

Quiz Questions

O'Reilly - Statistical Modeling and Inference with Python

Dr. Chester Ismay

August 2024

Week 1 Module 1

W1M1 - Question 1

What is the primary purpose of linear regression in predictive modeling?

- ▶ A) To determine the maximum value of a dataset.
- ▶ B) To fit a line that best represents the relationship between the independent variable and the dependent variable.
- ▶ C) To classify data into distinct categories.
- ▶ D) To find the median of the dependent variable.
- ▶ E) To sort data in ascending order.

W1M1 - Question 2

Which of the following best describes the dependent variable in the context of linear regression?

- ▶ A) It is the variable we manipulate to observe changes.
- ▶ B) It is always a binary variable.
- ▶ C) It is the outcome we are trying to predict or explain.
- ▶ D) It is the variable that remains constant.
- ▶ E) It is the mean of all observations.

W1M1 - Question 3

What assumption ensures that the residuals (differences between observed and predicted values) are normally distributed?

- ▶ A) Linearity
- ▶ B) Independence
- ▶ C) Normality
- ▶ D) Homoscedasticity
- ▶ E) Multicollinearity

W1M1 - Question 4

Why is it important to split the dataset into training and testing sets in machine learning?

- ▶ A) To reduce the size of the dataset.
- ▶ B) To test the model's performance on new, unseen data.
- ▶ C) To make the dataset more manageable.
- ▶ D) To increase the accuracy of predictions on the training data.
- ▶ E) To ensure all data points are used in training.

W1M1 - Question 5

Which assumption of linear regression requires that the variance of residuals is constant across all levels of the independent variable?

- ▶ A) Linearity
- ▶ B) Independence
- ▶ C) Normality
- ▶ D) Homoscedasticity
- ▶ E) Multicollinearity

W1M1 - Question 6

What does the assumption of linearity in linear regression imply about the relationship between the independent variable (x) and the dependent variable (y)?

- ▶ A) The relationship is nonlinear.
- ▶ B) The relationship is a straight line.
- ▶ C) The relationship involves multiple curves.
- ▶ D) The relationship is random.
- ▶ E) The relationship is quadratic.

W1M1 - Question 7

What does it mean if the residuals of a linear regression model are not normally distributed?

- ▶ A) The independent variable is not normally distributed.
- ▶ B) The dependent variable is not normally distributed.
- ▶ C) The model may not provide reliable statistical inferences.
- ▶ D) The variance of residuals is constant.
- ▶ E) The model is perfectly accurate.

W1M1 - Question 8

Which of the following best describes overfitting in the context of machine learning?

- ▶ A) The model performs well on both training and testing data.
- ▶ B) The model performs exceptionally well on training data but poorly on testing data.
- ▶ C) The model performs poorly on training data but well on testing data.
- ▶ D) The model performs equally well on all types of data.
- ▶ E) The model ignores the training data.

W1M1 - Question 9

What is the primary goal of minimizing the difference between observed and predicted values in linear regression?

- ▶ A) To maximize the difference between observed and predicted values.
- ▶ B) To reduce the dataset size.
- ▶ C) To find the optimal parameters for the regression line.
- ▶ D) To ensure the residuals are independent.
- ▶ E) To classify the data into categories.

W1M1 - Question 10

Which of the following best illustrates the concept of train-test splitting?

- ▶ A) Using the entire dataset for both training and testing.
- ▶ B) Splitting the dataset into two parts, one for training the model and one for testing its performance.
- ▶ C) Using only a small part of the dataset for testing.
- ▶ D) Randomly shuffling the data without any split.
- ▶ E) Ignoring the testing phase.

Week 1 Module 2

W1M2 - Question 1

What does a correlation coefficient of 0.8 indicate about the relationship between two variables?

- ▶ A) They have a weak negative correlation.
- ▶ B) They have a moderate positive correlation.
- ▶ C) They have a strong positive correlation.
- ▶ D) They have no correlation.
- ▶ E) They have a strong negative correlation.

W1M2 - Question 2

Which correlation coefficient value indicates that there is no recognizable pattern in the relationship between two variables?

- ▶ A) -1
- ▶ B) -0.5
- ▶ C) 0
- ▶ D) 0.5
- ▶ E) 1

W1M2 - Question 3

What type of correlation is observed when one variable increases while the other decreases?

- ▶ A) Positive correlation
- ▶ B) Negative correlation
- ▶ C) Zero correlation
- ▶ D) Strong correlation
- ▶ E) Moderate correlation

W1M2 - Question 4

How would you interpret a correlation coefficient of -0.9 ?

- ▶ A) Strong positive correlation
- ▶ B) Moderate positive correlation
- ▶ C) Strong negative correlation
- ▶ D) Weak negative correlation
- ▶ E) Zero correlation

W1M2 - Question 5

What does the absolute value of a correlation coefficient ($|r|$) tell us about the relationship between two variables?

- ▶ A) The direction of the relationship
- ▶ B) The strength of the relationship
- ▶ C) The cause of the relationship
- ▶ D) The independence of the relationship
- ▶ E) The variance of the relationship

W1M2 - Question 6

Which statement is true about the limitations of correlation coefficients?

- ▶ A) They capture non-linear relationships accurately.
- ▶ B) They are not affected by outliers.
- ▶ C) They only capture linear relationships.
- ▶ D) They measure the causality between variables.
- ▶ E) They are unaffected by the scale of the data.

W1M2 - Question 7

Why is it important to be cautious of outliers when interpreting correlation coefficients?

- ▶ A) Outliers do not affect correlation coefficients.
- ▶ B) Outliers can make the relationship appear weaker than it is.
- ▶ C) Outliers can make the relationship appear stronger or weaker than it is.
- ▶ D) Outliers indicate a strong relationship.
- ▶ E) Outliers are always removed from the data.

W1M2 - Question 8

If two variables have a correlation coefficient of -0.2 , how would you describe their relationship?

- ▶ A) Strong negative correlation
- ▶ B) Weak negative correlation
- ▶ C) Moderate positive correlation
- ▶ D) Strong positive correlation
- ▶ E) Zero correlation

W1M2 - Question 9

What is the range of possible values for a correlation coefficient?

- ▶ A) 0 to 1
- ▶ B) -1 to 0
- ▶ C) -2 to 2
- ▶ D) -1 to 1
- ▶ E) -0.5 to 0.5

W1M2 - Question 10

How would you interpret a correlation coefficient of 0.4?

- ▶ A) Strong positive correlation
- ▶ B) Moderate positive correlation
- ▶ C) Weak positive correlation
- ▶ D) Zero correlation
- ▶ E) Strong negative correlation

Week 1 Module 3

W1M3 - Question 1

What is the primary purpose of multiple regression in statistical modeling?

- ▶ A) To determine the median value of a dataset.
- ▶ B) To fit a line that represents the relationship between one independent variable and one dependent variable.
- ▶ C) To understand the combined effect of multiple independent variables on a single dependent variable.
- ▶ D) To classify data into distinct categories.
- ▶ E) To find the mode of the dependent variable.

W1M3 - Question 2

Which of the following best describes multiple regression?

- ▶ A) A statistical method to measure the relationship between one dependent variable and one independent variable.
- ▶ B) A technique used to classify data into different categories.
- ▶ C) A method to understand how several independent variables together influence a single dependent variable.
- ▶ D) A process to calculate the median value of multiple variables.
- ▶ E) A way to find the mode of a single dependent variable.

W1M3 - Question 3

Which assumption of multiple regression ensures that the residuals have a constant variance across all levels of the independent variables?

- ▶ A) Linearity
- ▶ B) Independence
- ▶ C) Normality
- ▶ D) Error homoscedasticity
- ▶ E) Non-multicollinearity

W1M3 - Question 4

Why is it important to check for multicollinearity in a multiple regression model?

- ▶ A) To ensure the residuals are normally distributed.
- ▶ B) To confirm the dependent variable is linear.
- ▶ C) To avoid redundancy and maintain the model's integrity.
- ▶ D) To ensure residuals have constant variance.
- ▶ E) To increase the complexity of the model.

W1M3 - Question 5

Which of the following best describes the difference between simple and multiple regression?

- ▶ A) Simple regression involves one dependent variable, while multiple regression involves multiple dependent variables.
- ▶ B) Simple regression involves one independent variable, while multiple regression involves multiple independent variables.
- ▶ C) Simple regression models non-linear relationships, while multiple regression models linear relationships.
- ▶ D) Simple regression is more complex than multiple regression.
- ▶ E) Simple regression always has higher accuracy than multiple regression.

W1M3 - Question 6

What is the significance of the coefficient β in a multiple regression model?

- ▶ A) It measures the correlation between the independent variables.
- ▶ B) It represents the mean of the dependent variable.
- ▶ C) It indicates the standard deviation of the residuals.
- ▶ D) It quantifies the effect of an independent variable on the dependent variable, holding other variables constant.
- ▶ E) It calculates the mode of the independent variables.

W1M3 - Question 7

In a study analyzing the factors that influence the academic performance of high school students, a researcher collects data on several variables including hours of study per week, attendance rate, extracurricular activities, and parental involvement. Which of the following would be the most appropriate target variable for a multiple regression analysis in this context?

- ▶ A) Hours of study per week
- ▶ B) Attendance rate
- ▶ C) Extracurricular activities
- ▶ D) Parental involvement
- ▶ E) Overall academic performance (e.g., GPA)

W1M3 - Question 8

What does the assumption of normality in multiple regression imply?

- ▶ A) The independent variables are normally distributed.
- ▶ B) The dependent variable is normally distributed.
- ▶ C) The residuals are normally distributed.
- ▶ D) The independent variables are independent of each other.
- ▶ E) The dependent variable has constant variance.

W1M3 - Question 9

How does multiple regression enhance the ability to make predictions compared to simple regression?

- ▶ A) By reducing the number of variables.
- ▶ B) By simplifying the relationship between variables.
- ▶ C) By incorporating multiple factors, leading to more accurate and comprehensive predictions.
- ▶ D) By ignoring the dependent variable.
- ▶ E) By focusing on non-linear relationships.

W1M3 - Question 10

What does the term “non-multicollinearity” refer to in the context of multiple regression assumptions?

- ▶ A) The residuals have constant variance.
- ▶ B) The relationship between variables is non-linear.
- ▶ C) The independent variables are not highly correlated with each other.
- ▶ D) The dependent variable is normally distributed.
- ▶ E) The residuals are independent.

Week 1 Module 4

W1M4 - Question 1

What is the primary purpose of logistic regression in binary classification problems?

- ▶ A) Predicting continuous outcomes
- ▶ B) Predicting the probability of an outcome belonging to one of two categories
- ▶ C) Reducing dimensionality of datasets
- ▶ D) Estimating the correlation between variables
- ▶ E) Visualizing data using scatter plots

W1M4 - Question 2

Which of the following is an appropriate application of logistic regression?

- ▶ A) Predicting the temperature for the next day
- ▶ B) Determining the likelihood of a patient having a disease (yes/no)
- ▶ C) Forecasting stock prices
- ▶ D) Analyzing the relationship between height and weight
- ▶ E) Reducing the number of variables in a dataset

W1M4 - Question 3

Which assumption must be met for logistic regression to be valid?

- ▶ A) The target variable must be continuous
- ▶ B) The predictors must be categorical
- ▶ C) The observations must be independent of each other
- ▶ D) The target variable must have more than two categories
- ▶ E) The predictors must be normally distributed

W1M4 - Question 4

How does logistic regression differ from linear regression?

- ▶ A) Logistic regression is used for predicting continuous outcomes
- ▶ B) Logistic regression can handle multiple dependent variables
- ▶ C) Logistic regression predicts probabilities of categorical outcomes
- ▶ D) Logistic regression reduces the dimensionality of the dataset
- ▶ E) Logistic regression visualizes linear relationships using scatter plots

W1M4 - Question 5

Why is it important for the predictors in logistic regression to have little or no multicollinearity?

- ▶ A) To ensure the target variable is binary
- ▶ B) To make sure the observations are independent
- ▶ C) To avoid redundancy and ensure reliable predictions
- ▶ D) To maintain a normal distribution of the predictors
- ▶ E) To visualize data using scatter plots

W1M4 - Question 6

Which of the following best describes the dependent/target variable in logistic regression?

- ▶ A) It can take any continuous value
- ▶ B) It must be binary, with two possible outcomes
- ▶ C) It must have three categories
- ▶ D) It is always a categorical variable with more than two levels
- ▶ E) It is used for reducing the dimensionality of the dataset

W1M4 - Question 7

What does the assumption of linearity of independent variables and the logit of the target variable imply in logistic regression?

- ▶ A) The independent variables must be linearly correlated with each other
- ▶ B) There should be a linear relationship between the independent variables and the log-odds of the target variable's success
- ▶ C) The target variable must be normally distributed
- ▶ D) The predictors must have constant variance
- ▶ E) The predictors must be categorical

W1M4 - Question 8

Which of the following is NOT an application of logistic regression?

- ▶ A) Medical diagnosis to predict disease presence
- ▶ B) Credit scoring to determine whether a loan defaults
- ▶ C) Marketing to classify whether to purchase or not
- ▶ D) Forecasting future stock prices
- ▶ E) Spam detection to classify emails into spam or not spam

W1M4 - Question 9

Why is it important to check for the assumption that observations are independent in logistic regression?

- ▶ A) To ensure the target variable is binary
- ▶ B) To reduce the dimensionality of the dataset
- ▶ C) To avoid bias and ensure reliable predictions
- ▶ D) To establish a linear relationship between variables
- ▶ E) To predict continuous outcomes

W1M4 - Question 10

What is the main difference between logistic regression and linear regression in terms of their target variables?

- ▶ A) Logistic regression has a continuous target variable, while linear regression has a binary target variable
- ▶ B) Logistic regression has a binary target variable, while linear regression has a continuous target variable
- ▶ C) Both logistic and linear regression have binary target variables
- ▶ D) Both logistic and linear regression have continuous target variables
- ▶ E) Logistic regression is used for categorical outcomes, while linear regression is used for ordinal outcomes

Week 1 Module 5

W1M5 - Question 1

What is the primary purpose of ANOVA (Analysis of Variance)?

- ▶ A) To compare the means of two groups
- ▶ B) To compare the means of three or more groups
- ▶ C) To reduce the dimensionality of datasets
- ▶ D) To predict continuous outcomes
- ▶ E) To visualize data using scatter plots

W1M5 - Question 2

Which of the following scenarios is most appropriate for using ANOVA?

- ▶ A) Predicting stock prices based on historical data
- ▶ B) Comparing the average heights of students in two different classes
- ▶ C) Analyzing the relationship between temperature and ice cream sales
- ▶ D) Comparing the test scores of students across three different teaching methods
- ▶ E) Determining the correlation between hours of study and exam scores

W1M5 - Question 3

What does a low p-value (typically less than 0.05) in ANOVA indicate?

- ▶ A) The differences in group means are likely due to chance
- ▶ B) The differences in group means are statistically significant
- ▶ C) The variances within each group are different
- ▶ D) The residuals are not normally distributed
- ▶ E) The observations are not independent

W1M5 - Question 4

What is the F-statistic in ANOVA?

- ▶ A) A measure of the total variance in the data
- ▶ B) The ratio of between-group variance to within-group variance
- ▶ C) The probability of observing the data if the null hypothesis is true
- ▶ D) The mean of the dependent variable
- ▶ E) The standard deviation of the residuals

W1M5 - Question 5

Which of the following best describes the assumption related to the distribution of residuals in the context of ANOVA, and why is it important?

- ▶ A) The residuals should have constant variance across all levels of the independent variable, ensuring homoscedasticity.
- ▶ B) The residuals should follow a bell-shaped curve to ensure that the F-statistic follows an F-distribution, making hypothesis testing valid.
- ▶ C) The residuals should show a linear relationship with the independent variable, confirming the linearity assumption.
- ▶ D) The residuals should be independent of each other to prevent autocorrelation, ensuring unbiased results.
- ▶ E) The residuals should be centered around zero to maintain the integrity of the mean comparison.

W1M5 - Question 6

What does the assumption of homogeneity of variances in ANOVA imply?

- ▶ A) The observations are independent
- ▶ B) The variances within each group are approximately equal
- ▶ C) The dependent variable is binary
- ▶ D) The independent variables are not highly correlated
- ▶ E) The residuals are normally distributed

W1M5 - Question 7

Which of the following best describes a one-way ANOVA?

- ▶ A) It compares the means of two groups based on one independent variable
- ▶ B) It compares the means of three or more groups based on one independent variable
- ▶ C) It examines the influence of two different categorical variables on one continuous dependent variable
- ▶ D) It reduces the number of variables in a dataset
- ▶ E) It predicts binary outcomes

W1M5 - Question 8

In the context of ANOVA, what is meant by "between-group variance"?

- ▶ A) Variance within each group
- ▶ B) Total variance observed in the data
- ▶ C) Variance in the mean scores between different groups
- ▶ D) Variance in the residuals of the model
- ▶ E) Variance in the predictors

W1M5 - Question 9

What is the main limitation of ANOVA in identifying which specific groups are different?

- ▶ A) It can only compare two groups at a time
- ▶ B) It is sensitive to violations of assumptions
- ▶ C) It requires normally distributed predictors
- ▶ D) It does not specify which groups are different
- ▶ E) It only works for binary dependent variables

W1M5 - Question 10

Why are post-hoc tests necessary after performing ANOVA?

- ▶ A) To verify the independence of observations
- ▶ B) To test for multicollinearity among predictors
- ▶ C) To check the normality of residuals
- ▶ D) To perform detailed pairwise comparisons between groups
- ▶ E) To predict continuous outcomes

Week 2 Module 1

W2M1 - Question 1

What is the primary advantage of non-parametric tests over parametric tests?

- ▶ A) They assume a specific distribution for the data.
- ▶ B) They are only applicable to normally distributed data.
- ▶ C) They can be used when data does not meet the assumptions of parametric tests.
- ▶ D) They require larger sample sizes.
- ▶ E) They are always more accurate than parametric tests.

W2M1 - Question 2

Which non-parametric test is a suitable alternative to the one-way ANOVA?

- ▶ A) Mann-Whitney U Test
- ▶ B) Kruskal-Wallis Test
- ▶ C) Wilcoxon Signed-Rank Test
- ▶ D) Chi-Square Test
- ▶ E) Spearman's Rank Correlation

W2M1 - Question 3

When should you use the Mann-Whitney U test?

- ▶ A) When comparing the means of three or more groups
- ▶ B) When comparing the medians of two independent groups with non-normal distributions
- ▶ C) When testing for correlation between two variables
- ▶ D) When comparing paired samples
- ▶ E) When analyzing categorical data

W2M1 - Question 4

In which scenario would the Kruskal-Wallis test be appropriate?

- ▶ A) Comparing the effectiveness of two marketing strategies on sales
- ▶ B) Comparing customer satisfaction ratings between two stores
- ▶ C) Comparing the effectiveness of different diets on weight loss across multiple groups with non-normally distributed data
- ▶ D) Testing the relationship between two continuous variables
- ▶ E) Comparing the before and after effects of a treatment on the same group of participants

W2M1 - Question 5

What is one of the main benefits of non-parametric tests?

- ▶ A) They are less sensitive to outliers and skewed data
- ▶ B) They assume the data is normally distributed
- ▶ C) They require equal variances among groups
- ▶ D) They are more complex to perform and interpret
- ▶ E) They are limited to small sample sizes

W2M1 - Question 6

Which of the following is NOT an assumption of ANOVA that can lead to the use of non-parametric tests when violated?

- ▶ A) Independence of observations
- ▶ B) Homogeneity of variances
- ▶ C) Normally distributed residuals
- ▶ D) Linear relationship between variables
- ▶ E) Equal sample sizes in each group

W2M1 - Question 7

Which of the following best describes the type of data suitable for non-parametric tests?

- ▶ A) Data that is normally distributed with equal variances
- ▶ B) Ordinal data or data that does not meet parametric test assumptions
- ▶ C) Data with a large sample size
- ▶ D) Data that is highly linear and homoscedastic
- ▶ E) Binary data with equal group sizes

W2M1 - Question 8

Why might you choose a non-parametric test over a parametric test?

- ▶ A) Because they are always more powerful
- ▶ B) Because they can handle larger datasets more efficiently
- ▶ C) Because they do not assume a specific distribution for the data
- ▶ D) Because they are simpler to calculate
- ▶ E) Because they only require categorical data

W2M1 - Question 9

Which test would you use to compare customer satisfaction ratings (not normally distributed) between two different stores?

- ▶ A) One-way ANOVA
- ▶ B) Two-way ANOVA
- ▶ C) Kruskal-Wallis Test
- ▶ D) Mann-Whitney U Test
- ▶ E) Paired t-test

W2M1 - Question 10

Which of the following is a characteristic of the Kruskal-Wallis test?

- ▶ A) It compares the means of three or more groups
- ▶ B) It requires normally distributed data
- ▶ C) It is less affected by outliers compared to one-way ANOVA
- ▶ D) It only compares two groups at a time
- ▶ E) It requires equal sample sizes in each group

Week 2 Module 2

W2M2 - Question 1

What is Spearman's Rank Correlation primarily used for?

- ▶ A) Measuring the strength and direction of a linear relationship between two variables
- ▶ B) Measuring the strength and direction of a monotonic relationship between two variables
- ▶ C) Reducing the dimensionality of datasets
- ▶ D) Predicting categorical outcomes
- ▶ E) Analyzing normally distributed data

W2M2 - Question 2

When is Kendall's Tau Correlation preferred over Spearman's Rank Correlation?

- ▶ A) When the data is normally distributed
- ▶ B) When there are many tied ranks in the dataset
- ▶ C) When analyzing continuous data
- ▶ D) When reducing dimensionality
- ▶ E) When predicting binary outcomes

W2M2 - Question 3

Which non-parametric test is ideal for ordinal data, such as ranking employees based on performance levels?

- ▶ A) Pearson's Correlation
- ▶ B) Linear Regression
- ▶ C) Spearman's Rank Correlation
- ▶ D) Chi-Square Test
- ▶ E) ANOVA

W2M2 - Question 4

What is the primary advantage of using the Theil-Sen Estimator in linear regression?

- ▶ A) It assumes normal distribution of errors
- ▶ B) It minimizes the impact of outliers
- ▶ C) It is simpler to calculate than ordinary least squares regression
- ▶ D) It requires normally distributed residuals
- ▶ E) It can only be used with categorical data

W2M2 - Question 5

In which scenario would you use Spearman's Rank Correlation instead of Pearson's Correlation?

- ▶ A) When the data is normally distributed and linear
- ▶ B) When the data has many outliers and is not normally distributed
- ▶ C) When you need to predict binary outcomes
- ▶ D) When analyzing categorical data
- ▶ E) When performing dimensionality reduction

W2M2 - Question 6

What type of relationship does Spearman's Rank Correlation measure?

- ▶ A) Linear relationship
- ▶ B) Exponential relationship
- ▶ C) Monotonic relationship
- ▶ D) Quadratic relationship
- ▶ E) Cubic relationship

W2M2 - Question 7

Why is the Theil-Sen Estimator useful for datasets with non-normal error distributions?

- ▶ A) It assumes a specific distribution for the data
- ▶ B) It remains effective without the assumption of normality
- ▶ C) It reduces the dimensionality of the dataset
- ▶ D) It predicts categorical outcomes
- ▶ E) It requires normally distributed predictors

W2M2 - Question 8

Which of the following scenarios is appropriate for using Kendall's Tau Correlation?

- ▶ A) Comparing the average test scores of students in three different classes
- ▶ B) Measuring the correlation between income and expenditure
- ▶ C) Assessing the relationship between ranks of competitors in multiple races with many tied ranks
- ▶ D) Predicting the likelihood of customer churn
- ▶ E) Analyzing the frequency of different categories in a survey

W2M2 - Question 9

What distinguishes Kendall's Tau Correlation from Spearman's Rank Correlation?

- ▶ A) Kendall's Tau is only used for binary data
- ▶ B) Kendall's Tau specifically accounts for the number of concordant and discordant pairs
- ▶ C) Kendall's Tau requires normally distributed data
- ▶ D) Kendall's Tau is a parametric test
- ▶ E) Kendall's Tau is simpler to calculate than Spearman's

W2M2 - Question 10

In the context of non-parametric correlation measures, what does a high positive value of Spearman's Rank Correlation indicate?

- ▶ A) A strong negative linear relationship
- ▶ B) A weak negative monotonic relationship
- ▶ C) A strong positive monotonic relationship
- ▶ D) No relationship between the variables
- ▶ E) A weak positive linear relationship

Week 2 Module 3

W2M3 - Question 1

What is the primary purpose of bootstrapping in statistical analysis?

- ▶ A) To assume a normal distribution for the data
- ▶ B) To reduce the sample size of the dataset
- ▶ C) To estimate the distribution of a statistic by resampling with replacement
- ▶ D) To predict future outcomes
- ▶ E) To perform dimensionality reduction

W2M3 - Question 2

Which of the following steps is NOT involved in the bootstrapping methodology?

- ▶ A) Resampling with replacement
- ▶ B) Calculating the statistic for each resample
- ▶ C) Generating confidence intervals
- ▶ D) Assuming a specific distribution for the population
- ▶ E) Repeating the resampling process many times

W2M3 - Question 3

What is the significance of generating resamples in the bootstrapping process?

- ▶ A) To reduce the computational complexity
- ▶ B) To assume normality in the data
- ▶ C) To create multiple hypothetical scenarios based on the original data
- ▶ D) To combine multiple datasets
- ▶ E) To ensure data normalization

W2M3 - Question 4

Why is bootstrapping considered a flexible method in statistical analysis?

- ▶ A) It reduces the dimensionality of the data
- ▶ B) It can be used for any type of data distribution
- ▶ C) It requires less computational power
- ▶ D) It provides exact predictions
- ▶ E) It always produces normal distributions

W2M3 - Question 5

What does resampling with replacement mean in the context of bootstrapping?

- ▶ A) Drawing samples without putting them back into the dataset
- ▶ B) Drawing samples and putting them back into the dataset before drawing again
- ▶ C) Drawing samples only once
- ▶ D) Replacing the original dataset with a new dataset
- ▶ E) Drawing samples without considering the sample size

W2M3 - Question 6

How does bootstrapping help in deriving confidence intervals for a statistic?

- ▶ A) By assuming the data follows a normal distribution
- ▶ B) By resampling with replacement to create a distribution of the statistic
- ▶ C) By reducing the sample size
- ▶ D) By performing linear regression on the data
- ▶ E) By normalizing the data

W2M3 - Question 7

In bootstrapping, why is it important to repeat the resampling and calculation process many times?

- ▶ A) To increase the sample size
- ▶ B) To ensure the statistic is normally distributed
- ▶ C) To build a robust estimate of the statistic's sampling distribution
- ▶ D) To reduce the variability in the data
- ▶ E) To perform dimensionality reduction

W2M3 - Question 8

Which of the following is a key advantage of using bootstrapping over traditional methods?

- ▶ A) It requires fewer computational resources
- ▶ B) It makes fewer assumptions about the data's underlying distribution
- ▶ C) It always produces more accurate results
- ▶ D) It reduces the need for large sample sizes
- ▶ E) It simplifies data preprocessing

W2M3 - Question 9

What does the distribution of resampled statistics represent in bootstrapping?

- ▶ A) The actual population distribution
- ▶ B) The variability of the statistic across different samples
- ▶ C) The mean of the original dataset
- ▶ D) The frequency distribution of the original data
- ▶ E) The standard deviation of the original dataset

W2M3 - Question 10

Which step in the bootstrapping process involves drawing samples with replacement from the original dataset?

- ▶ A) Calculation
- ▶ B) Normalization
- ▶ C) Resampling
- ▶ D) Prediction
- ▶ E) Aggregation

Week 2 Module 4

W2M4 - Question 1

What is the primary purpose of constructing confidence intervals in statistical analysis?

- ▶ A) To predict future outcomes
- ▶ B) To estimate the precision and uncertainty of a sample statistic
- ▶ C) To reduce the sample size of the dataset
- ▶ D) To assume a normal distribution for the data
- ▶ E) To perform dimensionality reduction

W2M4 - Question 2

Which of the following best describes a 95

- ▶ A) The interval will always contain the true population parameter
- ▶ B) The interval is guaranteed to contain the true parameter in any one sample
- ▶ C) Approximately 95 out of 100 such intervals will contain the true population parameter if the sampling process is repeated
- ▶ D) The interval contains 95
- ▶ E) The interval reduces the variability of the data

W2M4 - Question 3

What are the steps involved in constructing confidence intervals using bootstrapping?

1. Generate resamples from the original dataset with replacement.
2. Determine the percentiles of the resampled statistics.
3. Calculate the statistic of interest for each resample.

- ▶ A) 1, 2, 3
- ▶ B) 1, 3, 2
- ▶ C) 2, 1, 3
- ▶ D) 2, 3, 1
- ▶ E) 3, 2, 1

W2M4 - Question 4

What role does the Central Limit Theorem (CLT) play in constructing confidence intervals?

- ▶ A) It assumes the data is normally distributed
- ▶ B) It states that the sample mean will be equal to the population mean
- ▶ C) It justifies using a normal distribution for the sampling distribution of the sample mean as the sample size increases
- ▶ D) It reduces the sample size needed for accurate estimates
- ▶ E) It eliminates the need for resampling

W2M4 - Question 5

Why is the margin of error important in the context of confidence intervals?

- ▶ A) It reduces the sample size needed for analysis
- ▶ B) It represents the range within which the true population parameter is expected to lie
- ▶ C) It ensures the data follows a normal distribution
- ▶ D) It simplifies the calculation of the sample mean
- ▶ E) It determines the mode of the dataset

W2M4 - Question 6

How are the percentiles used in constructing a 95

- ▶ A) The 5th and 95th percentiles are used
- ▶ B) The 25th and 75th percentiles are used
- ▶ C) The 2.5th and 97.5th percentiles are used
- ▶ D) The 1st and 99th percentiles are used
- ▶ E) The 10th and 90th percentiles are used

W2M4 - Question 7

What does the confidence level of an interval indicate?

- ▶ A) The percentage of sample data points within the interval
- ▶ B) The probability that the interval contains the sample mean
- ▶ C) The degree of certainty that the interval contains the true population parameter
- ▶ D) The range of the sample data
- ▶ E) The mean of the resampled statistics

W2M4 - Question 8

In bootstrapping, why is it important to perform many resampling iterations (e.g., thousands)?

- ▶ A) To ensure the data becomes normally distributed
- ▶ B) To reduce the sample size needed for analysis
- ▶ C) To build a robust estimate of the statistic's sampling distribution
- ▶ D) To eliminate the need for calculating percentiles
- ▶ E) To ensure all original data points are used equally

W2M4 - Question 9

What does the Central Limit Theorem (CLT) state about the distribution of the sample mean?

- ▶ A) It follows a binomial distribution
- ▶ B) It approximates a normal distribution as the sample size increases
- ▶ C) It is always normally distributed regardless of sample size
- ▶ D) It is uniformly distributed
- ▶ E) It is unaffected by the sample size

W2M4 - Question 10

Why are confidence intervals preferred over point estimates in statistical analysis?

- ▶ A) They provide a single value estimate
- ▶ B) They assume the data is normally distributed
- ▶ C) They offer a range of values, reflecting the uncertainty and variability of the estimate
- ▶ D) They simplify the calculation process
- ▶ E) They eliminate the need for large sample sizes

Week 2 Module 5

W2M5 - Question 1

In the context of finance, how can bootstrapping be used for risk assessment?

- ▶ A) By assuming a normal distribution for asset returns
- ▶ B) By reducing the sample size of the dataset
- ▶ C) By sampling historical return data to generate a distribution of potential portfolio values
- ▶ D) By predicting future stock prices with certainty
- ▶ E) By simplifying the portfolio allocation process

W2M5 - Question 2

Which of the following is an application of bootstrapping in healthcare?

- ▶ A) Predicting future sales trends
- ▶ B) Estimating confidence intervals for treatment effects in clinical trials
- ▶ C) Measuring the effectiveness of marketing campaigns
- ▶ D) Reducing the dimensionality of patient data
- ▶ E) Performing linear regression on medical records

W2M5 - Question 3

How can bootstrapping be applied in marketing to evaluate the success of a campaign?

- ▶ A) By predicting future stock prices
- ▶ B) By reducing the variability in sales data
- ▶ C) By sampling conversion rate data to construct confidence intervals for key performance metrics
- ▶ D) By assuming a normal distribution for customer behavior
- ▶ E) By simplifying the calculation of customer lifetime value

W2M5 - Question 4

What is one of the key benefits of using bootstrapping in statistical analysis?

- ▶ A) It reduces the need for large sample sizes
- ▶ B) It assumes a specific distribution for the data
- ▶ C) It always provides more accurate results
- ▶ D) It simplifies the data preprocessing steps
- ▶ E) It provides exact predictions

W2M5 - Question 5

How is bootstrapping used in finance for portfolio optimization?

- ▶ A) By assuming that historical returns will repeat exactly in the future
- ▶ B) By reducing the number of assets in the portfolio
- ▶ C) By sampling historical returns to evaluate the performance of different portfolio compositions
- ▶ D) By predicting the exact future performance of each asset
- ▶ E) By simplifying the calculation of the average return

W2M5 - Question 6

In healthcare, how can bootstrapping enhance survival analysis?

- ▶ A) By assuming normal distribution of survival times
- ▶ B) By reducing the number of patients needed in a study
- ▶ C) By extracting patient survival times to estimate confidence intervals for survival rates
- ▶ D) By predicting the exact lifespan of each patient
- ▶ E) By simplifying the collection of patient data

W2M5 - Question 7

Why is bootstrapping particularly useful when dealing with small sample sizes?

- ▶ A) It requires fewer computational resources
- ▶ B) It provides exact predictions
- ▶ C) It helps to estimate variability and confidence intervals accurately without relying on large datasets
- ▶ D) It simplifies data analysis
- ▶ E) It assumes a normal distribution for the data

W2M5 - Question 8

Which of the following best describes the application of bootstrapping in marketing for customer segmentation?

- ▶ A) Predicting future stock prices
- ▶ B) Using demographic data to understand the variability and reliability of customer segments
- ▶ C) Simplifying the calculation of conversion rates
- ▶ D) Reducing the dimensionality of customer data
- ▶ E) Assuming a normal distribution for customer behavior

W2M5 - Question 9

How does bootstrapping help in estimating Value at Risk (VaR) for investment portfolios in finance?

- ▶ A) By assuming normal distribution of asset returns
- ▶ B) By using historical return data to generate a distribution of potential portfolio values
- ▶ C) By predicting future stock prices with certainty
- ▶ D) By simplifying the calculation of portfolio returns
- ▶ E) By reducing the number of assets in the portfolio

W2M5 - Question 10

What is a significant advantage of using bootstrapping over traditional parametric methods?

- ▶ A) It assumes a normal distribution for the data
- ▶ B) It provides exact predictions
- ▶ C) It does not require the data to follow a specific distribution
- ▶ D) It reduces the computational complexity
- ▶ E) It simplifies data collection