# ASSESSMENT-2                    SOHAN DAS

## 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

Ans. **Data cleaning and transformation:** The primary objective of data wrangling is data cleaning and transformation because it involves preparing raw data for analysis by removing inconsistencies, errors, and missing values, as well as transforming the data into a format suitable for analysis. While data visualization, statistical analysis, and machine learning modeling are important steps in the data analysis process, data wrangling is necessary to ensure that the data is accurate, complete, and structured properly for further analysis.

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

Ans. One popular method for transforming categorical data into numerical data is "one-hot encoding," sometimes known as "dummy encoding." This method uses binary vectors to represent each category inside a categorical variable. More specifically, a new binary column is produced for every category. A data point's matching binary column is set to 1 if it falls within that category, and to 0 otherwise. Through this procedure, the categorical data is converted into numerical data that machine learning algorithms can easily understand and handle. Imagine, for instance, that the category "Colour" has three categories: "Red," "Blue," and "Green." The variable would be represented by three binary columns using one-hot encoding: "Color_Red," "Color_Blue," and "Color_Green." Should a data point be "Red," the column labeled "Color_Red" would.

3. How does LabelEncoding differ from OneHotEncoding?

Ans.

| LabelEncoding | OneHotEncoding |
|---|---|
| 1. A distinct number label, or code, is assigned to each category within a categorical variable using LabelEncoding. | 1. In OneHotEncoding, each magnificence interior of a particular variable is represented as a binary vector. |
| 2. In essence, it substitutes an integer beginning with 0, 1, 2, and so on for each category. | 2. It creates a cutting-edge binary column for each category, and if a data element belongs to a category, the corresponding binary column is prepared to 1; otherwise, it`s a long way set to 0. |
| 3. It may not be appropriate for all categorical variables to use this encoding since it indicates an ordinal relationship between the categories, suggesting that one category is more or less than another. | 3. This encoding does now not impose any ordinal courting amongst training and is suitable for explicit variables without inherent order. |
| 4. LabelEncoding works well when there is a distinct order and the categorical variable has inherent ordinality, such as low, medium, and high. | 4. OneHotEncoding is useful even as there`s no ordinal relationship among training or even as the specific variable has nominal data. |
| 5. Example: ["Red", "Blue", "Green"] is probably encoded as [0, 1, 2]. | 5. Example: ["Red", "Blue", "Green"] is probably encoded as [[1, 0, 0], [0, 1, 0], [0, 0, 1]]. |

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

Ans. One commonly used method for detecting outliers in a dataset is the z-score method. The z-score measures how many standard deviations a data point is from the mean of the dataset.
1. Calculate the mean ($\mu$) and standard deviation ($\sigma$) of the dataset.
2. For each data point $x$ in the dataset, calculate its z-score using the formula:
3. Data points with z-scores that fall above or below a certain threshold (typically $\pm 3$) or $\pm 2$) standard deviations) are considered outliers.
**A. Data Quality Assurance:** Outliers may indicate errors in the data collection process, such as measurement errors or data entry mistakes. Identifying and correcting these errors improves the overall quality of the dataset.
**B. Statistical Analysis:** Outliers can significantly affect statistical measures such as the mean and standard deviation. By detecting and handling outliers appropriately, analysts can ensure that statistical analyses are more accurate and reliable.
**C. Model Performance:** Outliers can have a disproportionate impact on the performance of machine learning models. By detecting and removing outliers, models can be trained on cleaner data, leading to better performance and generalization to unseen data.
**D. Insight Generation:** Outliers may also contain valuable information about unusual phenomena or rare events in the data. By identifying outliers and investigating their underlying causes, analysts can gain insights into potential trends, anomalies, or patterns that may not be apparent from the rest of the data.

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

Ans. The quantile method, also known as the Interquartile Range (IQR) method, is a technique used to handle outliers in a data set.
**1. Calculate quartiles:** First, calculate the first quartile (Q1) and third quartile (Q3) of the data set. The first quartile (Q1) is the value below which 25% of the data falls and the third quartile (Q3) is the value below which 75% of the data falls.
**2. Calculate the interquartile range (IQR):** The interquartile range (IQR) is calculated as the difference between the third quartile (Q3) and the first quartile (Q1):
**3. Define outliers:** Set lower and upper limits for outliers. Typically, outliers are considered values less than $Q1 - 1.5 \times IQR$ or greater than $Q3 1.5 \times IQR$).
**4. Handling exceptions:** Exceptions can be handled in many different ways:- **Remove outliers**: Exclude data points that fall outside the lower and upper limits.- **Bound Exception**: Replace exceptions with the closest non-exceptional value within the limit.- **Data Transformation**: Transform data using techniques such as winterization, logarithmic transformation, or robust scaling to minimize the effect of outliers.
**5. Conduct analysis:** Once the outliers are handled, perform data analysis using the cleaned data set. Quantile methods are especially useful when the data distribution is skewed or non-normal. By focusing on the middle half of the data (interquartile range), this method provides a robust measure of the data distribution while ignoring outliers that could bias the analysis provides a more robust approach to detecting and handling outliers than methods based solely on the mean and standard deviation, especially when skewed or tailed distributions are present. heavy.

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

Ans. A box plot, also known as a box and whisker plot, is a valuable visualization tool used in data analysis to summarize the distribution of a data set and identify potential outliers. Here's how it makes data analysis and outlier detection easier:
**1. Visualize data distribution:** Box plots provide a visual summary of the central tendency, spread, and shape of a data set's distribution. It consists of a box representing the interquartile range (IQR), with a line inside the box representing the median. The "whisker" extends from the box to points that do not lie beyond the minimum and maximum values at some distance from the quartiles.
**2. Determine central tendency:** The position of the middle line in the box represents the central tendency of the data. If the median is closer to one end of the box, this indicates an asymmetric distribution.
**3. Evaluate spread and variability:** Box length (IQR) represents the spread or variability of the data. Longer boxes suggest more variation, while shorter boxes indicate less variation.

**4. Potential outlier detection:** Outliers, which are data points located significantly further from the tips of the whiskers, can be easily identified on Box Plot. In general, values greater than $Q3 + 1.5 \times IQR$ or less than $Q1 - 1.5 \times IQR$ are considered potential outliers. These data points are plotted individually as dots outside the whiskers, helping analysts spot them quickly.

**5. Group Comparison:** Box plots are especially useful for comparing the distribution of different groups or categories within a data set. By plotting multiple boxes side by side, analysts can visually compare the central tendency, spread, and variability of each group.

**6. Interpret outliers:** The presence of outliers can provide valuable information about the data, indicating potential errors, anomalies, or interesting phenomena that need further investigation. Box plots help interpret outliers by visually highlighting their location relative to the rest of the data distribution. In short, Box Plots are powerful data analysis tools for summarizing data distribution, identifying potential outliers, comparing groups, and better understanding data set characteristics. They provide a concise and intuitive way to visualize key summary statistics and evaluate the overall shape and variability of the data.

---

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

7. What type of regression is employed when predicting a continuous target variable?

Ans. When predicting a continuous target variable, linear regression is often used. Linear regression models the relationship between a target variable and one or more predictor variables by fitting a linear equation to observed data. The goal is to find the best-fit line (or hyperplane in higher dimensions) that minimizes the difference between the observed and predicted values of the target variable.

Linear regression is widely used in various fields such as statistics, economics, finance, engineering, and machine learning for tasks such as forecasting, inference, and modeling relationships between variables.

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

Ans. There are predominant kinds of regression:

**1. Linear Regression:**
- Linear regression fashions the connection among a structured variable (goal variable) and one or greater unbiased variables (predictor variables) with the aid of using a linear equation for Anti-information. - An easy linear regression equation with predictor variables is:
[ y = beta_0 + beta_1 x + epsilon ]
or:
- ( y ) is the structured variable.
- ( x ) is the unbiased variable.
- ( beta_0 ) is the y-intercept (intercept).
- ( beta_1 ) is the slope (coefficient),
- ( epsilon ) is the mistake period representing the distinction among the found and ana expressed immediate expected values. - Linear regression objectives to discover the best-becoming immediate line with the aid of minimizing the sum of squares of the variations among found and expected values (referred to as the residual sum of squares or RSS).

**2. Logistic regression:**
- Logistic regression is used whilst the structured variable is expressed, generally binary (yes/no, zero/1, etc.). - Despite its name, logistic regression is greater of a type set of rules than a regression set of rules as it predicts the opportunity that a statement belongs to a specific class. - Logistic regression fashions the connection between an express structured variable and one or greater unbiased variables with the aid of estimating the opportunity that a statement falls into a specific category. - The logistic function (sigmoid function) is used to transform a linear mixture of predictor variables into possibilities among zero and 1. The logistic regression equation is:
[ p(y=1 | x) = frac} ]

or:
- ( p(y=1 | x) ) is the opportunity that the structured variable is of kind 1, given the predictor variable ( x ).
- ( beta_0 ) is the intercept, on
- ( beta_1 ) is a coefficient.
- ( e ) is the bottom of the herbal logarithm. - Logistic regression estimates the coefficients that maximize the opportunity of staring at information for the version parameters.

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

Ans. Simple linear regression is generally used whilst you need to apprehend the connection among non-stop variables, wherein one variable (the predictor or unbiased variable) is used to expect the price of some other variable (the reaction or based variable). You could use easy linear regression whilst you accept as true with there may be a linear dating among the 2 variables, which means that modifications inside the predictor variable are related to modifications inside the reaction variable in a regular ratio. Suppose you need to research the connection between the range of hours studied and students` examination rankings. In this scenario, you may use easy linear regression to apprehend how the range of hours studied (predictor variable) impacts students' examination rankings (reaction variable). You hypothesize that there may be a linear dating among the 2 variables, which means that scholars who examine extra hours are anticipated to acquire better examination rankings and accumulate statistics at the range of hours every scholar studied and their corresponding examination rankings. By acting easy linear regression evaluation in this statistics, you may estimate the linear equation that high-quality describes the connection among the range of hours studied and examination rankings. This equation can then be used to expect examination rankings for college students primarily based totally on the range of hours they examine.

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

Ans. Simple linear regression is normally used when you need to recognize the connection among non-stop variables, in which one variable (the predictor or impartial variable) is used to expect the price of some other variable (the reaction or based variable). You might use easy linear regression while you consider there may be a linear dating among the 2 variables, which means that adjustments inside the predictor variable are related to adjustments inside the reaction variable in a consistent ratio. Suppose you need to research the connection between the range of hours studied and students` examination rankings. In this scenario, you may use easy linear regression to recognize how the range of hours studied (predictor variable) affects students' examination rankings (reaction variable). You hypothesize that there may be a linear dating among the 2 variables, which means that scholars who take a look at extra hours are anticipated to gain better You acquire information at the range of hours every scholar studied and their corresponding examination rankings. By appearing an easy linear regression evaluation in this information, you may estimate the linear equation that describes the connection between the range of hours studied and examination rankings. This equation can then be used to expect examination rankings for college kids primarily based totally on the range of hours they take a look at.

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

Ans. Polynomial regression needs to be applied whilst the connection between the impartial variable(s) and the based variable isn't linear however well-known shows a curved or non-linear pattern. Polynomial regression extends the idea of linear regression and suggests introducing polynomial phrases of the predictor variable(s) to seize non-linear relationships. A situation wherein polynomial regression might be most advantageous over easy linear regression is whilst reading records that suggest a non-linear fashion or pattern. Consider a state of affairs wherein you're reading the connection between the temperature (impartial variable) and the price of ice cream income (based variable). Initially, you may anticipate linear dating and use easy linear regression to version the records. However, upon plotting the records, you study that the connection between temperature and ice cream income isn't strictly linear however as an alternative well-known shows a curved pattern, wherein income begins with growth with temperature but then begins to lower after a positive factor because of excessive heat. In this situation, the usage of easy linear regression would possibly bring about bad health to the records and faulty predictions. Instead, polynomial regression may be hired to seize the non-linear dating greater accurately. By including polynomial phrases (e.g., quadratic or cubic phrases) inside the regression version, polynomial regression can higher seize the

curvature of the connection among temperature and ice cream in, main to advanced predictive overall performance and a greater correct illustration of the underlying records pattern.

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

Ans. In Polynomial Regression, the diploma of the polynomial represents the best electricity of the unbiased variable(s) inside the regression equation. A better diploma polynomial introduces extra curvature to the connection between the unbiased variable(s) and the established variable, permitting the version to seize greater complicated styles inside the facts.- A first-diploma polynomial represents a linear courting among the unbiased variable(s) and the established variable.- A second-diploma polynomial (quadratic) introduces a curve to the connection, making an allowance for a greater bendy in shape to the facts.- Higher-diploma polynomials (e.g., cubic, quartic, etc.) introduce an increasing number of complicated curves and bends to the connection, permitting the version to seize tricky styles withinside the facts. However, growing the diploma of the polynomial additionally will increase the version`s complexity. This elevated complexity may have each benefits and disadvantages:- Higher-diploma polynomials can seize greater complicated relationships inside the facts, permitting the version to shape the schooling facts greater closely.- They can offer a higher shape to nonlinear facts styles, mainly to stepped forward predictive performance.- Increased complexity can cause overfitting, in which the version learns to seize noise and random fluctuations inside the schooling facts in preference to the underlying actual courting. This can bring about negative generalization to new, unseen facts.- Higher-diploma polynomials require greater parameters to be predicted from the facts, main to a better danger of overfitting, especially while the quantity of facts factors is confined relative.

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

Ans. The key difference between Multi Linear Regression and Polynomial Regression lies withinside the shape of dating they model a number of the impartial variables and the established variable:- In Multi Linear Regression, the model assumes a linear dating of a number of the established variables and more than one impartial variables.- It extends the concept of clean linear regression to embody more than one predictor variable.- The regression equation in Multi Linear Regression takes the form:[ $y = beta\_0 beta\_1 x\_1 beta\_2 x\_2$ lots $beta\_n x\_n epsilon$ ]- ( $y$ )is the established variable,- ( $x\_1, x\_2$, lots, $x\_n$ ) are the impartial variables,- ( $beta\_0, beta\_1, beta\_2$, lots, $beta\_n$ ) are the coefficients representing the effect of each impartial variable on the established variable,- ( $epsilon$ ) is the error term.- In Polynomial Regression, the model permits for non-linear relationships between a number of the established variable and the impartial variable(s) with the useful resource of the usage of introducing polynomial terms the- It captures the curvature and non-linearity inside the statistics with the useful resource of the usage turning into a polynomial characteristic to the decided statistics.- The regression equation in Polynomial Regression can take the form of a polynomial characteristic, such as:[ $y = beta\_0 beta\_1x beta\_2 x^2$ lots $beta\_d x^d epsilon$ ]- ( $y$ ) is the established variable,- ( $x$ ) is the impartial variable,- ( $beta\_0, beta\_1, beta\_2$, lots, $beta\_d$ ) are the coefficients of the polynomial terms,- ( $d$ ) is the degree of the polynomial, representing the nice electricity of ( $x$ ) withinside the equation,- ( $epsilon$ ).

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

Ans. Multi Linear Regression is maximum suitable in eventualities in which there are a couple of impartial variables that together impact the structured variable. Here`s a state of affairs in which Multi Linear Regression is suitable: Suppose you're an actual property agent tasked with figuring out the promotion rate of homes primarily based totally on different factors together with the scale of the residence, the variety of bedrooms, the wide variety of bathrooms, and the neighbor hood's crime rate. In this state of affairs, there are a couple of impartial variables (size, wide variety of bedrooms, wide variety of bathrooms, crime rate) that together have an effect Using Multi Linear Regression, you could expand a version that carries a lot of these impartial variables to are expecting the promoting rate of homes. The regression equation could take the form:[ text = $beta\_0 beta\_1$ times text $beta\_2$ times text $beta\_3$ times text $beta\_4$ times text epsilon ]- ( text ) is the structured variable,- ( text ), ( text ), ( text ), and ( text ) are the impartial variables,- ( $beta\_0, beta\_1, beta\_2, beta\_3, beta\_4$ ) are the coefficients representing the impact of every impartial variable at the promoting rate,- ( $epsilon$ ) is the mistake term. In this state of affairs, Multi Linear Regression permits you to research how every impartial variable (size, wide variety of bedrooms, wide variety of

bathrooms, crime rate) contributes to the general promoting rate of homes at the same time as controlling for the results of different variables. It affords complete expertise of the elements influencing residence fees and enables making knowledgeable selections inside the actual property market.

15. What is the primary goal of regression analysis?

Ans. The number one purpose of regression evaluation is to apprehend and quantify the connection among one or greater unbiased variables (predictors) and a based variable (outcome) if you want to make predictions, infer patterns, and draw insights from the data.1. **Model the Relationship**: Regression evaluation seeks to version the purposeful courting among the unbiased variables and the based variable. It aims to discover how adjustments inside the unbiased variables are related to adjustments inside the based variable.2. **Predict Outcomes**: Regression fashions may be used to expect the price of the based variable primarily based totally on the values of the unbiased variables. These predictions may be precious for making forecasts, estimating destiny trends, or assessing the effect of interventions.3. **Inferential Purposes**: Regression evaluation may be used for inferential purposes, together with checking out hypotheses approximately the connection among variables, assessing the statistical importance of predictors, and making inferences approximately populace parameters primarily based totally on pattern data.4. **Control and Adjustment**: Regression evaluation lets in for controlling and adjusting for ability confounding variables or covariates that could have an impact on the connection among the unbiased and based variables. This enables in keeping apart the precise results of hobby and lowering bias inside the estimates.