

Applied Data Science Questions

What is the process of cleaning, transforming, and enriching raw data into a usable form for analysis?

- a) Data Wrangling
- b) Data Visualization
- c) Data Mining
- d) Data Analysis

Answer: a) Data Wrangling

Which programming language is widely used for data science and machine learning?

- a) Java
- b) C++
- c) Python
- d) Ruby

Answer: c) Python

What type of machine learning algorithm is used for predicting numeric values, such as house prices?

- a) Classification
- b) Clustering
- c) Regression
- d) Reinforcement Learning

Answer: c) Regression

Which technique is used to find the best set of hyperparameters for a machine learning model?

- a) Gradient Descent
- b) Feature Engineering
- c) Hyperparameter Optimization
- d) K-nearest Neighbors

Answer: c) Hyperparameter Optimization

What is the term used for machine learning algorithms that learn from labeled data to make predictions or decisions?

- a) Unsupervised Learning
- b) Reinforcement Learning
- c) Supervised Learning
- d) Semi-Supervised Learning

Answer: c) Supervised Learning

Which evaluation metric is commonly used for classification problems and represents the ratio of correctly predicted instances to the total instances?

- a) Mean Squared Error (MSE)
- b) Area Under the Curve (AUC)
- c) F1 Score
- d) Accuracy

Answer: d) Accuracy

Which cloud platform is known for providing AI services like natural language processing and computer vision?

- a) Google Cloud Platform
- b) Microsoft Azure
- c) IBM Cloud
- d) Amazon Web Services (AWS)

Answer: c) IBM Cloud

Which module in data science deals with analyzing and interpreting data to extract useful insights?

- a) Data Visualization
- b) Data Wrangling Techniques
- c) Model Evaluation Metrics
- d) Hyper-parameter Optimization

Answer: a) Data Visualization

Which module in data science focuses on finding patterns and relationships in data without using labeled examples?

- a) Unsupervised Learning
- b) Supervised Learning - Classification
- c) Model Evaluation Metrics
- d) Hyper-parameter Optimization

Answer: a) Unsupervised Learning

Which technique is used to replace missing values in a dataset with appropriate values?

- a) Outlier Detection
- b) Feature Scaling
- c) Data Imputation
- d) Feature Engineering

Answer: c) Data Imputation

In machine learning, what term is used for the dataset used to test the model's performance after training on the training dataset?

- a) Validation set
- b) Test set
- c) Training set
- d) Unlabeled set

Answer: b) Test set

Which algorithm is commonly used for clustering similar data points together?

- a) K-nearest Neighbors (KNN)
- b) Decision Trees
- c) K-means
- d) Support Vector Machines (SVM)

Answer: c) K-means

Which data type is used for categorical variables that have no intrinsic ordering?

- a) Integer
- b) String
- c) Float
- d) Boolean

Answer: b) String

What is the purpose of feature scaling in data preprocessing?

- a) To remove outliers from the data
- b) To convert categorical features into numerical format
- c) To normalize the data to a similar scale
- d) To handle missing values in the dataset

Answer: c) To normalize the data to a similar scale

Which machine learning algorithm is used for anomaly detection and novelty detection?

- a) Naive Bayes
- b) Random Forest
- c) Support Vector Machines (SVM)
- d) Isolation Forest

Answer: d) Isolation Forest

Which statistical measure gives an idea of how much the values in a dataset vary from the mean?

- a) Mean Absolute Deviation (MAD)
- b) Standard Deviation
- c) Variance
- d) Median Absolute Deviation (MAD)

Answer: b) Standard Deviation

Which Python library is commonly used for data manipulation and analysis?

- a) TensorFlow
- b) Keras
- c) Pandas
- d) Scikit-learn

Answer: c) Pandas

Which Python library is widely used for data visualization?

- a) NumPy
- b) Matplotlib
- c) Seaborn
- d) SciPy

Answer: b) Matplotlib

Which statement best describes the term "overfitting" in machine learning?

- a) The model performs well on the training data but poorly on unseen data.
- b) The model performs equally well on both the training and test data.
- c) The model cannot capture complex patterns in the data.
- d) The model is under-trained and lacks accuracy.

Answer: a) The model performs well on the training data but poorly on unseen data.

Which technique is used for reducing the dimensionality of data while preserving its variance?

- a) Principal Component Analysis (PCA)
- b) K-means Clustering
- c) Decision Trees
- d) Ridge Regression

Answer: a) Principal Component Analysis (PCA)

Which machine learning algorithm is best suited for image recognition tasks?

- a) K-nearest Neighbors (KNN)
- b) Support Vector Machines (SVM)
- c) Convolutional Neural Networks (CNN)
- d) Decision Trees

Answer: c) Convolutional Neural Networks (CNN)

Which method can be used to handle imbalanced datasets in classification problems?

- a) Randomly removing samples from the majority class
- b) Using accuracy as the evaluation metric
- c) Oversampling the minority class
- d) Ignoring the class imbalance and training the model as usual

Answer: c) Oversampling the minority class

Which algorithm is commonly used for text classification tasks, such as spam detection?

- a) Linear Regression
- b) K-nearest Neighbors (KNN)
- c) Naive Bayes
- d) Gradient Boosting Machines (GBM)

Answer: c) Naive Bayes

What is the main advantage of using cloud services for AI and ML applications?

- a) Lower computational power
- b) Reduced cost and scalability
- c) Restricted access to AI models
- d) Limited storage capabilities

Answer: b) Reduced cost and scalability

Which Python library is commonly used for creating interactive data visualizations?

- a) Matplotlib
- b) Seaborn
- c) Plotly
- d) Bokeh

Answer: c) Plotly

What type of data preprocessing technique is used to convert text data into numerical format for machine learning algorithms?

- a) Data Imputation
- b) Feature Scaling
- c) Feature Engineering
- d) Text Encoding

Answer: d) Text Encoding

Which supervised learning algorithm is used for making predictions with discrete or categorical target variables?

- a) Linear Regression
- b) Decision Trees
- c) Logistic Regression
- d) K-means Clustering

Answer: c) Logistic Regression

What is the purpose of cross-validation in machine learning?

- a) To evaluate the model's performance on unseen data
- b) To compare different machine learning algorithms
- c) To handle missing values in the dataset
- d) To increase the model's complexity

Answer: a) To evaluate the model's performance on unseen data

Which evaluation metric is used to assess the performance of a regression model by measuring the average difference between predicted and actual values?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Absolute Error (MAE)
- d) Precision

Answer: c) Mean Absolute Error (MAE)

Which algorithm is used for imputing missing values in a dataset based on the relationship between variables?

- a) Linear Regression
- b) K-means Clustering
- c) Multiple Imputation by Chained Equations (MICE)
- d) Decision Trees

Answer: c) Multiple Imputation by Chained Equations (MICE)

Which Python library provides tools for data manipulation and analysis, as well as mathematical functions and arrays?

- a) Scikit-learn
- b) NumPy
- c) Pandas
- d) Matplotlib

Answer: b) NumPy

Which machine learning algorithm can be used for both classification and regression tasks and is based on an ensemble of decision trees?

- a) K-nearest Neighbors (KNN)
- b) Random Forest
- c) Naive Bayes
- d) Support Vector Machines (SVM)

Answer: b) Random Forest

In a confusion matrix, which term represents the number of correctly predicted positive instances?

- a) True Positive (TP)
- b) False Positive (FP)
- c) True Negative (TN)
- d) False Negative (FN)

Answer: a) True Positive (TP)

Which unsupervised learning algorithm is used to find the optimal number of clusters in a dataset?

- a) K-means Clustering
- b) Hierarchical Clustering
- c) Density-Based Spatial Clustering of Applications with Noise (DBSCAN)
- d) Elbow Method

Answer: d) Elbow Method

Which Python library provides efficient implementations of many machine learning algorithms and is widely used for building predictive models?

- a) TensorFlow
- b) Keras
- c) Scikit-learn
- d) PyTorch

Answer: c) Scikit-learn

Which data preprocessing technique is used to transform categorical variables into numerical form by assigning each category a unique integer?

- a) Data Normalization
- b) Feature Scaling
- c) One-Hot Encoding
- d) Label Encoding

Answer: d) Label Encoding

Which machine learning algorithm is used for reducing the dimensionality of data while preserving its non-linear structure?

- a) Principal Component Analysis (PCA)
 - b) t-distributed Stochastic Neighbor Embedding (t-SNE)
 - c) Linear Discriminant Analysis (LDA)
 - d) K-means Clustering
- Answer: b) t-distributed Stochastic Neighbor Embedding (t-SNE)

Which technique is used to handle missing values by estimating the missing values based on the values of other variables in the dataset?

- a) Data Imputation
 - b) Data Normalization
 - c) Data Encoding
 - d) Data Scaling
- Answer: a) Data Imputation

Which machine learning algorithm is used for making predictions with continuous target variables?

- a) Logistic Regression
 - b) Decision Trees
 - c) K-means Clustering
 - d) Linear Regression
- Answer: d) Linear Regression

Which evaluation metric is commonly used for imbalanced datasets and represents the ability of a model to identify positive instances correctly?

- a) F1 Score
 - b) Accuracy
 - c) Precision
 - d) Area Under the Curve (AUC)
- Answer: a) F1 Score

Which technique is used for feature selection to identify the most relevant features in a dataset?

- a) Lasso Regression
 - b) Ridge Regression
 - c) Recursive Feature Elimination (RFE)
 - d) Principal Component Analysis (PCA)
- Answer: c) Recursive Feature Elimination (RFE)

Which Python library is used for creating deep learning models and neural networks?

- a) TensorFlow
 - b) Keras
 - c) Pandas
 - d) NumPy
- Answer: b) Keras

Which technique is used to handle the class imbalance problem by generating synthetic samples for the minority class?

- a) Ensemble Learning
 - b) SMOTE (Synthetic Minority Over-sampling Technique)
 - c) Ridge Regression
 - d) Recursive Feature Elimination (RFE)
- Answer: b) SMOTE (Synthetic Minority Over-sampling Technique)

Which machine learning algorithm is used for finding patterns and relationships in data using reinforcement signals?

- a) K-means Clustering
- b) Decision Trees
- c) Reinforcement Learning
- d) Linear Regression

Answer: c) Reinforcement Learning

Which method is used to split the dataset into training and testing sets while preserving the original class distribution?

- a) K-fold Cross-Validation
- b) Hold-out Validation
- c) Stratified Sampling
- d) Random Sampling

Answer: c) Stratified Sampling

Which Python library is used for creating and training deep learning models with efficient numerical computations?

- a) NumPy
- b) Matplotlib
- c) TensorFlow
- d) Seaborn

Answer: c) TensorFlow

Which evaluation metric is used to assess the performance of a classification model by measuring the trade-off between precision and recall?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Absolute Error (MAE)
- d) Area Under the Curve (AUC)

Answer: a) F1 Score

Which unsupervised learning algorithm is used for finding patterns in data based on the concept of "association" between items?

- a) Apriori Algorithm
- b) Hierarchical Clustering
- c) k-Nearest Neighbors (k-NN)
- d) Principal Component Analysis (PCA)

Answer: a) Apriori Algorithm

Which technique is used for tuning hyperparameters by searching through different combinations to find the best model performance?

- a) Grid Search
- b) Random Search
- c) Gradient Descent
- d) Stochastic Optimization

Answer: a) Grid Search

Which evaluation metric is used to assess the performance of a classification model by measuring the area under the Receiver Operating Characteristic (ROC) curve?

- a) F1 Score
- b) R-squared (R²) Score

- c) Mean Absolute Error (MAE)
- d) Area Under the Curve (AUC)

Answer: d) Area Under the Curve (AUC)

Which Python library is used for statistical computations and hypothesis testing?

- a) Pandas
- b) NumPy
- c) SciPy
- d) Matplotlib

Answer: c) SciPy

Which machine learning algorithm is used for predicting categorical target variables with more than two classes?

- a) Decision Trees
- b) K-means Clustering
- c) Naive Bayes
- d) Random Forest

Answer: d) Random Forest

Which data preprocessing technique is used to scale the features to a specific range, such as [0, 1] or [-1, 1]?

- a) Feature Scaling
- b) Data Normalization
- c) Data Imputation
- d) Label Encoding

Answer: a) Feature Scaling

Which Python library is used for creating interactive visualizations for exploratory data analysis?

- a) Seaborn
- b) Plotly
- c) Matplotlib
- d) Pandas

Answer: b) Plotly

Which technique is used to handle the curse of dimensionality by projecting the data into a lower-dimensional space?

- a) Ridge Regression
- b) Principal Component Analysis (PCA)
- c) Recursive Feature Elimination (RFE)
- d) Gradient Boosting Machines (GBM)

Answer: b) Principal Component Analysis (PCA)

Which evaluation metric is used to assess the performance of a regression model by measuring the proportion of variance in the target variable explained by the model?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Absolute Error (MAE)
- d) Precision

Answer: b) R-squared (R²) Score

Which unsupervised learning algorithm is used for grouping similar data points into clusters based on their distance from cluster centers?

- a) K-means Clustering
- b) Decision Trees

- c) Principal Component Analysis (PCA)
- d) k-Nearest Neighbors (k-NN)

Answer: a) K-means Clustering

Which data preprocessing technique is used to scale the features to have a mean of zero and a standard deviation of one?

- a) Feature Scaling
- b) Data Normalization
- c) Data Imputation
- d) Label Encoding

Answer: a) Feature Scaling

Which machine learning algorithm is used for predicting categorical target variables with two classes?

- a) Decision Trees
- b) K-means Clustering
- c) Naive Bayes
- d) Support Vector Machines (SVM)

Answer: d) Support Vector Machines (SVM)

Which technique is used for selecting the best features in a dataset based on their importance in predicting the target variable?

- a) Lasso Regression
- b) Ridge Regression
- c) Recursive Feature Elimination (RFE)
- d) Principal Component Analysis (PCA)

Answer: c) Recursive Feature Elimination (RFE)

Which Python library is used for creating and training deep learning models with a focus on simplicity and ease of use?

- a) NumPy
- b) Matplotlib
- c) Keras
- d) TensorFlow

Answer: c) Keras

Which evaluation metric is used to assess the performance of a classification model by measuring the proportion of true positive predictions out of all positive instances?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Absolute Error (MAE)
- d) Recall

Answer: d) Recall

Which machine learning algorithm is used for finding patterns in data by dividing the dataset into subsets using a series of binary decisions?

- a) K-nearest Neighbors (KNN)
- b) Decision Trees
- c) Random Forest
- d) Support Vector Machines (SVM)

Answer: b) Decision Trees

Which method is used to split the dataset into training, validation, and testing sets to assess model performance effectively?

- a) K-fold Cross-Validation
- b) Hold-out Validation
- c) Stratified Sampling
- d) Random Sampling

Answer: a) K-fold Cross-Validation

Which Python library is used for creating and training deep learning models with a focus on GPU acceleration?

- a) NumPy
- b) Matplotlib
- c) PyTorch
- d) Keras

Answer: c) PyTorch

Which evaluation metric is used to assess the performance of a regression model by measuring the average squared difference between predicted and actual values?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Squared Error (MSE)
- d) Precision

Answer: c) Mean Squared Error (MSE)

Which unsupervised learning algorithm is used for grouping data points based on their similarity to a given number of cluster centroids?

- a) K-means Clustering
- b) Hierarchical Clustering
- c) t-distributed Stochastic Neighbor Embedding (t-SNE)
- d) Principal Component Analysis (PCA)

Answer: a) K-means Clustering

Which data preprocessing technique is used to convert categorical variables into numerical form while creating binary columns for each category?

- a) Data Normalization
- b) Feature Scaling
- c) One-Hot Encoding
- d) Label Encoding

Answer: c) One-Hot Encoding

Which machine learning algorithm is used for predicting categorical target variables with more than two classes, often in the context of decision-making?

- a) Decision Trees
- b) K-means Clustering
- c) Naive Bayes
- d) Gradient Boosting Machines (GBM)

Answer: a) Decision Trees

Which technique is used to handle the class imbalance problem by combining the predictions of multiple models?

- a) Ensemble Learning
- b) SMOTE (Synthetic Minority Over-sampling Technique)
- c) Ridge Regression
- d) Recursive Feature Elimination (RFE)

Answer: a) Ensemble Learning

Which Python library is used for creating and training deep learning models with a focus on efficiency and speed?

- a) NumPy
- b) Matplotlib
- c) TensorFlow
- d) Keras

Answer: c) TensorFlow

Which evaluation metric is used to assess the performance of a classification model by measuring the ability to correctly identify negative instances?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Absolute Error (MAE)
- d) Specificity

Answer: d) Specificity

Which machine learning algorithm is used for predicting continuous target variables based on an ensemble of decision trees?

- a) Decision Trees
- b) K-means Clustering
- c) Random Forest
- d) Linear Regression

Answer: c) Random Forest

Which method is used to split the dataset into training, validation, and testing sets while preserving the original class distribution and considering the imbalance in the target variable?

- a) K-fold Cross-Validation
- b) Hold-out Validation
- c) Stratified Sampling
- d) Random Sampling

Answer: c) Stratified Sampling

Which Python library is used for creating and training deep learning models with a focus on GPU acceleration and distributed computing?

- a) NumPy
- b) Matplotlib
- c) PyTorch
- d) Keras

Answer: c) PyTorch

Which evaluation metric is used to assess the performance of a regression model by measuring the proportion of variance in the target variable not explained by the model?

- a) F1 Score
- b) R-squared (R²) Score
- c) Mean Squared Error (MSE)
- d) Precision

Answer: b) R-squared (R²) Score

Which unsupervised learning algorithm is used for grouping data points into clusters based on their similarity in a hierarchical manner?

- a) K-means Clustering
- b) Hierarchical Clustering
- c) t-distributed Stochastic Neighbor Embedding (t-SNE)

d) Principal Component Analysis (PCA)

Answer: b) Hierarchical Clustering

Which data preprocessing technique is used to convert categorical variables into numerical form by replacing each category with its corresponding frequency in the dataset?

a) Data Normalization

b) Feature Scaling

c) Frequency Encoding

d) Label Encoding

Answer: c) Frequency Encoding

Which machine learning algorithm is used for predicting categorical target variables with more than two classes, often in the context of probability estimation?

a) Decision Trees

b) K-means Clustering

c) Naive Bayes

d) Logistic Regression

Answer: d) Logistic Regression

Which technique is used to handle the class imbalance problem by generating synthetic samples for the minority class and merging them with the original dataset?

a) Ensemble Learning

b) SMOTE (Synthetic Minority Over-sampling Technique)

c) Ridge Regression

d) Recursive Feature Elimination (RFE)

Answer: b) SMOTE (Synthetic Minority Over-sampling Technique)

Which Python library is used for creating and training deep learning models with a focus on flexibility and customization?

a) NumPy

b) Matplotlib

c) TensorFlow

d) PyTorch

Answer: d) PyTorch

Which evaluation metric is used to assess the performance of a classification model by measuring the ability to correctly identify positive instances?

a) F1 Score

b) R-squared (R²) Score

c) Mean Absolute Error (MAE)

d) Sensitivity

Answer: d) Sensitivity

Which machine learning algorithm is used for predicting continuous target variables based on an ensemble of decision trees with regularization?

a) Decision Trees

b) K-means Clustering

c) Random Forest

d) Ridge Regression

Answer: d) Ridge Regression

Which method is used to split the dataset into training, validation, and testing sets by randomly assigning instances to each set?

- a) K-fold Cross-Validation
- b) Hold-out Validation
- c) Stratified Sampling
- d) Random Sampling

Answer: d) Random Sampling

Which Python library is used for creating and training deep learning models with a focus on efficient computation on CPUs and GPUs?

- a) NumPy
- b) Matplotlib
- c) TensorFlow
- d) Keras

Answer: c) TensorFlow

Which evaluation metric is used to assess the performance of a regression model by measuring the average absolute difference between predicted and actual values?

- a) F1 Score
- b) R-squared (R^2) Score
- c) Mean Absolute Error (MAE)
- d) Precision

Answer: c) Mean Absolute Error (MAE)

Which unsupervised learning algorithm is used for projecting high-dimensional data into a lower-dimensional space while preserving the pairwise distances between data points?

- a) K-means Clustering
- b) Hierarchical Clustering
- c) t-distributed Stochastic Neighbor Embedding (t-SNE)
- d) Principal Component Analysis (PCA)

Answer: c) t-distributed Stochastic Neighbor Embedding (t-SNE)

Which data preprocessing technique is used to convert categorical variables into numerical form by assigning each category a unique integer while considering the order of categories?

- a) Data Normalization
- b) Feature Scaling
- c) Ordinal Encoding
- d) Label Encoding

Answer: c) Ordinal Encoding

Which machine learning algorithm is used for predicting categorical target variables with two classes, often in the context of probability estimation?

- a) Decision Trees
- b) K-means Clustering
- c) Naive Bayes
- d) Logistic Regression

Answer: d) Logistic Regression

Which technique is used to handle the class imbalance problem by generating synthetic samples for the minority class and merging them with the original dataset using weighted averages?

- a) Ensemble Learning
- b) SMOTE (Synthetic Minority Over-sampling Technique)
- c) Ridge Regression
- d) Recursive Feature Elimination (RFE)

Answer: b) SMOTE (Synthetic Minority Over-sampling Technique)

Which Python library is used for creating and training deep learning models with a focus on ease of use and seamless integration with TensorFlow?

- a) NumPy
 - b) Matplotlib
 - c) Keras
 - d) PyTorch
- Answer: c) Keras

Which evaluation metric is used to assess the performance of a classification model by measuring the trade-off between true positive rate and false positive rate?

- a) F1 Score
 - b) R-squared (R²) Score
 - c) Mean Absolute Error (MAE)
 - d) Receiver Operating Characteristic (ROC) Curve
- Answer: d) Receiver Operating Characteristic (ROC) Curve

Which machine learning algorithm is used for predicting continuous target variables based on an ensemble of decision trees with regularization and feature selection?

- a) Decision Trees
 - b) K-means Clustering
 - c) Random Forest
 - d) Lasso Regression
- Answer: d) Lasso Regression

Which method is used to split the dataset into training, validation, and testing sets while preserving the original class distribution and considering the stratification of the target variable?

- a) K-fold Cross-Validation
 - b) Hold-out Validation
 - c) Stratified Sampling
 - d) Random Sampling
- Answer: a) K-fold Cross-Validation

Which Python library is used for creating and training deep learning models with a focus on flexibility and ease of use for researchers and practitioners?

- a) NumPy
 - b) Matplotlib
 - c) TensorFlow
 - d) PyTorch
- Answer: d) PyTorch

Which evaluation metric is used to assess the performance of a regression model by measuring the proportion of variance in the target variable explained by the model, adjusted for the number of features?

- a) F1 Score
 - b) Adjusted R-squared Score
 - c) Mean Squared Error (MSE)
 - d) Precision
- Answer: b) Adjusted R-squared Score

Which unsupervised learning algorithm is used for projecting high-dimensional data into a lower-dimensional space by preserving the pairwise distances between data points and emphasizing on the global structure of the data?

- a) K-means Clustering

- b) Hierarchical Clustering
- c) t-distributed Stochastic Neighbor Embedding (t-SNE)
- d) Principal Component Analysis (PCA)

Answer: c) t-distributed Stochastic Neighbor Embedding (t-SNE)

Which data preprocessing technique is used to convert categorical variables into numerical form by replacing each category with its corresponding mean or median value?

- a) Data Normalization
- b) Feature Scaling
- c) Mean Encoding
- d) Label Encoding

Answer: c) Mean Encoding

Which machine learning algorithm is used for predicting categorical target variables with two classes, often in the context of probabilistic classification?

- a) Decision Trees
- b) K-means Clustering
- c) Naive Bayes
- d) Logistic Regression

Answer: d) Logistic Regression

Which technique is used to handle the class imbalance problem by generating synthetic samples for the minority class and merging them with the original dataset using a weighted average based on the distance of the nearest neighbors?

- a) Ensemble Learning
- b) SMOTE (Synthetic Minority Over-sampling Technique)
- c) Ridge Regression
- d) K-nearest Neighbors (KNN)

Answer: b) SMOTE (Synthetic Minority Over-sampling Technique)