

Answer 1: Collinearity can negatively impact statistical models by introducing instability, making it difficult to interpret the individual effects of variables and potentially leading to overfitting. By reducing the dimensionality of the data, collinear variables can be combined or eliminated, reducing the multicollinearity issue and improving the model's performance and interpretability.

The correct answer is d) Collinearity.

Answer2: The correct answer is b) Random Forest.

Random Forest is a machine learning algorithm that is based on the idea of bagging (bootstrap aggregating). Bagging is a technique that involves creating multiple subsets of the original dataset through resampling, training individual models on these subsets, and then combining their predictions to make the final prediction.

Answer3: The correct answer is c) Decision trees are prone to overfitting.

Decision trees have a disadvantage of being prone to overfitting. Overfitting occurs when a decision tree model becomes too complex and captures noise or irrelevant patterns in the training data, leading to poor generalization performance on unseen data.

Answer4: The term known as on which machine learning algorithms build a model based on sample data is c) Training data.

Training data refers to the labeled or annotated data that is used to train a machine learning model. It consists of input samples or instances along with their corresponding output labels or target values. The purpose of using training data is to enable the machine learning algorithm to learn patterns, relationships, and rules from the data in order to build a predictive or descriptive model.

Answer5: The correct answer is c) Anomaly detection.

Anomaly detection is a machine learning technique specifically designed to detect outliers or anomalies in data. It aims to identify instances that deviate significantly from the norm or expected patterns within a dataset.

Answer6: "Case based" is not a recognized numerical function representation in the context of machine learning. It may refer to case-based reasoning, which is a technique where new problems are solved by reusing solutions from similar past cases. However, it is not specifically associated with a numerical function representation.

The incorrect numerical function in the various function representation of machine learning is c) Case based.

Answer7: The correct answer is d) Both a and b.

Analysis of machine learning algorithms requires both statistical learning theory and computational learning theory.

The correct answer is b) SVG.

SVG is not a recognized machine learning algorithm. It may be a typographical error or a term that is not commonly associated with machine learning.

Answer8: The correct answer is c) Both a and b.

The k-nearest neighbor (k-NN) algorithm has certain difficulties, which are:

a) Curse of dimensionality: The curse of dimensionality refers to the problem where the performance of k-NN degrades as the number of dimensions or features in the dataset increases. As the dimensionality increases, the available data becomes sparse, and the distances between data points become less informative. This can lead to decreased accuracy and increased computational complexity.

b) Calculation of distance for all training cases: The k-NN algorithm calculates the distance between a test case and all the training cases in order to determine the k nearest neighbors. This can be computationally expensive, especially for large datasets with many training examples. As the dataset grows, the time required for distance calculations increases, making the algorithm slower.

Therefore, both the curse of dimensionality and the computational burden of calculating distances for all training cases are difficulties associated with the k-nearest neighbor algorithm.

Answer9: The correct answer is b) 2.

Answer10: The correct answer is d) KMeans.

KMeans is not a supervised learning algorithm. It is an unsupervised learning algorithm used for clustering, which is a task of grouping similar data points together based on their similarity or distance in the feature space.

Answer11: The correct answer is c) Neither feature nor number of groups is known.

Unsupervised learning is a type of machine learning where the algorithm learns patterns, relationships, or structures in data without any explicit supervision or labeled examples. In unsupervised learning, the data does not have predefined labels or target variables.

Answer12: The correct answer is b) SVG.

SVG is not a recognized machine learning algorithm. It may be a typographical error or a term that is not commonly associated with machine learning.

Answer13: The correct answer is b) Underfitting.

Underfitting occurs when a machine learning model is too simple or lacks complexity to capture the underlying trend or patterns in the input data. In such cases, the model fails to learn the training data well and performs poorly on both the training set and unseen data.

Answer14: The correct answer is a) Reinforcement learning.

Reinforcement learning is a type of machine learning that focuses on learning through interactions with an environment. It is particularly well-suited for scenarios where an agent learns to make sequential decisions in order to maximize a cumulative reward signal. The applications mentioned, such as real-time decisions, game AI, learning tasks, skill acquisition, and robot navigation, often involve dynamic decision-making in complex environments, making reinforcement learning a suitable approach.

Answer15: The correct answer is b) Mean squared error (MSE).

Mean squared error is a common evaluation metric used in machine learning to measure the average squared difference between the predicted output of a classifier or regression model and the actual output or ground truth values. It is calculated by taking the average of the squared differences between the predicted and actual values.

Answer 16: The correct answer is a) Linear, binary.

Logistic regression is a linear regression technique that is commonly used to model data with a binary outcome or response variable. The binary outcome refers to a categorical variable that can take one of two possible values, such as "yes" or "no," "true" or "false," or 0 or 1.

Answer 17: The correct answer is A. supervised learning.

Classifying reviews of a new Netflix series involves training a machine learning model using labeled data where the reviews are marked as positive, negative, or neutral. The labeled data serves as the training set, and the goal is to build a model that can accurately predict the sentiment or classification of new, unseen reviews.

Answer18: The correct answer is C. both a and b.

Both Euclidean distance and Manhattan distance are powerful distance metrics used by geometric models.

Answer19: The correct answer is D. none of these.

The given options, A, B, and C, are not techniques for dimensionality reduction. Instead, they describe different strategies for feature selection or data preprocessing. While these strategies may be useful for data cleaning and preparing the data before applying dimensionality reduction techniques, they do not directly reduce the dimensions of the dataset.

Answer20: The correct answer is C. input attribute.

Both supervised learning and unsupervised clustering require input attributes.

Answer21: The correct answer is (A) SVM allows very low error in classification.

In Support Vector Machines (SVM), the term "hard margin" refers to the condition where the SVM algorithm aims to find a decision boundary with the maximum possible margin between classes while allowing very low error in classification. In other words, the SVM tries to find a hyperplane that separates the classes with the least possible misclassification.

Answer22: The correct answer is (B) Only 2: Depth of Tree.

Increasing the depth of the trees in a Random Forest model can result in overfitting. Overfitting occurs when the model becomes too complex and captures the noise or irrelevant patterns in the training data, leading to poor generalization on unseen data.

Answer23: The correct answer is (A)  $-(6/10 \log(6/10) + 4/10 \log(4/10))$ .

Entropy is a measure of the impurity or disorder in a set of values. In the case of a binary target variable with two possible outcomes (0 and 1), the entropy can be calculated using the formula:

$$\text{Entropy} = - (p_0 \log(p_0) + p_1 \log(p_1))$$

where  $p_0$  is the proportion of values with outcome 0 and  $p_1$  is the proportion of values with outcome 1.

In this case, we have 6 occurrences of 0 and 4 occurrences of 1 out of a total of 10 values.

So,  $p_0 = 6/10 = 0.6$  and  $p_1 = 4/10 = 0.4$ .

Plugging these values into the entropy formula:

$$\text{Entropy} = - (0.6 \log(0.6) + 0.4 \log(0.4))$$

Simplifying this expression gives us:

$$\text{Entropy} = -(6/10 \log(6/10) + 4/10 \log(4/10))$$

Therefore, the correct answer is (A)  $-(6/10 \log(6/10) + 4/10 \log(4/10))$ .

Answer24: The correct answer is (A) weights are regularized with the l1 norm.

Lasso, which stands for "Least Absolute Shrinkage and Selection Operator," is a regularization technique used in linear regression. It can be interpreted as a least-squares linear regression model with an additional regularization term.

Answer25: The correct answer is (D) Perceptron.

The Perceptron algorithm is a linear classification algorithm that aims to find a decision boundary that separates the classes. It updates the weights based on misclassified points until all points are correctly classified or a maximum number of iterations is reached. If the training set is linearly separable, the Perceptron algorithm is guaranteed to converge and find a separating hyperplane.

Answer26: The correct answer is (D) Either 2 or 3.

When dealing with multi-collinear features (features that are highly correlated with each other), it is necessary to take appropriate actions to address the collinearity and avoid potential issues in the model.

Answer27: The correct answer is (B) increase by 5 pounds.

In the given least squares line equation,  $y = 120 + 5x$ , the coefficient of  $x$  is 5. This coefficient represents the rate of change of the weight ( $y$ ) with respect to the height ( $x$ ).

According to the equation, for every one-unit increase in height ( $x$ ), the weight ( $y$ ) is expected to increase by the coefficient of  $x$ , which is 5.

Therefore, if the height is increased by one inch, the weight should increase by 5 pounds.

Hence, the correct answer is (B) increase by 5 pounds.

Answer28: The correct answer is (D) Minimize the squared distance from the points.

The linear regression equation, specifically the Ordinary Least Squares (OLS) method, aims to find the best-fitting line that minimizes the sum of the squared differences between the predicted values and the actual values of the dependent variable.

In other words, the line described by the linear regression equation seeks to minimize the squared distance (also known as the residual or error) between the observed data points and the predicted values on the line. This approach ensures that the line is as close as possible to the data points, providing the best linear approximation to the relationship between the independent variable(s) and the dependent variable.

Hence, the correct answer is (D) Minimize the squared distance from the points.

Answer29: The correct answer is (B) As the value of one attribute increases, the value of the second attribute also increases.

The correlation coefficient is a measure of the strength and direction of the linear relationship between two variables. It ranges from -1 to +1, where a value of +1 indicates a perfect positive linear relationship, 0 indicates no linear relationship, and -1 indicates a perfect negative linear relationship.

Answer30: The correct answer is (B) Convolutional Neural Network.

When dealing with image identification or recognition problems, such as recognizing a dog in a photo, Convolutional Neural Networks (CNNs) are the most suited neural network architecture.