

What does one mean by the term "machine learning"?

Machine learning is a field of computer science that aims to teach computers how to learn and act without being explicitly programmed. Machine learning involves the construction of algorithms that adapt their models to improve their ability to make predictions.

Can you name 4 types of problems where it shines?

Sentimental Analysis, Demand analysis, spam classification, Recommendation engine.

What is a labeled training set?

In machine learning, a labeled training set is a dataset that has been annotated with corresponding labels or output values.

In supervised learning, the goal is to learn a function that maps input data to output values, and a labeled training set provides examples of the input-output pairs that the machine learning algorithm should learn from.

For example, in a computer vision task, a labeled training set might consist of images of animals with corresponding labels indicating the type of animal in the image (e.g., "cat", "dog", "horse", etc.).

What are the two most common supervised tasks?

Classification: Classification is a supervised learning task where the goal is to predict a categorical or discrete output variable based on input variables.

Regression: Regression is a supervised learning task where the goal is to predict a continuous output variable based on input variables.

Can you name 4 common unsupervised tasks?

Clustering, Dimensionality Reduction, Anomaly Detection, Association Rule Learning

What type of Machine Learning algorithm would you use to allow a robot to walk in various unknown terrains?

One type of Machine Learning algorithm that could be used to allow a robot to walk in various unknown terrains is Reinforcement Learning, which enables the robot to learn through trial and error by receiving rewards or punishments based on its actions.

What type of algorithm would you use to segment your customers into multiple groups?

One type of algorithm that can be used to segment customers into multiple groups is Clustering, specifically K-means clustering, which is an unsupervised learning technique that partitions the customers into K distinct clusters based on their similarities and differences in terms of various features or attributes.

Would you frame the problem of spam detection as a supervised learning problem or an unsupervised learning problem?

Spam detection is typically framed as a supervised learning problem, where the algorithm is trained on labeled data (emails labeled as spam or not spam) to learn a decision boundary that can classify new emails as spam or not spam.

Unsupervised learning techniques, such as clustering or anomaly detection, can also be used in some cases, but they are less commonly used for spam detection.

What is an online learning system?

An online learning system is a type of machine learning system that can learn and adapt to new data in real-time as it becomes available, without requiring retraining on the entire dataset.

In contrast to batch learning, where the algorithm is trained on a fixed dataset offline and then applied to new data, online learning algorithms continuously update their model parameters based on new observations or data streams, allowing them to adapt and improve over time.

This makes online learning well-suited for applications where the data is constantly changing, such as in recommendation systems, fraud detection, or predictive maintenance.

What is out-of-core learning?

Out-of-core learning is a type of machine learning technique that is used when the size of the dataset is too large to fit into the memory of a computer, requiring the algorithm to process the data in smaller, manageable chunks.

Instead of loading the entire dataset into memory, out-of-core learning algorithms read and process small batches of data at a time, often using techniques such as stochastic gradient descent (SGD) to optimize the model parameters on each batch.

What type of learning algorithm relies on a similarity measure to make predictions?

The type of learning algorithm that relies on a similarity measure to make predictions is called Instance-based learning, also known as memory-based learning. In this type of learning, the algorithm stores the entire training dataset in memory and makes predictions based on the similarity of new input instances to the examples in the training data.

What is the difference between a model parameter and a learning algorithm's hyperparameter?

In machine learning, a model parameter is a variable that is learned by the algorithm during the training process, while a learning algorithm's hyperparameter is a setting that is chosen by the developer or user before training the model.

A model parameter is a feature or coefficient that determines the behavior of the model and is learned by the algorithm to minimize the loss function on the training data. For example, in linear regression, the model parameters are the coefficients of the linear equation that best fits the data.

On the other hand, a learning algorithm's hyperparameters are settings that are chosen before training the model and determine the behavior of the algorithm itself. Examples of hyperparameters include the learning rate, regularization strength, number of layers in a neural network, or number of trees in a random forest. These hyperparameters affect how the algorithm learns from the data and how well it generalizes to new, unseen data.

What do model based learning algorithms search for? What is the most common strategy they use to succeed? How do they make predictions?

Model-based learning algorithms search for a model that can generalize well to new, unseen data. The most common strategy they use to succeed is to learn a function that minimizes the loss function on the training data while avoiding overfitting.

Model-based learning algorithms make predictions by using the learned function to map input features to output values. For example, in linear regression, the learned function is a linear equation that maps input features to a continuous output value. In classification problems, the learned function is a decision boundary that separates different classes of data.

To make predictions on new, unseen data, the model-based learning algorithm applies the learned function to the input features and produces an output value. The quality of the prediction depends on how well the learned function generalizes to new data and how well the training data represents the underlying distribution of the data.

Can you name 4 of the main challenges in Machine Learning?

Data quality and quantity: Obtaining sufficient amounts of high-quality data that are representative of the problem at hand can be a major challenge.

Overfitting: Learning models that are overly complex can lead to overfitting, where the model performs well on the training data but poorly on new, unseen data.

Selection of appropriate algorithms and hyperparameters: Selecting the most appropriate algorithm and tuning its hyperparameters can have a significant impact on the model's performance.

Interpretability: Understanding how a model makes its predictions and why it behaves in a certain way is important for building trust and ensuring that the model is making decisions in a fair and ethical manner.

If your model performs great on the training data but generalizes poorly to new instances, what is happening? Can you name 3 possible solutions?

If a model performs well on the training data but generalizes poorly to new instances, it is likely overfitting, which means that the model has learned the noise and specific features of the training data instead of generalizing to new data.

Three possible solutions to this problem are:

Collect more data: If the model is overfitting due to a lack of data, collecting more data can help the model learn more representative features of the problem and avoid overfitting.

Simplify the model: Reducing the complexity of the model by removing features, decreasing the number of layers in a neural network, or increasing the regularization strength can help to prevent overfitting.

Use cross-validation: Cross-validation can help to evaluate the model's performance on new, unseen data by splitting the available data into training and validation sets, training the model on the training set, and evaluating its performance on the validation set. This can help to identify overfitting and select the best model that generalizes well to new data.

What is a test set and why would you want to use it?

In machine learning, a test set is a set of data that is held out from the training data and is used to evaluate the performance of the trained model on new, unseen data.

The purpose of using a test set is to assess the model's ability to generalize to new, unseen data and to estimate its performance in the real world. By evaluating the model on a separate set of data, we can avoid overfitting to the training data and obtain a more accurate estimate of the model's performance.

The test set is typically used after the model has been trained and tuned using the training set and a separate validation set. The test set is kept completely separate from the training and validation data and is only used once to assess the final performance of the trained model.

It is important to use a test set to avoid overfitting and to obtain a more accurate estimate of the model's performance on new data. Without a test set, the model may appear to perform well on the training and validation data but may not generalize well to new, unseen data.

What is the purpose of a validation set?

In machine learning, a validation set is a set of data that is used to evaluate the performance of a model during training and to tune its hyperparameters.

The purpose of using a validation set is to assess the model's performance on new, unseen data and to prevent overfitting. During training, the model is trained on the training set and its performance is evaluated on the validation set. By monitoring the performance of the model on the validation set, we can adjust the hyperparameters of the model to improve its performance.

The validation set is used to estimate the generalization error of the model and to select the best model among a set of candidate models. It is important to use a separate validation set instead of the test set because the test set should only be used once to assess the final performance of the trained model, while the validation set is used multiple times during training to select the best model.

Overall, the purpose of a validation set is to help prevent overfitting and to select the best model among a set of candidate models based on their performance on new, unseen data.

What can go wrong if you tune hyperparameters using the test set?

Tuning hyperparameters using the test set can lead to overfitting and biased performance estimates.

What is cross-validation and why would you prefer it to a validation set?

Cross-validation is a technique used in machine learning to evaluate the performance of a model and to select the best model among a set of candidate models. It involves partitioning the available data into multiple subsets, called folds, and using each fold in turn as a validation set while the model is trained on the remaining folds.

