# CertNexus Certified Data Science Practitioner (CDSP) Practice Exam (Sample)

**Study Guide**



BY EXAMZIFY

## Everything you need from our exam experts!

# **Questions**

1. **What term is used to describe a model that is deemed useful for its intended task?**

   A. Overfitting

   B. Skillful

   C. Underfitting

   D. Selection Bias

2. **What does the coefficient of determination ($R^2$) indicate?**

   A. Effectiveness of a model

   B. Independence of variables

   C. Variability explained by a model

   D. Distribution of data

3. **What is meant by collinearity in regression analysis?**

   A. Non-linearity between independent variables

   B. Relationship between dependent and independent variables

   C. When two features exhibit a linear relationship

   D. Multicollinearity issues

4. **In data science, what method is used when a data example can only be classified as a 1 or 0?**

   A. Multiclass classification

   B. Binary classification

   C. Regression

   D. Time series analysis

5. **What is defined as a value that deviates significantly from the main distribution of values?**

   A. Average

   B. Outliers

   C. Median

   D. Standard deviation

6. What process can help in making raw data more understandable and usable?

   A. Transformation

   B. Parsing

   C. Aggregation

   D. Normalizing

7. Which set of statistical parameters is used to measure a distribution?

   A. Median, Mode, Range, Variance

   B. Mean, Variance, Skewness, Kurtosis

   C. Frequency, Probability, Standard Deviation, Variance

   D. Mean, Median, Skewness, Standard Deviation

8. What is the term for a type of data analysis that quantitatively summarizes patterns and relationships in a dataset?

   A. Qualitative analysis

   B. Inferential statistics

   C. Summary statistics

   D. Exploratory data analysis

9. What k-fold cross-validation method uses all data points in the dataset as folds?

   A. LOOCV

   B. Bootstrap Sampling

   C. LPOCV

   D. Stratified K-Fold

10. Data that holds categorical values is referred to as what type of data?

   A. Quantitative data

   B. Descriptive data

   C. Qualitative data

   D. Ordinal data

# **Answers**

**1. B**
**2. C**
**3. C**
**4. B**
**5. B**
**6. B**
**7. B**
**8. C**
**9. C**
**10. C**

# Explanations

## 1. What term is used to describe a model that is deemed useful for its intended task?

A. Overfitting

**B. Skillful**

C. Underfitting

D. Selection Bias

The term used to describe a model that is considered useful for its intended task is "skillful." A skillful model effectively captures the underlying patterns in the data and provides accurate predictions or classifications for new, unseen data. This concept emphasizes the model's ability to generalize well, meaning it performs satisfactorily not just on the training dataset but also on test datasets, reflecting its utility in practical applications. In contrast, overfitting refers to a situation where a model learns the training data too well, including its noise and outliers, which results in poor performance on new data. Underfitting occurs when a model is too simplistic to capture the data's underlying structure, leading to poor performance on both training and test datasets. Selection bias is a situation where certain individuals or groups are systematically excluded from the sample, leading to skewed or invalid results. Thus, "skillful" is the term that accurately reflects a model's effectiveness for the tasks it is designed to address.

## 2. What does the coefficient of determination ($R^2$) indicate?

A. Effectiveness of a model

B. Independence of variables

**C. Variability explained by a model**

D. Distribution of data

The coefficient of determination, commonly denoted as $R^2$, is a statistical measure that offers insight into how well a regression model explains and predicts future outcomes. Specifically, $R^2$ quantifies the proportion of the variance in the dependent variable that can be attributed to the independent variables within the model. When $R^2$ is expressed as a percentage, it reflects the extent to which the model accounts for the variability observed in the data. For example, an $R^2$ value of 0.85 would indicate that 85% of the variance in the dependent variable can be explained by the independent variables in the model, suggesting a strong relationship between them. Thus, $R^2$ serves as a key indicator of the model's explanatory power, allowing analysts to assess the effectiveness of the model in capturing the underlying patterns of the data. The other options, while relevant to data analysis, do not accurately describe the specific purpose of $R^2$. The effectiveness of a model can depend on various other factors, such as model assumptions and validation metrics. The independence of variables relates more closely to multicollinearity concerns in regression analysis, and distribution of data refers to how data points are spread or clustered around a central value, which is not what $R^2$

## 3. What is meant by collinearity in regression analysis?

**A. Non-linearity between independent variables**

**B. Relationship between dependent and independent variables**

**C. When two features exhibit a linear relationship**

**D. Multicollinearity issues**

Collinearity in regression analysis refers to the scenario where two or more features (independent variables) exhibit a strong linear relationship with each other. This means that one independent variable can be predicted from the other(s) with a high degree of accuracy. In practice, this can create challenges in regression, as it complicates the estimation of the individual effects of each feature on the dependent variable. When collinearity is present, it becomes difficult to assess the contribution of each independent variable to the prediction because their effects are intertwined. Understanding collinearity is crucial for data scientists and statisticians as it can lead to inflated standard errors of the coefficients, making hypothesis tests unreliable. While this term is often mentioned in relation to multicollinearity, which refers specifically to the situation where multiple independent variables are correlated, collinearity broadly addresses the linear relationship between two features. In contrast, non-linearity between independent variables does not directly relate to collinearity, as collinearity specifically involves linear relationships. The relationship between dependent and independent variables pertains more to the overall regression model rather than the interrelationships of the independent variables themselves. Multicollinearity issues are a specific scenario under the umbrella of collinearity but pertain more to situations where three

## 4. In data science, what method is used when a data example can only be classified as a 1 or 0?

**A. Multiclass classification**

**B. Binary classification**

**C. Regression**

**D. Time series analysis**

Binary classification is the method used when a data example can be categorized into one of two distinct classes, typically represented as 0 or 1. This type of classification is essential for scenarios where the outcome can only fall into one of two possible categories, such as whether an email is spam or not, or whether a patient has a certain disease. In binary classification, algorithms are trained on labeled data to understand the relationship between the input features and the binary outcome. The model outputs a probability that is then mapped to one of the two categories based on a chosen threshold, often set at 0.5. The other options refer to different approaches not suited for binary outcomes. Multiclass classification is used when there are more than two classes to choose from. Regression is aimed at predicting continuous outcomes rather than categorical ones. Time series analysis deals with data points indexed in time order and is utilized for forecasting rather than binary classification tasks. Thus, binary classification is the appropriate term for this context.

## 5. What is defined as a value that deviates significantly from the main distribution of values?

**A. Average**

**B. Outliers**

**C. Median**

**D. Standard deviation**

The term that describes a value that deviates significantly from the main distribution of values is outliers. Outliers are observations that fall far away from the overall pattern of data, often located in the tails of the distribution. They can arise from variability in the data, measurement errors, or other external factors that do not fit within the general trend of the dataset. Identifying outliers is crucial in data analysis because they can significantly influence statistical measures, such as the mean and standard deviation, leading to potential misinterpretation of the data. In many data analyses, outliers may be investigated further to determine their cause and whether they should be included or excluded from analysis, as they provide insight into data quality and underlying processes. The other options, such as average, median, and standard deviation, relate to central tendency and dispersion measures rather than representing deviations from the norm. While the average and median describe typical values around which data might cluster, and standard deviation measures the spread or dispersion of data points, they do not specifically refer to extreme values or deviations like outliers do.

## 6. What process can help in making raw data more understandable and usable?

**A. Transformation**

**B. Parsing**

**C. Aggregation**

**D. Normalizing**

Parsing is a crucial process in data handling that involves breaking down raw data into more manageable and understandable components. It refers to the technique of analyzing data structures and extracting useful pieces of information based on predetermined rules. For instance, when dealing with a complex dataset, parsing can simplify it by separating individual data fields, allowing for easier analysis and interpretation. This process lays the groundwork for further data manipulation and is essential for cleaning and organizing data so it can be effectively used for analysis. By transforming raw data into a structured format, parsing enhances the usability of the data and prepares it for further processes such as aggregation or normalization. All these factors contribute to making the raw data more understandable for data scientists and analysts, ultimately leading to better insights and decisions. While transformation, aggregation, and normalizing are all important concepts within the data preparation domain, parsing specifically focuses on the initial breakdown of raw data, making it the most direct answer to the question posed.

## 7. Which set of statistical parameters is used to measure a distribution?

A. Median, Mode, Range, Variance

**B. Mean, Variance, Skewness, Kurtosis**

C. Frequency, Probability, Standard Deviation, Variance

D. Mean, Median, Skewness, Standard Deviation

**The set of statistical parameters that accurately measures a distribution includes the mean, variance, skewness, and kurtosis.   The mean provides a measure of central tendency, illustrating the average of the data points in the distribution. Variance indicates how spread out the values are around the mean, giving an insight into the variability of the data. Skewness assesses the asymmetry of the distribution, revealing whether the data leans towards one side (left or right) of the distribution. Lastly, kurtosis measures the "tailedness" of the probability distribution, indicating the presence of outliers and how peaked the distribution is compared to a normal distribution.  Together, these parameters offer a comprehensive overview of a distribution's shape, center, and spread, making them essential for effective statistical analysis.   Other options either include parameters that are not primarily used to quantify distribution properties (such as range or frequency) or do not include a complete set that captures both central tendency and the shape of the distribution effectively.**

## 8. What is the term for a type of data analysis that quantitatively summarizes patterns and relationships in a dataset?

A. Qualitative analysis

B. Inferential statistics

**C. Summary statistics**

D. Exploratory data analysis

**The term for a type of data analysis that quantitatively summarizes patterns and relationships in a dataset is indeed summary statistics. This approach involves calculating measures such as mean, median, mode, variance, and standard deviation, which provide concise information about the central tendency, dispersion, and overall distribution of the data. Summary statistics allow analysts to understand the general characteristics of the dataset in a straightforward manner.  While exploratory data analysis encompasses a broader approach that includes visualizations and various techniques to uncover underlying patterns, summary statistics specifically focus on quantifying those relationships and patterns. Both inferential statistics and qualitative analysis serve different purposes; inferential statistics is used to draw conclusions about a population based on sample data, and qualitative analysis involves non-numerical data and interpretations. Therefore, summary statistics is the most appropriate term for the type of quantitative analysis being described in the question.**

## 9. What k-fold cross-validation method uses all data points in the dataset as folds?

### A. LOOCV

### B. Bootstrap Sampling

### C. LPOCV

### D. Stratified K-Fold

The correct answer is LOOCV, which stands for Leave-One-Out Cross-Validation. This method involves using all but one data point in the dataset as the training set and the single remaining data point as the validation set. This process is repeated so that each data point in the dataset serves as the validation set once, effectively creating as many folds as there are data points in the dataset. LOOCV is particularly beneficial in situations where the dataset is small, as it maximizes the training data for each iteration, thus providing a better understanding of the model's performance by making full use of the available data. This results in a very reliable estimation of the model's predictive performance. In contrast, bootstrap sampling typically involves resampling data with replacement to create multiple training sets and leave out various portions of the data without a fixed fold approach. Stratified K-Fold is a method that divides the dataset into k folds, ensuring that each fold has a representative distribution of the target variable, but it does not use every data point as a fold. LPOCV (Leave-P-Out Cross-Validation) is another variation like LOOCV but removes p data points at a time for validation, rather than just one, thus it does not utilize all data

## 10. Data that holds categorical values is referred to as what type of data?

### A. Quantitative data

### B. Descriptive data

### C. Qualitative data

### D. Ordinal data

Data that holds categorical values is referred to as qualitative data. This type of data is characterized by non-numeric categories or labels that are used to describe characteristics or traits. For instance, data representing colors, names, or types of animals falls into this category because it does not measure quantities but rather classifies or describes a quality or characteristic. Qualitative data can further be categorized into nominal and ordinal data. Nominal data has no specific order (e.g., types of fruits), while ordinal data has a defined order or ranking (e.g., survey responses like "satisfied," "neutral," or "dissatisfied"). However, the primary characteristic that distinguishes qualitative data is that it deals with categories rather than numbers. In contrast, quantitative data pertains to numerical values and can be measured, allowing for calculations and statistical analysis. Descriptive data generally refers to data used to summarize or describe various attributes, rather than specifying a type of data. Therefore, for data that consists of categorical values, qualitative data is the appropriate term.