

# ASSESSMENT-2

## Data Wrangling and Regression Analysis

### Section A: Data Wrangling

1. What is the primary objective of data wrangling?

- ▣ a) Data visualization
- ▣ b) Data cleaning and transformation
- ▣ c) Statistical analysis
- ▣ d) Machine learning modeling

A. The primary objective of data wrangling is:

- b) Data cleaning and transformation

2. Explain the technique used to convert categorical data into numerical data. How does it help in data analysis?

A. The technique used to convert categorical data into numerical data is called "encoding." Encoding is essential because many machine learning algorithms and statistical models require numerical input. There are different methods for encoding categorical data, and two common approaches are:

Label Encoding:

In label encoding, each category is assigned a unique numerical label.

This is suitable for ordinal categorical data where there is an inherent order among the categories.

For example, if you have categories like "low," "medium," and "high," you could assign them labels like 0, 1, and 2, respectively.

One-Hot Encoding:

One-hot encoding creates binary columns for each category and indicates the presence of the category with a 1 or 0.

This is suitable for nominal categorical data where there is no inherent order among the categories.

For example, if you have categories like "red," "green," and "blue," one-hot encoding would create three binary columns, one for each color.

How it helps in data analysis:

**Compatibility with Algorithms:** Many machine learning algorithms and statistical models, such as linear regression or support vector machines, require numerical input. Encoding allows you to use categorical data in these algorithms.

**Improved Model Performance:** Properly encoding categorical data can improve the performance of machine learning models. Algorithms often perform better when they can interpret the encoded information accurately.

**Statistical Analysis:** For statistical analysis, numerical representations of categorical data are often necessary. It enables the application of statistical techniques that rely on numerical values.

**Handling Non-Numeric Features:** Data wrangling, including encoding, is crucial for handling non-numeric features in a dataset. It prepares the data for various analyses and modeling tasks.

### 3. How does LabelEncoding differ from OneHotEncoding?

A. Label Encoding and One-Hot Encoding are two different techniques used to convert categorical data into a numerical format. Here's how they differ:

Label Encoding:

In Label Encoding, each category is assigned a unique numerical label.

The numerical labels are often assigned in a way that preserves the ordinal relationship if there is one. For example, if you have categories like "low," "medium," and "high," Label Encoding might assign them labels 0, 1, and 2, respectively.

Label Encoding is suitable for ordinal categorical data where there is a meaningful order among the categories.

Example:

Categories: ["low", "medium", "high"]

Label Encoding: [0, 1, 2]

One-Hot Encoding:

One-Hot Encoding creates binary columns for each category and represents the presence or absence of a category with a 1 or 0.

It is suitable for nominal categorical data where there is no inherent order among the categories. Each category gets its own binary column, and only one of these columns will have a 1 for a given data point.

One-Hot Encoding avoids introducing ordinal relationships that might not exist in the original data.

Example:

Categories: ["red", "green", "blue"]

One-Hot Encoding:

red	green	blue
1	0	0
0	1	0
0	0	1

Key Differences:

Representation:

Label Encoding assigns a single numerical label to each category.

One-Hot Encoding creates binary columns, with each column representing a category.

Ordinality:

Label Encoding considers the ordinal relationship between categories.

One-Hot Encoding treats categories as nominal, without imposing any ordinal relationship.

Number of Columns:

Label Encoding results in a single column of numerical labels.

One-Hot Encoding results in multiple binary columns, one for each category.

Suitability:

Label Encoding is suitable for ordinal categorical data.

One-Hot Encoding is suitable for nominal categorical data.

The choice between Label Encoding and One-Hot Encoding depends on the nature of the categorical data and the requirements of the analysis or modeling task.

4. Describe a commonly used method for detecting outliers in a dataset. Why is it important to identify outliers?

A. One commonly used method for detecting outliers in a dataset is the "IQR (Interquartile Range) method." The IQR is a measure of statistical dispersion, representing the range between the first quartile (Q1) and the third quartile (Q3) of the data. Outliers are often identified based on their position relative to the IQR.

Here's how the IQR method for outlier detection works:

Calculate the IQR:

Compute the first quartile (Q1) and the third quartile (Q3) of the dataset.

Calculate the IQR as the difference between Q3 and Q1:  $IQR = Q3 - Q1$ .

Define Lower and Upper Bounds:

Define a lower bound as  $Q1 - 1.5 * IQR$ .

Define an upper bound as  $Q3 + 1.5 * IQR$ .

Identify Outliers:

Any data point that falls below the lower bound or above the upper bound is considered an outlier.

Optional: Adjust the Multiplier:

The multiplier (1.5 in the standard method) can be adjusted based on the desired sensitivity to outliers. A larger multiplier increases the range and identifies fewer outliers, while a smaller multiplier identifies more outliers.

Why is it important to identify outliers?

Impact on Descriptive Statistics:

Outliers can significantly affect descriptive statistics such as the mean and standard deviation. Identifying and handling outliers is crucial for obtaining accurate measures of central tendency and dispersion.

Model Performance:

Outliers can distort the results of statistical models. Removing or properly handling outliers is important for improving the performance and accuracy of predictive models.

Data Quality and Anomalies:

Outliers may indicate errors in data collection or measurement. Identifying outliers helps in maintaining data quality and uncovering potential anomalies.

Biases in Inferences:

Outliers can introduce biases in statistical inferences. It is essential to address outliers to make more reliable and robust conclusions from the data.

Data Understanding:

Identifying outliers provides insights into the characteristics of the data. It helps in understanding the distribution and behavior of the data points.

Robustness of Analyses:

Analyses and statistical tests are often more robust when conducted on datasets without significant outlier influence. Detecting and handling outliers contribute to the robustness of analyses.

Identifying outliers is crucial for maintaining data quality, ensuring the accuracy of statistical analyses, and improving the performance of predictive models. It is an essential step in the data preprocessing phase before performing further analyses or building models.

5. Explain how outliers are handled using the Quantile Method.

A. The Quantile Method, often based on the Interquartile Range (IQR), is a common approach for handling outliers in a dataset. The process involves identifying outliers based on the IQR and then taking appropriate actions. Here's a step-by-step explanation:

Calculate the Interquartile Range (IQR):

Compute the first quartile (Q1) and the third quartile (Q3) of the dataset.

Calculate the IQR as the difference between Q3 and Q1:  $IQR = Q3 - Q1$ .

Define Lower and Upper Bounds:

Define a lower bound as  $Q1 - 1.5 * IQR$ .

Define an upper bound as  $Q3 + 1.5 * IQR$ .

Identify Outliers:

Any data point that falls below the lower bound or above the upper bound is considered an outlier.

Handle Outliers:

Once outliers are identified, there are several ways to handle them:

**Removal:** Exclude the outliers from the dataset. This may be appropriate if the outliers are likely due to errors or anomalies.

**Transformation:** Transform the values of the outliers to bring them within a reasonable range. Common transformations include log transformations.

**Imputation:** Replace outlier values with the mean, median, or a more suitable imputation method.

Optional: Adjust the Multiplier:

The multiplier (1.5 in the standard method) can be adjusted based on the desired sensitivity to outliers. A larger multiplier increases the range and identifies fewer outliers, while a smaller multiplier identifies more outliers.

Example:

Suppose you have a dataset and you apply the Quantile Method with a multiplier of 1.5. You find that any data point below  $Q1 - 1.5 * IQR$  or above  $Q3 + 1.5 * IQR$  is an outlier. You identify a few data points as outliers and decide to either remove them, transform them, or impute them based on the nature of the data and the analysis goals.

Why Handle Outliers:

Handling outliers improves the accuracy of descriptive statistics and prevents them from being overly influenced by extreme values.

It contributes to the robustness of statistical analyses and machine learning models.

Outliers can introduce biases, and handling them ensures more reliable inferences.

It enhances the overall quality and integrity of the dataset.

The Quantile Method provides a systematic way to identify outliers based on the IQR, and subsequent actions are taken to handle these outliers based on the specific characteristics of the data and the objectives of the analysis.

6. Discuss the significance of a Box Plot in data analysis. How does it aid in identifying potential outliers?

A. A Box Plot, also known as a Box-and-Whisker Plot, is a graphical representation of the distribution of a dataset. It is a valuable tool in data analysis for several reasons and plays a key role in identifying potential outliers. Here are the significant aspects of a Box Plot:

Visualizing Data Distribution:

A Box Plot provides a visual summary of the distribution of data, including information about the median, quartiles, and potential outliers.

It consists of a box that represents the interquartile range (IQR), with a line inside marking the median. Whiskers extend from the box to the minimum and maximum values within a certain range.

Identifying Central Tendency:

The position of the median within the box indicates the central tendency of the data.

The length of the box (IQR) gives a sense of the spread of the central portion of the data.

Detecting Skewness and Symmetry:

The symmetry or skewness of the distribution can be observed by the position of the median within the box and the lengths of the whiskers.

Highlighting Potential Outliers:

Box Plots include individual data points beyond the whiskers, which are considered potential outliers.

Outliers are identified as points beyond a certain distance from the median or quartiles, often determined using a multiplier of the IQR (Interquartile Range). Points outside this range are displayed individually.

Comparing Multiple Distributions:

Box Plots are effective for comparing the distributions of different groups or categories in a dataset.

Side-by-side Box Plots make it easy to observe variations and differences in central tendency, spread, and the presence of outliers across groups.

Resistant to Extremes:

Box Plots are resistant to extreme values, making them suitable for visualizing datasets with outliers without being overly influenced by them.

How Box Plot Aids in Identifying Potential Outliers:

**Outlier Visualization:** The individual data points beyond the whiskers are explicitly shown in the Box Plot, making it easy to visually identify potential outliers.

**Quantitative Thresholds:** The whiskers are often defined based on statistical thresholds, such as 1.5 times the IQR. Data points beyond these thresholds are flagged as potential outliers.

**Comparison of Spread:** By comparing the length of the whiskers to the box length, one can quickly identify whether the spread of the data is influenced by potential outliers.

**Outlier Patterns:** The Box Plot can reveal patterns in the distribution of potential outliers, such as whether they are evenly distributed or concentrated in specific regions.

In summary, the significance of a Box Plot in data analysis lies in its ability to provide a concise summary of the distribution of data, visually highlight central tendency, detect skewness, and, importantly, identify potential outliers in a way that is easily interpretable. It aids in understanding the overall shape of the data and is particularly useful when exploring and comparing multiple datasets.

## Section B: Regression Analysis

7. What type of regression is employed when predicting a continuous target variable? A. When predicting a continuous target variable, linear regression is commonly employed. Linear regression is a type of regression analysis that models the relationship between a dependent variable (target) and one or more independent variables (features) by fitting a linear equation to observed data. The goal is to find the best-fitting linear line that minimizes the difference between the predicted values and the actual values of the target variable.

The general form of a simple linear regression equation for predicting a continuous target variable  $Y$  based on a single independent variable  $X$  is:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where  $\beta_0$  is the y intercept,  $\beta_1$  is the slope of the line, and  $\varepsilon$  is a random error term. In the traditional regression model, values of  $X$ -variable are assumed to be fixed by the experimenter. The model is still valid if  $X$  is random (as is more commonly the case), but only if  $X$  is measured without error.

In the case of multiple linear regression, where there are multiple independent variables, the equation becomes:

$$Y = \beta_0. + \beta_1 x1. + \beta_2 x2. + \dots + \beta_p x p. + \varepsilon$$

where  $E(\varepsilon) = 0$  and  $Var(\varepsilon) = \sigma^2$ .

Linear regression is widely used in various fields for tasks such as predicting house prices, sales, stock prices, and many other continuous numerical outcomes. It assumes a linear relationship between the

independent variables and the target variable, and its simplicity and interpretability make it a popular choice for regression analysis.

8. Identify and explain the two main types of regression.

A. The two main types of regression are:

Linear Regression:

**Definition:** Linear regression models the relationship between a dependent variable (target) and one or more independent variables (features) by fitting a linear equation to the observed data.

**Equation (Simple Linear Regression):** For a simple linear regression with one independent variable:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where  $\beta_0$  is the y intercept,  $\beta_1$  is the slope of the line, and  $\varepsilon$  is a random error term. In the traditional regression model, values of X-variable are assumed to be fixed by the experimenter. The model is still valid if X is random (as is more commonly the case), but only if X is measured without error.

**Equation (Multiple Linear Regression):** For multiple linear regression with multiple independent variables:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

**Objective:** Minimize the difference between predicted and actual values by adjusting the coefficients

**Use Cases:** Predicting continuous numerical outcomes, such as house prices, sales, or any variable with a linear relationship.

Logistic Regression:

**Definition:** Despite its name, logistic regression is used for binary classification problems, not regression. It models the probability of a binary outcome (0 or 1) based on one or more independent variables.

**Equation:** The logistic regression model transforms a linear combination of input features using the logistic function (sigmoid function) to produce a probability between 0 and 1:

$$P = \frac{1}{1 + e^{-z}}$$



$$1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_n X_n)}$$

**Use Cases:** Binary classification problems, such as spam detection, fraud detection, or any task where the outcome is categorical with two classes.

While linear regression is suitable for predicting continuous numerical outcomes, logistic regression is employed when dealing with binary classification tasks. Both types of regression involve estimating coefficients to describe the relationship between variables, but their applications and underlying assumptions differ. Linear regression assumes a linear relationship, while logistic regression models the probability of an event occurring.

9. When would you use Simple Linear Regression? Provide an example scenario.

A. Simple Linear Regression is used when there is a linear relationship between two variables, and you want to predict the value of one variable based on the values of another. Specifically, it is suitable when you have a single independent variable (predictor) and a single dependent variable (response) and assume that the relationship between them is linear.

Example Scenario:

Suppose you are a data analyst working for a real estate agency, and you want to understand the relationship between the square footage of a house (independent variable) and its selling price (dependent variable). You collect data on various houses, recording the square footage and the corresponding selling prices. Your goal is to build a model that can predict the selling price of a house based on its square footage.

In this scenario:

**Independent Variable (X):** Square footage of the house.

**Dependent Variable (Y):** Selling price of the house.

You can use Simple Linear Regression to create a model of the form:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where  $\beta_0$  is the y intercept,  $\beta_1$  is the slope of the line, and  $\varepsilon$  is a random error term. In the traditional regression model, values of X-variable are assumed to be fixed by the experimenter. The model is still valid if X is random (as is more commonly the case), but only if X is measured without error.

Usage of Simple Linear Regression in this scenario allows you to understand the linear relationship between square footage and selling price, make predictions, and potentially provide insights to clients or stakeholders in the real estate industry. Keep in mind that the assumption of linearity should be reasonable for accurate model predictions. If the relationship is more complex, you might consider multiple linear regression or other modeling techniques.

10. In Multi Linear Regression, how many independent variables are typically involved?

A. In Multiple Linear Regression, there are typically two or more independent variables involved. The term "multiple" in multiple linear regression refers to the fact that there are multiple predictors or independent variables used to model the relationship with a single dependent variable.

The general form of the multiple linear regression equation is:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

where  $E(\varepsilon) = 0$  and  $\text{Var}(\varepsilon) = \sigma^2$ .

Linear regression is widely used in various fields for tasks such as predicting house prices, sales, stock prices, and many other continuous numerical outcomes. It assumes a linear relationship between the independent variables and the target variable, and its simplicity and interpretability make it a popular choice for regression analysis.

Multiple Linear Regression is a powerful tool when dealing with situations where the relationship between the dependent variable and the outcome is influenced by multiple factors. It allows for the consideration of the combined effect of multiple predictors on the response variable.

11. When should Polynomial Regression be utilized? Provide a scenario where Polynomial Regression would be preferable over Simple Linear Regression.

A. Polynomial Regression should be utilized when the relationship between the independent variable and the dependent variable is not linear but exhibits a more complex, nonlinear pattern. In Polynomial Regression, the relationship is modeled as an nth-degree polynomial equation, allowing for curved or non-linear fits.

Scenario where Polynomial Regression is preferable over Simple Linear Regression:

Example Scenario: Predicting Nonlinear Relationships

Suppose you are analyzing the growth of a plant over time, and you want to predict the plant's height based on the number of days since planting.

In this case, the relationship may not be linear; instead, the plant growth might follow a curved pattern.

**Independent Variable (X):** Number of days since planting.

**Dependent Variable (Y):** Plant height.

In a scenario like this, Simple Linear Regression might not capture the underlying pattern accurately. The growth of the plant may not follow a straight line; it could exhibit a curve. Polynomial Regression can be employed to model the relationship more flexibly.

The Polynomial Regression equation takes the form:

$$Y = \beta_0. + \beta_1 x_1. + \beta_2 x_2. + \dots + \beta_p x_p. + \varepsilon$$

Polynomial Regression is preferable over Simple Linear Regression when the relationship between the variables is not linear, and a more flexible model is needed to capture non-linear patterns in the data. It is essential to use Polynomial Regression judiciously and consider the trade-off between model complexity and the risk of overfitting to the specific dataset.

12. What does a higher degree polynomial represent in Polynomial Regression? How does it affect the model's complexity?

A. In Polynomial Regression, the degree of the polynomial represents the highest power of the independent variable in the regression equation. A higher degree polynomial includes more terms with higher powers of the independent variable, allowing the model to capture more complex, non-linear patterns in the data.

The general form of a polynomial regression equation is:

$$Y = \beta_0. + \beta_1 x_1. + \beta_2 x_2. + \dots + \beta_p x_p. + \varepsilon$$

Effect of Higher Degree Polynomial:

Increased Flexibility:

A higher degree polynomial provides the model with increased flexibility to fit complex patterns in the data. It allows the regression line to curve and capture non-linear relationships.

Overfitting Risk:

While higher flexibility is beneficial for capturing intricate patterns, it also increases the risk of overfitting. Overfitting occurs when the model fits the training data too closely, capturing noise and idiosyncrasies that may not generalize well to new data.

Model Complexity:

The model becomes more complex with higher-degree polynomials. Model complexity refers to the number of parameters (coefficients) in the model. As the degree increases, the number of parameters also increases, making the model more intricate.

Balancing Complexity and Generalization:

Choosing the appropriate degree for the polynomial involves a trade-off between capturing complex patterns in the training data and maintaining the model's ability to generalize to new, unseen data. A balance must be struck to prevent overfitting.

Visualizing Model Complexity:

Visualizing the regression line or curve for different degrees helps in understanding how the model's complexity changes. Higher-degree polynomials may result in curves with more twists and turns.

Hyperparameter Tuning:

Practitioners often perform hyperparameter tuning to find the optimal degree of the polynomial. Techniques such as cross-validation can help assess the model's performance on different subsets of the data.

A higher degree polynomial in Polynomial Regression allows the model to capture more complex patterns, but it comes with the risk of overfitting. Careful consideration and validation are needed to choose an appropriate degree that balances model complexity and generalization to new data.

13. Highlight the key difference between Multi Linear Regression and Polynomial Regression.

A. The key difference between Multi Linear Regression and Polynomial Regression lies in the nature of the relationships they can model:

Multi Linear Regression:

**Nature of Relationship:** Multi Linear Regression models the relationship between a dependent variable and two or more independent variables in a linear fashion.

**Equation:** The equation for Multi Linear Regression is a linear combination of independent variables with respective coefficients.

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

**Usage:** Suitable for scenarios where there are multiple predictors, each with a linear relationship with the dependent variable.

Polynomial Regression:

**Nature of Relationship:** Polynomial Regression models the relationship between a dependent variable and an independent variable in a non-linear fashion, allowing for more complex patterns.

**Equation:** The equation for Polynomial Regression includes higher-order terms of the independent variable, creating a polynomial equation.

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

**Usage:** Suitable for scenarios where the relationship between the dependent variable and the independent variable exhibits a curved or non-linear pattern.

Summary:

Multi Linear Regression focuses on modeling linear relationships between multiple independent variables and a dependent variable.

Polynomial Regression allows for the modeling of non-linear relationships by incorporating higher-order terms of the independent variable in the regression equation.

In essence, while Multi Linear Regression assumes a linear relationship with multiple predictors, Polynomial Regression provides more flexibility to capture non-linear patterns in the data by introducing polynomial terms. The choice between these regression techniques depends on the nature of the data and the underlying relationship that needs to be modeled.

14. Explain the scenario in which Multi Linear Regression is the most appropriate regression technique.

A. Multi Linear Regression is the most appropriate regression technique in scenarios where there is a need to model the relationship between a dependent variable and two or more independent variables, and the relationship is assumed to be linear. It is suitable for situations where multiple factors or predictors contribute to the variation in the dependent variable.

Key Characteristics of Scenarios for Multi Linear Regression:

Multiple Predictors:

There are two or more independent variables available for predicting the dependent variable. Each independent variable is expected to contribute independently to the variation in the dependent variable.

Linear Relationship:

The assumption is that the relationship between the dependent variable and the independent variables is linear. This means that a change in any of the independent variables is associated with a constant change in the dependent variable.

Quantitative Data:

The variables involved in the regression analysis are quantitative (continuous or discrete numerical variables). Multi Linear Regression is not suitable for categorical predictors.

### Independence of Predictors:

The independent variables should be reasonably independent of each other. High multicollinearity (strong correlation) between predictors can affect the stability and interpretability of the regression coefficients.

### Example Scenario:

Suppose you are working on predicting the sales of a product in a retail setting. You have data on various factors that might influence sales, such as advertising expenditure, pricing, and the store's location. In this scenario:

Dependent Variable (Y): Sales of the product.

**Independent Variables (X1, X2, ...):** Advertising expenditure, pricing, store's location, and potentially other factors.

The goal is to model how changes in advertising expenditure, pricing, and location are associated with changes in product sales. A Multi Linear Regression model can be used to estimate the coefficients for each predictor, providing insights into the relative impact of each factor on sales.

### Steps in Multi Linear Regression:

#### Formulation of the Model:

Develop the Multi Linear Regression equation based on the available independent variables.

#### Estimation of Coefficients:

Use statistical methods to estimate the coefficients of the independent variables.

#### Model Validation:

Evaluate the model's performance using techniques like hypothesis testing, residual analysis, and model fit statistics.

#### Interpretation of Results:

Interpret the coefficients to understand the impact of each independent variable on the dependent variable.

In summary, Multi Linear Regression is suitable when dealing with multiple predictors and assuming a linear relationship between these predictors and the dependent variable. It is commonly used in various fields, including economics, finance, marketing, and social sciences, to analyze and understand complex relationships in data.

### 15. What is the primary goal of regression analysis?

A. The primary goal of regression analysis is to examine and model the relationship between a dependent variable (or response variable) and one or more independent variables (or predictor variables). The fundamental objective is to understand how changes in the independent variables are associated with changes in the dependent variable.

Regression analysis aims to quantify and describe this relationship,

allowing for prediction, inference, and insights into the underlying patterns in the data.

Key goals of regression analysis include:

Modeling Relationships:

Identify and model the relationship between the dependent variable and the independent variables. This involves determining the functional form of the relationship (e.g., linear, non-linear) and estimating the parameters of the model.

Prediction:

Use the established regression model to make predictions of the dependent variable's values based on the values of the independent variables. This is particularly valuable for forecasting and decision-making.

Inference:

Draw inferences about the population parameters from the sample data. Inferential statistics help assess the reliability and significance of the relationships observed in the data.

Understanding Variable Contributions:

Understand the individual contributions of each independent variable to the variation in the dependent variable. This helps identify which factors are more influential and provides insights into their impact.

Model Assessment:

Evaluate the quality and performance of the regression model. This involves checking assumptions, assessing goodness of fit, and identifying any potential issues such as multicollinearity or outliers.

Variable Selection:

Identify which independent variables are statistically significant and contribute meaningfully to the model. Variable selection helps simplify the model and improve its interpretability.

Quantifying Relationships:

Quantify the strength and direction of the relationships between variables using regression coefficients. These coefficients indicate the change in the dependent variable associated with a one-unit change in the corresponding independent variable.

Assessing Model Fit:

Evaluate how well the regression model fits the observed data. Good model fit ensures that the model accurately represents the patterns and trends present in the dataset.

Overall, regression analysis serves as a powerful statistical tool for exploring, modeling, and understanding the relationships between variables. Whether used for prediction, hypothesis testing, or decision support, regression analysis is widely applied in various fields, including economics, finance, biology, psychology, and engineering.

