**AI / ML Training**

**Assignment: Data Wrangling and Regression Analysis**

**Instructions:** Answer the following questions to the best of your ability. Provide concise explanations where necessary.

**Section A: Data Wrangling (Questions 1-6)**

1. What is the primary objective of data wrangling?
   - a) Data visualization
   - b) Data cleaning and transformation
   - c) Statistical analysis
   - d) Machine learning modeling

   Answer:

   1. **Data Visualization:**
      - **Objective:** To present data in a graphical or visual format.
      - **Purpose:** Facilitates easier understanding of patterns, trends, and insights within the data.

   2. **Statistical Analysis:**
      - **Objective:** To analyze and interpret patterns and trends within the data using statistical methods.
      - **Purpose:** Aims to draw meaningful conclusions, make predictions, or identify relationships in the data.

   3. **Machine Learning Modeling:**
      - **Objective:** To develop models that can learn patterns from data and make predictions or decisions without explicit programming.
      - **Purpose:** Enables automation of tasks, predictions, or decision-making based on patterns identified in the data.

   4. **Data Cleaning and Transformation (Data Wrangling):**
      - **Objective:** To clean and transform raw data into a usable and structured format.
      - **Purpose:** Prepares the data for analysis by handling missing values, outliers, and ensuring it is in a suitable format for statistical analysis or machine learning modeling.

   5. **Data Wrangling:**
      - **Objective:** The broader process of acquiring, cleaning, and transforming raw data.
      - **Purpose:** Ensures that data is in a suitable format for analysis, whether it be for statistical analysis or machine learning modeling.

2. Explain the technique used to convert categorical data into numerical data. Howdoes it help in data analysis?
   Answer:

The technique used to convert categorical data into numerical data is known as **Encoding**. It involves representing categorical variables with numerical values. There are two common methods for encoding categorical data: **Label Encoding** and **One-Hot Encoding**.

1. **Label Encoding:**
   - In Label Encoding, each category is assigned a unique numerical label or code.
   - It is suitable for ordinal data, where the categories have a meaningful order.
   - For example, if we have categories like "Low," "Medium," and "High," Label Encoding might assign them numerical labels like 0, 1, and 2.
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 data, where categories don't have a meaningful order.
   - Each category gets its own column, and the presence of the category is indicated by a 1 in the respective column.
   - For example, if we have categories like "Red," "Blue," and "Green," One-Hot Encoding might create three columns with binary values indicating the presence of each color.

**How it helps in data analysis:**

- **Numeric Input for Algorithms:** Many machine learning algorithms require numeric input. By converting categorical data into numerical format, we enable the use of these algorithms.
- **Improved Model Performance:** Machine learning models often perform better when working with numerical data. Encoding allows us to utilize a broader range of models and techniques.
- **Handling Categorical Variables:** Some statistical techniques and algorithms can only handle numerical values. Encoding provides a way to include categorical variables in these analyses.
- **Preserving Ordinal Information:** In the case of Label Encoding, when there is an ordinal relationship among categories, the encoded values can preserve that order.

3. How does LabelEncoding differ from OneHotEncoding?
   Answer:

**Difference between LabelEncoding and OneHotEncoding:**

- **LabelEncoding:** It involves assigning a unique numerical label to each category. It is suitable for ordinal data but may introduce ordinal relationships.
- **OneHotEncoding:** It represents each category with a binary vector. It is suitable for nominal data and avoids introducing ordinal relationships.

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

**Detecting Outliers:**

- A commonly used method is the "IQR (Interquartile Range) Method." Outliers are identified as data points that fall below Q1 - 1.5 * IQR or above Q3 + 1.5 * IQR, where Q1 and Q3 are the first and third quartiles.

**Importance:**

- Outliers can skew statistical measures and impact the accuracy of predictive models.
- Identifying outliers is crucial for maintaining data integrity and making informed decisions.

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

**Handling Outliers using the Quantile Method:**

- Outliers can be capped or winsorized using the values of Q1 - k * IQR or Q3 + k * IQR, where k is a constant.
- This helps in mitigating the impact of outliers on statistical measures.

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

**Significance of a Box Plot:**

- A Box Plot visually represents the distribution of data, showing the median, quartiles, and potential outliers.
- It aids in identifying the spread of data, central tendency, and skewness.
- Outliers can be visually identified as points beyond the "whiskers" of the box plot.

**Section B: Regression Analysis (Questions 7-15)**

7. What type of regression is employed when predicting a continuous targetvariable?
    Answer:

**Type of Regression for Predicting Continuous Target Variable:**

- The type of regression employed when predicting a continuous target variable is **Linear Regression.**

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

**Two Main Types of Regression:**

- The two main types of regression are:
    - **Simple Linear Regression:** Involves one independent variable to predict the dependent variable.
    - **Multiple Linear Regression:** Involves two or more independent variables to predict the dependent variable.

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

**Use of Simple Linear Regression:**

- Simple Linear Regression is used when there is a linear relationship between the independent variable and the dependent variable. It is suitable when you have only one independent variable.
- Example Scenario: Predicting the score of a student based on the number of hours they studied.

10. In Multi Linear Regression, how many independent variables are typicallyinvolved?
    Answer:

**Number of Independent Variables in Multi Linear Regression:**

- In Multi Linear Regression, there are typically two or more independent variables involved

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

**Use of Polynomial Regression:**

- Polynomial Regression is utilized when the relationship between the independent variable and the dependent variable is nonlinear.
- Scenario: Predicting the sales of a product based on historical data, where the relationship exhibits curves or bends.

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

**Higher Degree Polynomial in Polynomial Regression:**

- A higher degree polynomial in Polynomial Regression represents a more complex model.

- The degree of the polynomial determines the number of bends or curves the model can fit.
- Higher degree polynomials can lead to overfitting if not carefully chosen, making the model too flexible and fitting noise in the data.

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

**Key Difference between Multi Linear Regression and Polynomial Regression:**
- The key difference lies in the form of the relationship:
  - Multi Linear Regression assumes a linear relationship.
  - Polynomial Regression allows for nonlinear relationships by introducing polynomial terms.

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

**Scenario for Multi Linear Regression:**
- Multi Linear Regression is appropriate when there are multiple independent variables influencing the dependent variable simultaneously.
- Example: Predicting the price of a house based on features like square footage, number of bedrooms, and location.

15. What is the primary goal of regression analysis?
Answer:

**Primary Goal of Regression Analysis:**
- The primary goal of regression analysis is to model the relationship between the dependent variable and one or more independent variables. It aims to understand and quantify how changes in independent variables are associated with changes in the dependent variable.