning model can deliver predictions regarding the data. In naïve words, "Regression shows a line or curve that passes through all the data points on a target-predictor graph in such a way that the vertical distance between the data points and the regression line is minimum." It is used principally for prediction, forecasting, time series modeling, and determining the causal-effect relationship between variables.

Linear Regression

Linear regression is a quiet and simple statistical regression method used for predictive analysis and shows the relationship between the continuous variables. Linear regression shows the linear relationship between the independent variable (X-axis) and the dependent variable (Y-axis), consequently called linear regression. If there is a single input variable (x), such linear regression is called simple linear regression. And if there is more than one input variable, such linear regression is called multiple linear regression. The linear regression model gives a sloped straight line describing the relationship within the variables.



The above graph presents the linear relationship between the dependent variable and independent variables. When the value of x (**independent variable**) increases, the value of y (**dependent variable**) is likewise increasing. The red line is referred to as the best fit straight line. Based on the given data points, we try to plot a line that models the points the best.

To calculate best-fit line linear regression uses a traditional slope-intercept form.

$$y = mx + b \implies y = a_0 + a_1x$$

y= Dependent Variable.

x= Independent Variable.

a₀= intercept of the line.

a₁ = Linear regression coefficient.

Need of a Linear regression

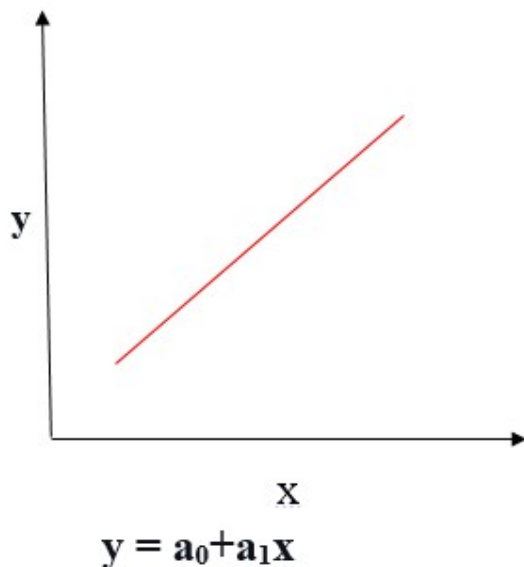
As mentioned above, Linear regression estimates the relationship between a dependent variable and an independent variable. Let's understand this with an easy example:

Let's say we want to estimate the salary of an employee based on year of experience. You have the recent company data, which indicates that the relationship between experience and salary. Here year of experience is an independent variable, and the salary of an employee is a dependent variable, as the salary of an employee is dependent on the experience of an employee. Using this insight, we can predict the future salary of the employee based on current & past information.

A regression line can be a Positive Linear Relationship or a Negative Linear Relationship.

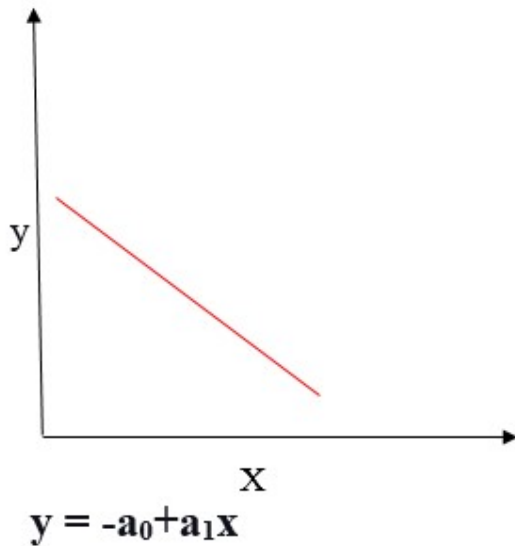
Positive Linear Relationship

If the dependent variable expands on the Y-axis and the independent variable progress on X-axis, then such a relationship is termed a Positive linear relationship.



Negative Linear Relationship

If the dependent variable decreases on the Y-axis and the independent variable increases on the X-axis, such a relationship is called a negative linear relationship.



The goal of the linear regression algorithm is to get the best values for a_0 and a_1 to find the best fit line. The best fit line should have the least error means the error between predicted values and actual values should be minimized.

Cost function

The cost function helps to figure out the best possible values for a_0 and a_1 , which provides the best fit line for the data points.

Cost function optimizes the regression coefficients or weights and measures how a linear regression model is performing. The cost function is used to find the accuracy of the **mapping function** that maps the input variable to the output variable. This mapping function is also known as **the Hypothesis function**.

In Linear Regression, **Mean Squared Error (MSE)** cost function is used, which is the average of squared error that occurred between the predicted values and actual values.

By simple linear equation $y=mx+b$ we can calculate MSE as:

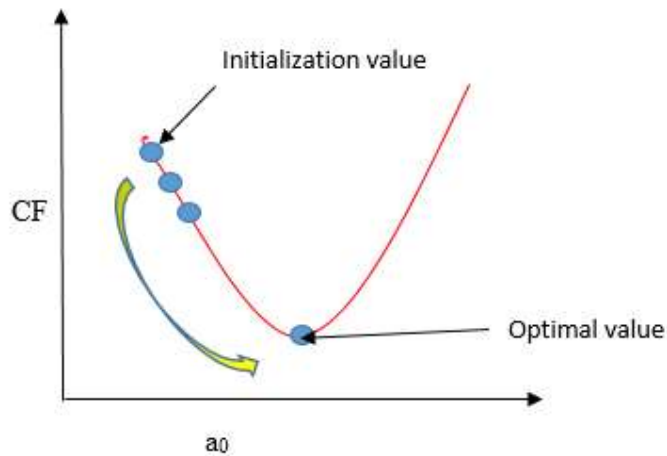
Let's y = actual values, y_i = predicted values

$$MSE = \frac{1}{N} \sum_{i=1}^n (y_i - (mx_i + b))^2$$

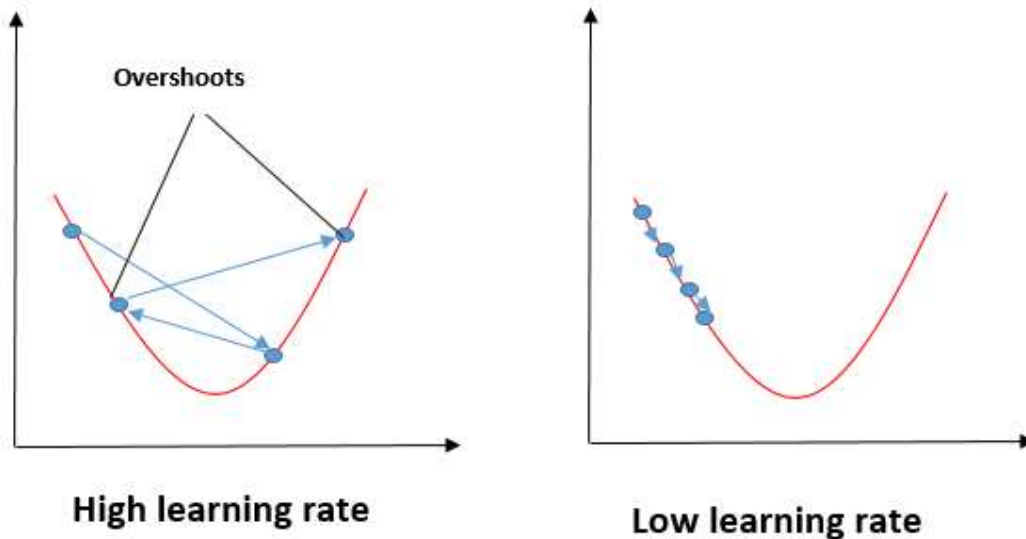
Using the MSE function, we will change the values of a_0 and a_1 such that the MSE value settles at the minima. Model parameters **x_i , b (a_0, a_1)** can be manipulated to minimize the cost function. These parameters can be determined using the gradient descent method so that the cost function value is minimum.

Gradient descent

Gradient descent is a method of updating a_0 and a_1 to minimize the cost function (MSE). A regression model uses gradient descent to update the coefficients of the line ($a_0, a_1 \Rightarrow x_i, b$) by reducing the cost function by a random selection of coefficient values and then iteratively update the values to reach the minimum cost function.



Imagine a pit in the shape of U. You are standing at the topmost point in the pit, and your objective is to reach the bottom of the pit. There is a treasure, and you can only take a discrete number of steps to reach the bottom. If you decide to take one footstep at a time, you would eventually get to the bottom of the pit but, this would take a longer time. If you choose to take longer steps each time, you may get to sooner but, there is a chance that you could overshoot the bottom of the pit and not near the bottom. In the gradient descent algorithm, the number of steps you take is the learning rate, and this decides how fast the algorithm converges to the minima.



To update a_0 and a_1 , we take gradients from the cost function. To find these gradients, we take partial derivatives for a_0 and a_1 .

$$J = \frac{1}{n} \sum_{i=1}^n (a_0 + a_1 \cdot x_i - y_i)^2$$

$$\frac{\partial J}{\partial a_0} = \frac{2}{n} \sum_{i=1}^n (a_0 + a_1 \cdot x_i - y_i)$$

$$\frac{\partial J}{\partial a_1} = \frac{2}{n} \sum_{i=1}^n (a_0 + a_1 \cdot x_i - y_i) \cdot x_i$$

$$\frac{\partial J}{\partial a_0} = \frac{2}{n} \sum_{i=1}^n (pred_i - y_i)$$

$$\frac{\partial J}{\partial a_1} = \frac{2}{n} \sum_{i=1}^n (pred_i - y_i) \cdot x_i$$

$$a_0 = a_0 - \alpha \cdot \frac{2}{n} \sum_{i=1}^n (pred_i - y_i)$$

$$a_1 = a_1 - \alpha \cdot \frac{2}{n} \sum_{i=1}^n (pred_i - y_i) \cdot x_i$$

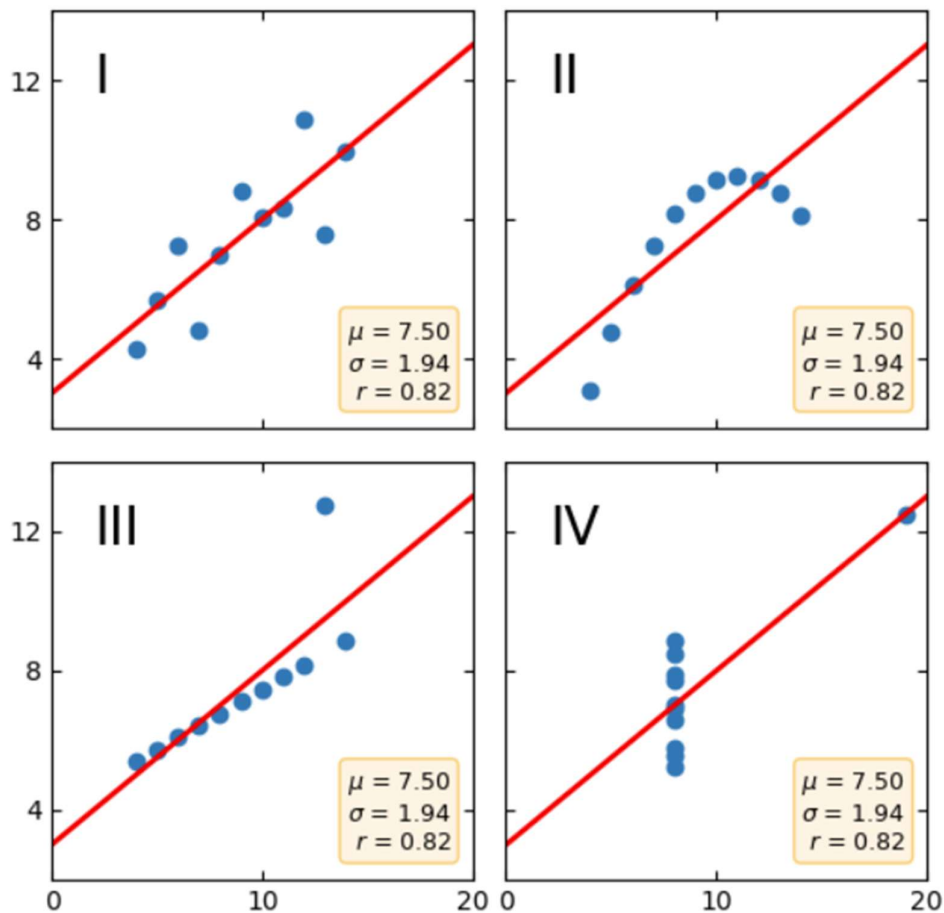
partial derivatives are the gradients and they are used to update the

The partial derivatives are the gradients, and they are used to update the values of a_0 and a_1 . Alpha is the learning rate.

2. Explain the Anscombe's quartet in detail. (3 marks)

Ans- 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.

Anscombe's quartet tells us about the importance of visualizing data before applying various algorithms to build models. This suggests the data features must be plotted to see the distribution of the samples that can help you identify the various anomalies present in the data (outliers, diversity of the data, linear separability of the data, etc.). Moreover, the linear regression can only be considered a fit for the data with linear relationships and is incapable of handling any other kind of data set.



We can describe the four data sets as:

Data Set 1: fits the linear regression model pretty well.

Data Set 2: cannot fit the linear regression model because the data is non-linear.

Data Set 3: shows the outliers involved in the data set, which cannot be handled by the linear regression model.

Data Set 4: shows the outliers involved in the data set, which also cannot be handled by the linear regression model.

Anscombe's quartet helps us to understand the importance of data visualization and how easy it is to fool a regression algorithm. So, before attempting to interpret and model the data or implement any machine learning algorithm, we first need to visualize the data set in order to help build a well-fit model.

3. What is Pearson's R? (3 marks)

Ans- The Pearson correlation coefficient is a descriptive statistic, meaning that it summarizes the characteristics of a dataset. Specifically, it describes the strength and direction of the linear relationship between two quantitative variables.

Although interpretations of the relationship strength (also known as effect size) vary between disciplines, the table below gives general rules of thumb:

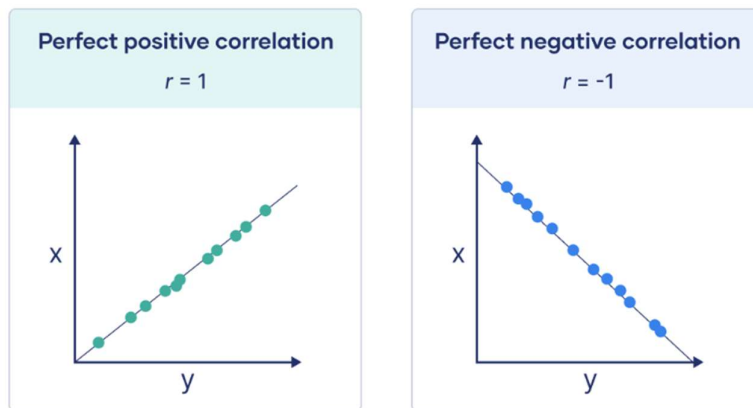
Pearson correlation coefficient (r) value	Strength	Direction
Greater than .5	Strong	Positive
Between .3 and .5	Moderate	Positive
Between 0 and .3	Weak	Positive
0	None	None
Between 0 and $-.3$	Weak	Negative
Between $-.3$ and $-.5$	Moderate	Negative
Less than $-.5$	Strong	Negative

Visualizing the Pearson correlation coefficient:-

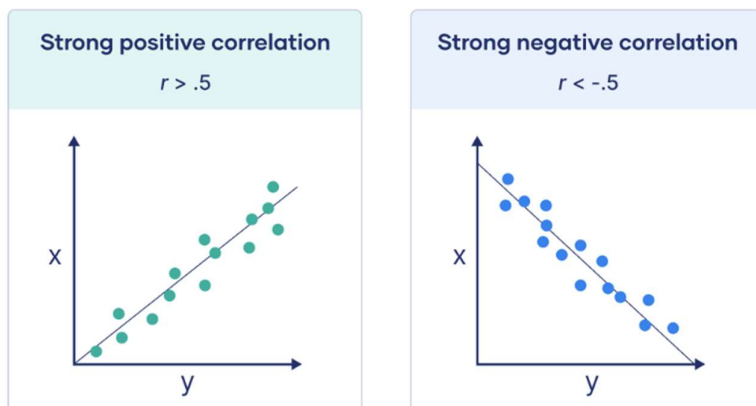
Another way to think of the Pearson correlation coefficient (r) is as a measure of how close the observations are to a line of best fit.

The Pearson correlation coefficient also tells you whether the slope of the line of best fit is negative or positive. When the slope is negative, r is negative. When the slope is positive, r is positive.

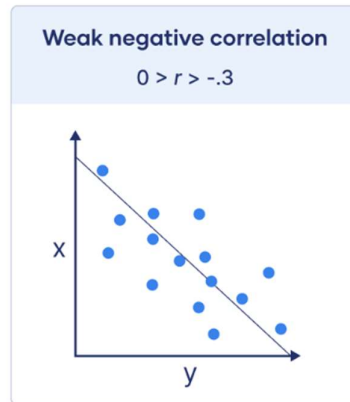
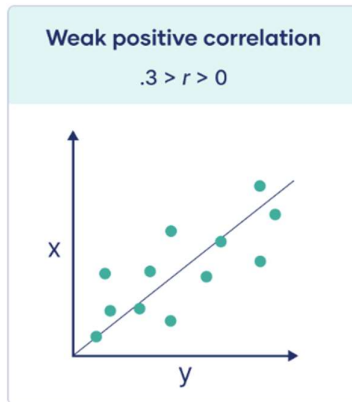
When r is 1 or -1 , all the points fall exactly on the line of best fit:



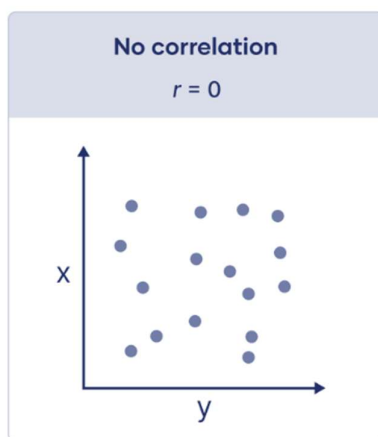
When r is greater than .5 or less than $-.5$, the points are close to the line of best fit:



When r is between 0 and .3 or between 0 and $-.3$, the points are far from the line of best fit:



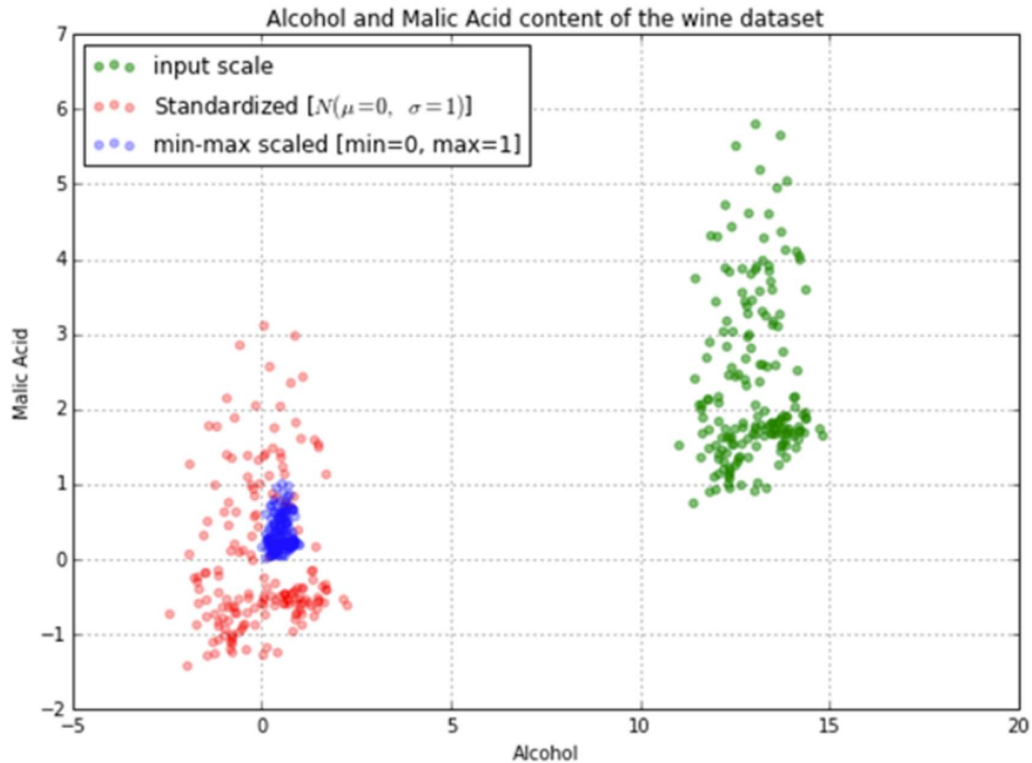
When r is 0, a line of best fit is not helpful in describing the relationship between the variables:



4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (3 marks)

Ans- Scaling is a method used to normalize the range of independent variables or features of data. In data processing, it is also known as data normalization and is generally performed during the data preprocessing step. Just to give you an example — if you have multiple independent variables like age, salary, and height; With their range as (18–100 Years), (25,000–75,000 Euros), and (1–2 Meters) respectively, feature scaling would help them all to be in the same range, for example- centered around 0 or in the range (0,1) depending on the scaling technique.

In the below figure we can see the impact of Standardization and Normalisation Scaling on the data set-



Normalization:-

Also known as min-max scaling or min-max normalization, it is the simplest method and consists of rescaling the range of features to scale the range in [0, 1]. The general formula for normalization is given as:

$$x' = a + \frac{(x - \min(x))(b - a)}{\max(x) - \min(x)}$$

Here, $\max(x)$ and $\min(x)$ are the maximum and the minimum values of the feature respectively.

We can also do a normalization over different intervals, e.g. choosing to have the variable laying in any $[a, b]$ interval, a and b being real numbers. To rescale a range between an arbitrary set of values $[a, b]$, the formula becomes:

Standardization:-

Feature standardization makes the values of each feature in the data have zero mean and unit variance. The general method of calculation is to determine the distribution mean and standard deviation for each feature and calculate the new data point by the following formula:

$$x' = \frac{x - \bar{x}}{\sigma}$$

Here, σ is the standard deviation of the feature vector, and \bar{x} is the average of the feature vector.

5. You might have observed that sometimes the value of VIF is infinite. Why does this happen? (3 marks)

Ans- VIF (Variance Inflation Factor) is defined as-

$$VIF = 1/(1 - R^2)$$

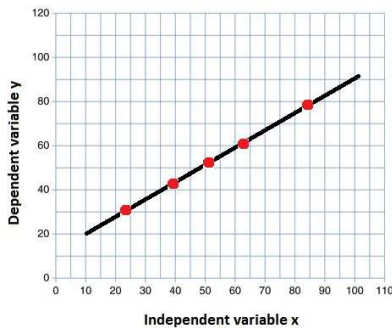
VIF will be infinite when R^2 tends to 1.

R squared or Coefficient of determination, or R^2 is a measure that provides information about the goodness of fit of the regression model. In simple terms, it is a statistical measure that tells how well the plotted regression line fits the actual data. R squared measures how much the variation is there in predicted and actual values in the regression model.

- R-squared values range from 0 to 1, usually expressed as a percentage from 0% to 100%.
- And this value of R square tells you how well the data fits the line you've drawn.
- The higher the model's R-Squared value, the better the regression line fits the data.

$R^2=1$, All the variation in the y values is accounted for by the x values.

Graph-



When R^2 is 1, it means that data is fitting perfectly. **This is simply the overfitting and must be avoided.**

6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (3 marks)

Ans- The quantile-quantile (q-q) plot is a graphical technique for determining if two data sets come from populations with a common distribution.

A q-q plot is a plot of the quantiles of the first data set against the quantiles of the second data set. By a quantile, we mean the fraction (or percent) of points below the given value. That is, the 0.3 (or 30%) quantile is the point at which 30% percent of the data fall below and 70% fall above that value.

A 45-degree reference line is also plotted. If the two sets come from a population with the same distribution, the points should fall approximately along this reference line. The greater the departure from this reference line, the greater the evidence for the conclusion that the two data sets have come from populations with different distributions.

The advantages of the q-q plot are:

- The sample sizes do not need to be equal.
- Many distributional aspects can be simultaneously tested. For example, shifts in location, shifts in scale, changes in symmetry, and the presence of outliers can all be detected from this plot. For example, if the two data sets come from populations whose distributions differ only by a shift in location, the points should lie along a straight line that is displaced either up or down from the 45-degree reference line.

Here is Q-Q plot from the bike assignment project-

```
In [79]: #Q-Q plot for count variable
import scipy.stats as stats
import pylab

stats.probplot(bikes['cnt'], dist="norm", plot=pylab)
pylab.show()
```

