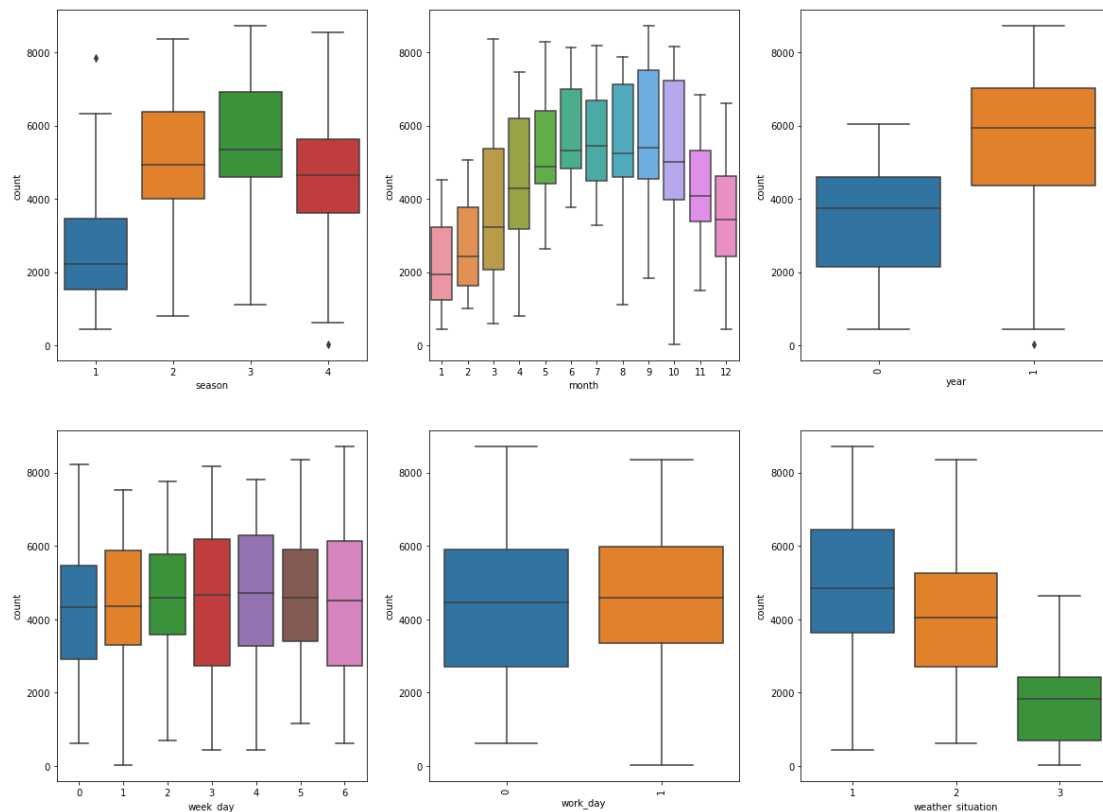


Assignment-based Subjective Questions

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:



There were 6 categorical variables in the dataset viz, season, month, year, week_day, work_day and weather_situation.

We used Box plot (refer the fig above) to study their effect on the dependent variable ('cnt').

Inference:

- **Season:** Almost 32% of the bike booking were happening in season3 with a median of over 5000 booking (for the period of 2 years). This was followed by season2 & season4 with 27% & 25% of total booking. This indicates, season can be a good predictor for the dependent variable.
- **Month:** Almost 10% of the bike booking were happening in the months 5,6,7,8 & 9 with a median of over 4000 booking per month. This indicates, month has some trend for bookings and can be a good predictor for the dependent variable.
- **Year:** For year 2019, count of users is higher than year 2018.
- **Week day:** The bike demand is almost even throughout the week.
- **Work day:** The median count of users is constant almost throughout the week.
- **Weather Situation:** Almost 67% of the bike booking were happening during clear weather with a median of close to 5000 booking (for the period of 2 years). This was followed by cloudy-

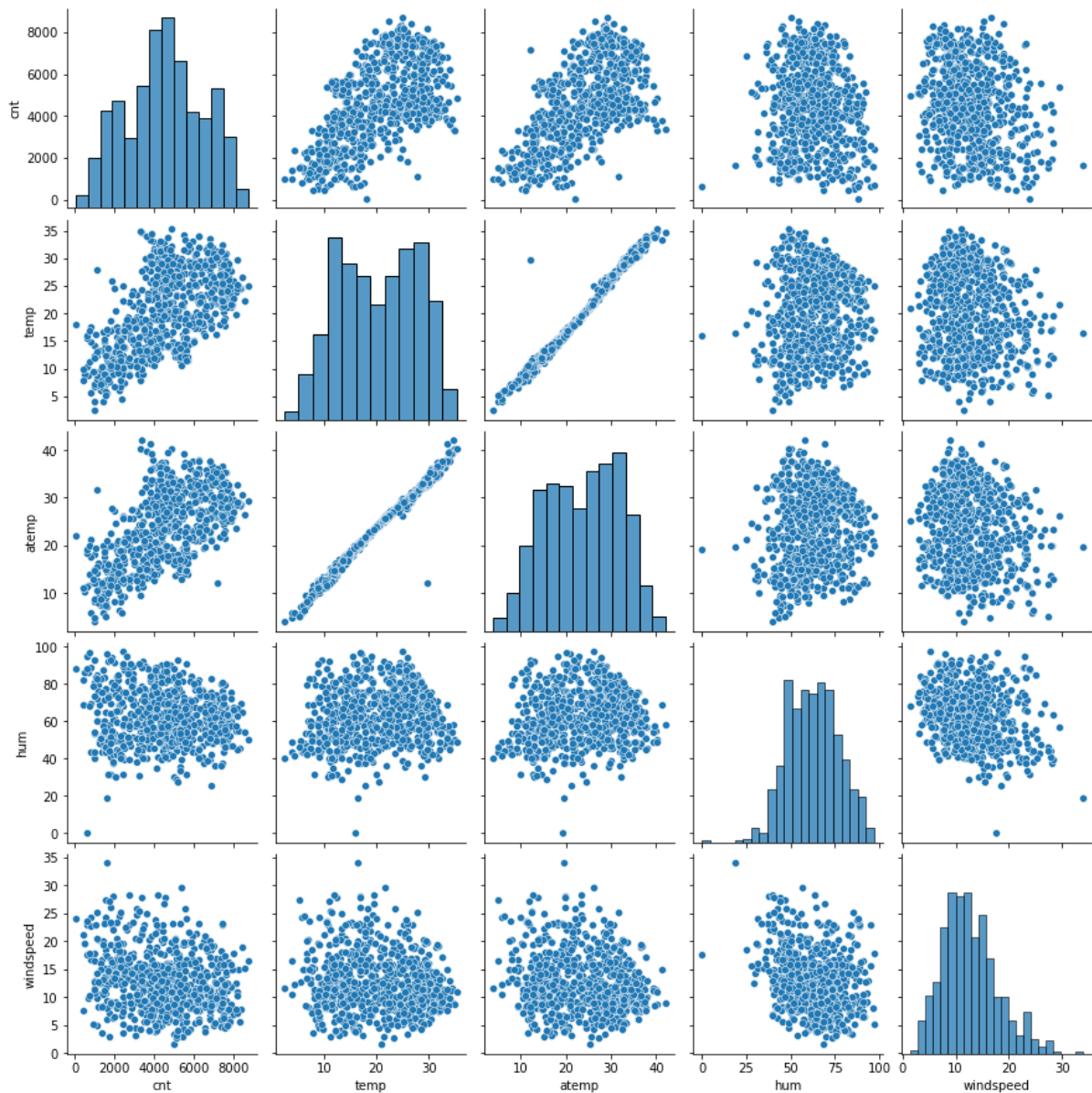
misty day with 30% of total booking. This indicates, weather situation does show some trend towards the bike bookings can be a good predictor for the dependent variable.

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

Ans: It is important to use drop_first=True because it helps in reducing the extra column created during dummy variable creation. Hence it reduces the correlations created among dummy variables. It gives smaller set of data for regression analysis. This will reduce the amount of time required by the algorithm to find the best fit line with the given data points.

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

Ans:

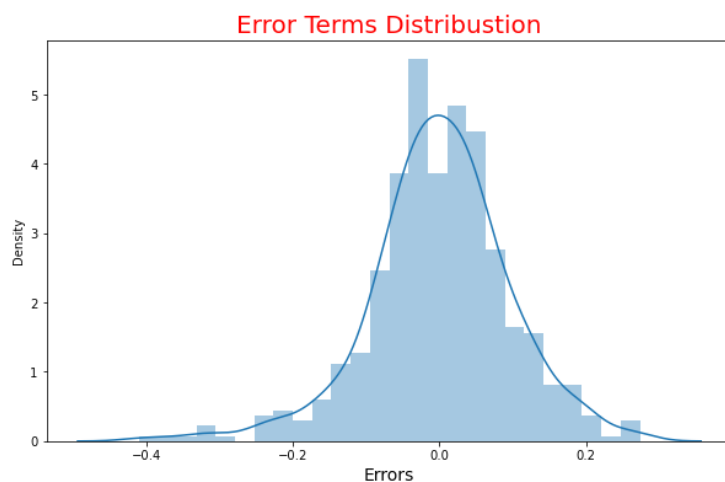


Using the below pairplot it can be seen that , “temp” and “atemp” are the two numerical variables which are highly correlated with the target variable (cnt).

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

Ans: The following tests were done to validate the assumptions of linear regression:

- First, linear regression needs the relationship between the independent and dependent variables to be linear. We visualised the numeric variables using a pairplot to see if the variables are linearly related or not. Refer to the notebook for more details.
- Secondly, Residuals distribution should follow normal distribution and centred around 0 (mean = 0). We validated this assumption about residuals by plotting a distplot of residuals and saw if residuals are following normal distribution or not. The diagram below shows that the residuals are distributed about mean = 0.
- Thirdly, linear regression assumes that there is little or no multicollinearity in the data. Multicollinearity occurs when the independent variables are too highly correlated with each other. We calculated the VIF (Variance Inflation Factor) to get the quantitative idea about how much the feature variables are correlated with each other in the new model. Refer to the notebook for more details.



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

Ans: Top 3 features are:

Feature	Coefficient
Temperature	0.6292
Year	0.2309
Wind speed	-0.1496

General Subjective Questions

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

Ans: Linear Regression is a type of supervised Machine Learning algorithm that is used for the prediction of numeric values. Linear Regression is the most basic form of regression analysis. Regression is the most commonly used predictive analysis model.

Linear regression is based on the popular equation “ $y = mx + c$ ”.

It assumes that there is a linear relationship between the dependent variable(y) and the predictor(s)/independent variable(x).

In regression, we calculate the best fit line which describes the relationship between the independent and dependent variable. Regression is performed when the dependent variable is of continuous data type and Predictors or independent variables could be of any data type like continuous, nominal/categorical etc. Regression method tries to find the best fit line which shows the relationship between the dependent variable and predictors with least error.

In regression, the output/dependent variable is the function of an independent variable and the coefficient and the error term.

Regression is broadly divided into simple linear regression and multiple linear regression.

1. Simple Linear Regression: SLR is used when the dependent variable is predicted using only one independent variable. The equation for SLR will be:

The diagram shows the equation $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ with the following labels and arrows:

- Dependent Variable** points to Y_i .
- Population Y intercept** points to β_0 .
- Population Slope Coefficient** points to β_1 .
- Independent Variable** points to X_i .
- Random Error term** points to ϵ_i .

Below the equation, two blue curly braces group the terms:

- A brace under $\beta_0 + \beta_1 X_i$ is labeled **Linear component**.
- A brace under ϵ_i is labeled **Random Error component**.

2. Multiple Linear Regression: MLR is used when the dependent variable is predicted using multiple independent variables. The equation for MLR will be:

$$\text{observed data} \rightarrow y = b_0 + b_1x_1 + b_2x_2 + \dots + b_px_p + \varepsilon$$

$$\text{predicted data} \rightarrow y' = b_0 + b_1x_1 + b_2x_2 + \dots + b_px_p$$

$$\text{error} \rightarrow \varepsilon = y - y'$$

B0 = intercept (constant term)

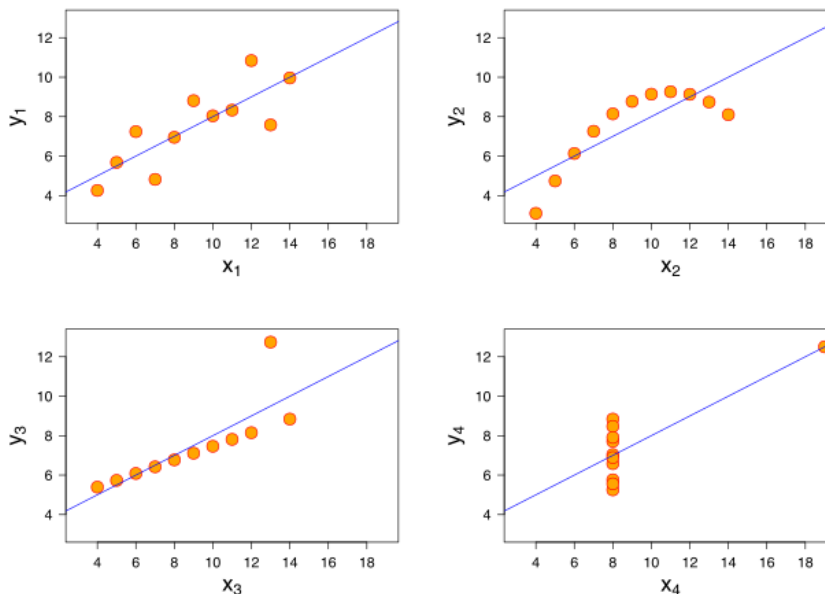
B1 = coefficient for X1 variable

B2 = coefficient for X2 variable

B3 = coefficient for X3 variable and so on...

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

Ans: Anscombe's Quartet was developed by statistician Francis Anscombe. It includes four data sets that have almost identical statistical features, but they have a very different distribution and look totally different when plotted on a graph. It was developed to emphasize both the importance of graphing data before analysing it and the effect of outliers and other influential observations on statistical properties.



- The first scatter plot (top left) appears to be a simple linear relationship, corresponding to two variables correlated where y could be modelled as gaussian with mean linearly dependent on x.

- The second graph (top right); while a relationship between the two variables is obvious, it is not linear, and the Pearson correlation coefficient is not relevant. A more general regression and the corresponding coefficient of determination would be more appropriate.
- In the third graph (bottom left), the modelled relationship is linear, but should have a different regression line (a robust regression would have been called for). The calculated regression is offset by the one outlier which exerts enough influence to lower the correlation coefficient from 1 to 0.816.
- Finally, the fourth graph (bottom right) shows an example when one high-leverage point is enough to produce a high correlation coefficient, even though the other data points do not indicate any relationship between the variables.

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

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.

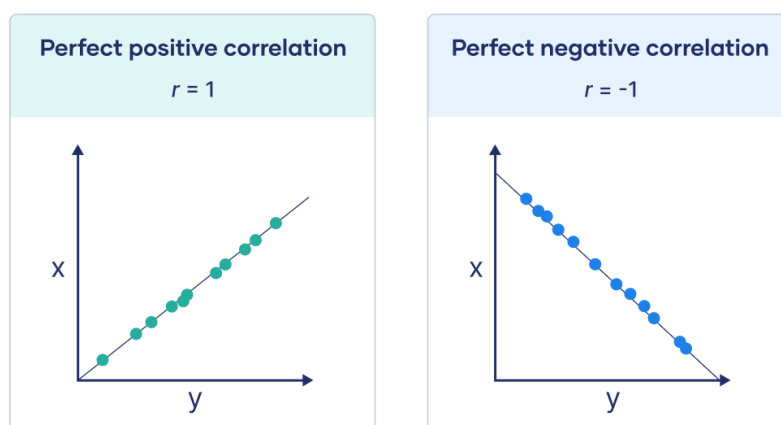
It is a measure of linear correlation between two sets of data.

Pearson correlation coefficient when:

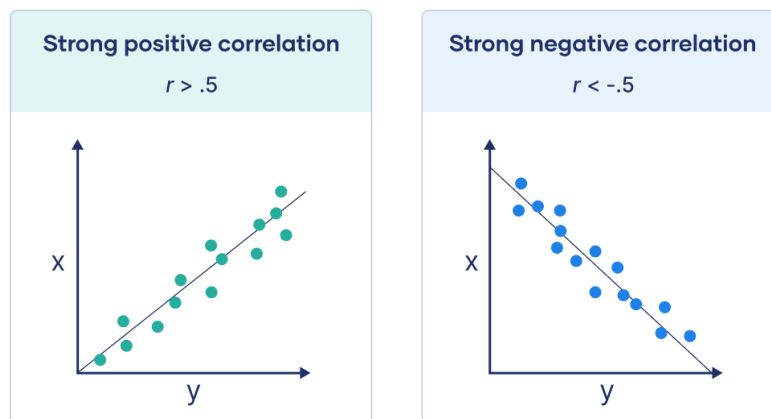
- the relationship is linear
- both variables are quantitative
- normally distributed
- have no outliers.

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.

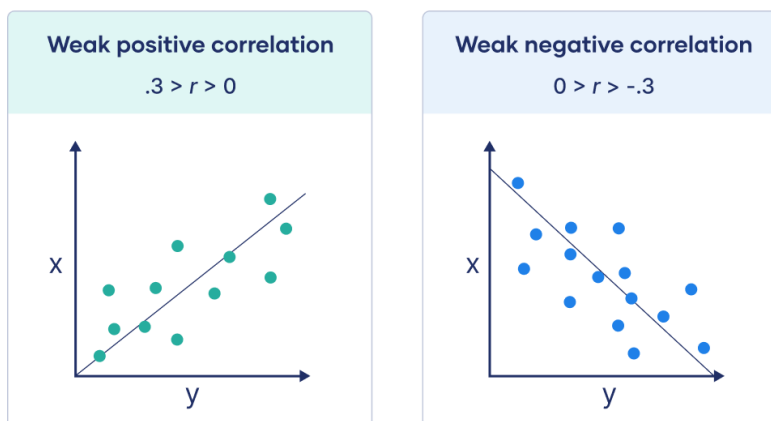
When r is 1 or -1 , all the points fall exactly on the line of best fit:



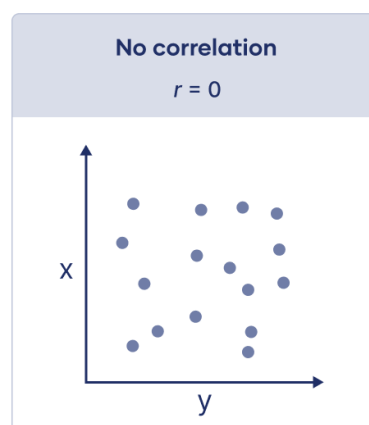
When r is greater than .5 or less than $-.5$, the points are close to the line of best fit:



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:



Formula:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

r	=	correlation coefficient
x_i	=	values of the x-variable in a sample
\bar{x}	=	mean of the values of the x-variable
y_i	=	values of the y-variable in a sample
\bar{y}	=	mean of the values of the y-variable

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

Ans: Feature scaling is one of the most important data pre-processing step in machine learning. Algorithms that compute the distance between the features are biased towards numerically larger values if the data is not scaled.

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 typically means rescales the values into a range of [0,1]. Standardization typically means rescales data to have a mean of 0 and a standard deviation of 1 (unit variance).

Difference between Normalization and Standardization

Normalization	Standardization
Minimum and maximum value of features are used for scaling	Mean and standard deviation is used for scaling.
It is used when features are of different scales.	It is used when we want to ensure zero mean and unit standard deviation.
Scales values between [0, 1] or [-1, 1].	It is not bounded to a certain range.
It is really affected by outliers.	It is much less affected by outliers.
Scikit-Learn provides a transformer called MinMaxScaler for Normalization.	Scikit-Learn provides a transformer called StandardScaler for standardization.
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.
It is useful when we don't know about the distribution	It is useful when the feature distribution is Normal or Gaussian.
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? (3 marks)

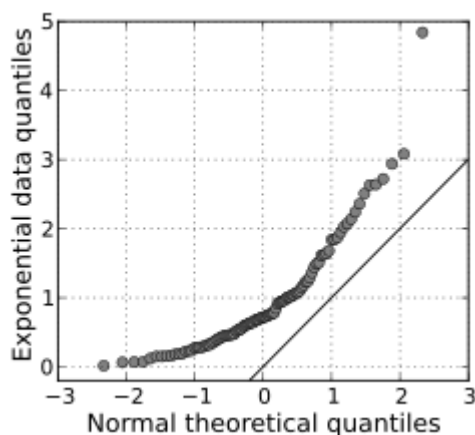
Ans: If there is perfect correlation, then $VIF = \infty$. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get $R^2 = 1$, which leads to $1/(1-R^2)$ infinity. To solve this problem we need to drop one of the variables from the dataset which is 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. (3 marks)

Ans: 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.