

1. **What is bias?**
In machine learning, bias refers to the error introduced by approximating a real-world problem, which may be complex, by a simplified model.
2. **What is variance?**
Variance is the amount by which a model's prediction can change for different training datasets. It measures the model's sensitivity to the training data.
3. **Discuss bias-variance tradeoff in machine learning:**
The bias-variance tradeoff is the balance between a model's ability to capture the underlying patterns in the data (bias) and its sensitivity to variations in the training set (variance). A model with high bias may oversimplify the data, while a high variance model may overfit to noise. The goal is to find the right level of complexity to generalize well to unseen data.
4. **What is supervised, unsupervised, and reinforcement learning?**
Supervised learning involves training a model on a labeled dataset. Unsupervised learning involves finding patterns in unlabeled data. Reinforcement learning involves training a model to make sequences of decisions by interacting with an environment.
5. **What is classification? How is it different from regression?**
Classification is a task where the goal is to predict the categorical class labels of new instances. Regression predicts continuous numerical values. The main difference lies in the type of output – discrete classes for classification and continuous values for regression.
6. **Define softmax function and its use:**
The softmax function converts a vector of real numbers into a probability distribution. It is often used as the output activation function in a neural network for multi-class classification problems.
7. **What is an activation function? Discuss various activation functions:**
An activation function introduces non-linearity to the neural network. Examples include sigmoid, hyperbolic tangent (tanh), and Rectified Linear Unit (ReLU).
8. **How is sigmoid different than ReLU?**
Sigmoid outputs values between 0 and 1, suitable for binary classification. ReLU outputs the input for positive values and zero for negative values, making it computationally efficient.
9. **What is a perceptron?**
A perceptron is the simplest form of a neural network, a binary classifier that takes multiple binary inputs and produces a single binary output.
10. **What is a multilayer perceptron (MLP)? How is it better than a perceptron?**
An MLP is a neural network with multiple layers, including input, hidden, and output layers. It can learn complex patterns and relationships, unlike a perceptron.
11. **What is data preprocessing?**
Data preprocessing is the process of cleaning and transforming raw data into a suitable format for machine learning. It includes handling missing values, scaling, encoding categorical variables, etc.
12. **Discuss various data preprocessing steps:**
Steps include handling missing data, dealing with outliers, scaling features, encoding categorical variables, and splitting data into training and testing sets.

13. **What is standardization, normalization, and binarization?**
Standardization scales features to have a mean of 0 and a standard deviation of 1. Normalization scales features to a range of 0 to 1. Binarization converts numerical values into binary values based on a threshold.
14. **What is feature scaling and feature extraction?**
Feature scaling ensures that features are on a similar scale. Feature extraction involves selecting a subset of relevant features from the original set.
15. **Define min-max scaling:**
Min-max scaling transforms features to a specific range (e.g., 0 to 1) by subtracting the minimum value and dividing by the range.
16. **What is a label encoder?**
A label encoder converts categorical labels into numerical values, often used for encoding target variables in classification tasks.
17. **What is one-hot encoding? What is the significance of using it?**
One-hot encoding converts categorical variables into binary vectors. It is significant as it allows algorithms to work with categorical data without assuming any ordinal relationship.
18. **What is a confusion matrix and classification report?**
A confusion matrix summarizes the performance of a classification algorithm, while a classification report provides metrics such as precision, recall, and F1-score for each class.
19. **What is precision, recall, and F1-score?**
Precision is the ratio of correctly predicted positive observations to the total predicted positives. Recall is the ratio of correctly predicted positive observations to the all observations in actual class. F1-score is the weighted average of precision and recall.
20. **What is MLPClassifier and MLPRegressor? Define its methods from sklearn:**
MLPClassifier is a multi-layer perceptron classifier, and MLPRegressor is a multi-layer perceptron regressor in scikit-learn. Methods include fit, predict, and score.
21. **What is the pandas library? Discuss important functionalities:**
Pandas is a Python library for data manipulation and analysis. Key functionalities include data structures like DataFrames, handling missing data, and providing powerful methods for data manipulation and analysis.
22. **What is the use of NumPy and scikit-learn?**
NumPy is used for numerical operations in Python, providing support for large, multi-dimensional arrays and matrices. Scikit-learn is a machine learning library that includes tools for classification, regression, clustering, and more.
23. **What is Bayes' theorem?**
Bayes' theorem is a mathematical formula that describes the probability of an event, based on prior knowledge of conditions that might be related to the event.
24. **What is a Bayesian classifier?**
A Bayesian classifier is a probabilistic model that makes predictions using Bayes' theorem. It calculates the probability of each class and selects the class with the highest probability.

25. **What is a decision tree?**
A decision tree is a tree-like model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility.
26. **Define ID3, CART, C4.5 decision trees:**
ID3, CART, and C4.5 are algorithms for building decision trees. ID3 uses information gain, CART uses Gini impurity, and C4.5 uses gain ratio for attribute selection.
27. **What is information gain and entropy?**
Information gain measures the reduction in entropy or disorder in a system. Entropy measures the impurity or disorder of a set.
28. **What is splitinfo and gain ratio in C4.5?**
Splitinfo is the measure of the amount of information required to specify the chosen attribute. Gain ratio is the ratio of information gain to split information.
29. **What is Gini Index?**
Gini Index is a measure of impurity used in decision tree algorithms. It quantifies how often a randomly chosen element would be incorrectly classified.
30. **What is dimensionality reduction?**
Dimensionality reduction is the process of reducing the number of features in a dataset while preserving its essential information.
31. **What is eigenvalue and eigenvector?**
In linear algebra, an eigenvalue is a scalar that represents how an eigenvector is scaled during a linear transformation.
32. **Discuss PCA algorithm:**
Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms data into a new coordinate system where variables are uncorrelated (principal components).
33. **Discuss LDA algorithm:**
Linear Discriminant Analysis (LDA) is a dimensionality reduction technique used in the context of pattern classification to find a linear combination of features that characterizes or separates classes.
34. **What is SVM?**
Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification and regression tasks. It finds the hyperplane that best separates classes in a high-dimensional space.
35. **What is soft margin and margin?**
The soft margin in SVM allows for some misclassification to achieve a better overall fit. The margin is the distance between the hyperplane and the nearest data point from either class.
36. **What is clustering?**
Clustering is the task of grouping a set of objects in such a way that objects in the same group (cluster) are more similar to each other than to those in other groups.

37. **What is the k-means algorithm? What are its drawbacks?**

K-means is a clustering algorithm that partitions data into k clusters. Drawbacks include sensitivity to initial cluster centers and difficulty with non-convex shapes.

38. **What is the elbow method? How to use it:**

The elbow method is a technique to find the optimal number of clusters in k-means clustering by plotting the cost (sum of squared distances) as a function of the number of clusters and selecting the "elbow" point where the rate of decrease sharply changes.

39. **What is KNN? What are its drawbacks?**

K-nearest neighbors (KNN) is a simple, instance-based learning algorithm. Drawbacks include sensitivity to irrelevant features and high computational cost for large datasets.

40. **Define CNN:**

Convolutional Neural Network (CNN) is a type of neural network designed for image recognition and processing. It uses convolutional layers to automatically and adaptively learn spatial hierarchies of features.

41. **What is deep learning?**

Deep learning is a subfield of machine learning that involves neural networks with many layers (deep neural networks). It excels at learning hierarchical representations of data.

42. **What is regularization?**

Regularization is a technique used to prevent overfitting in machine learning models by adding a penalty term to the cost function.

43. **What are training set, validation set, and test set?**

The training set is used to train the model, the validation set is used to fine-tune hyperparameters, and the test set is used to evaluate the model's performance on unseen data.

44. **What is overfitting and underfitting?**

Overfitting occurs when a model is too complex and fits the training data too closely, performing poorly on new, unseen data. Underfitting occurs when a model is too simple and fails to capture the underlying patterns.

45. **What is a scatter plot and box plot?**

A scatter plot displays individual data points in two dimensions. A box plot (box-and-whisker plot) displays the distribution of a dataset and identifies outliers.

46. **What is an outlier? How to identify an outlier?**

An outlier is an observation that lies an abnormal distance from other values in a random sample. Identification methods include statistical measures like Z-score or visualization techniques like box plots.

47. **Discuss techniques to rectify outliers:**

Techniques include removing outliers, transforming variables, or imputing values based on the distribution.

48. **What is cross-validation?**

Cross-validation is a resampling technique used to evaluate machine learning models by partitioning the data into subsets, training the model on some, and testing it on the others.

49. **What is k-fold cross-validation?**

K-fold cross-validation divides the data into k subsets (folds), trains the model on k-1 folds, and tests it on the remaining fold. This process is repeated k times, with each fold serving as the test set exactly once.

50. **How to handle missing, null, or NaN values in the dataset?**

Strategies include removing rows or columns with missing values, imputing missing values using the mean or median, or using more advanced imputation techniques.

51. **How to check the number of null values in a dataset?**

In Python, you can use the `isnull()` method along with `sum()` to count the number of null values in each column of a DataFrame.

Code:

```
import pandas as pd
# Assuming df is your DataFrame
null_counts = df.isnull().sum()
print(null_counts)
```