

# Machine Learning #01

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

Machine learning is a subfield of artificial intelligence (AI) that involves developing algorithms and statistical models that enable computers to learn from and make predictions or decisions based on data, without being explicitly programmed to do so. In other words, it's a way of teaching computers to learn from examples and experience, instead of being told exactly what to do.

## ▼ Can you think of 4 distinct types of issues where it shines?

1. **Image and speech recognition:** Machine learning algorithms can be trained to recognize images and speech with a high degree of accuracy. This is useful in applications such as **facial recognition, object detection, and voice-controlled devices**.
2. **Natural language processing:** Machine learning can be used to analyze and understand human language, including **sentiment analysis, text classification, and language translation**. This is useful in applications such as **chatbots, language translation tools, and content moderation**.
3. **Predictive analytics:** Machine learning algorithms can be used to analyze large amounts of data and make predictions about future events or outcomes. This is useful in applications such as **fraud detection, customer churn prediction, and demand forecasting**.
4. **Recommender systems:** Machine learning can be used to build personalized recommendation systems that suggest products, services, or content to users based on their preferences and behavior. This is useful in applications such as **e-commerce, streaming services, and social media platforms**.

## ▼ What is a labeled training set, and how does it work?

A labeled training set **is a dataset that includes examples of input data, along with their corresponding output or "label" values**. In other words, each example in the dataset has been labeled with the correct answer or outcome that the machine learning algorithm is trying to learn.

For example, a labeled training set for a spam email filter might include thousands of emails that have been labeled as either "spam" or "not spam." **Each email is an example of input data, and its label indicates whether it is spam or not.**

**During the training phase, the machine learning algorithm uses the labeled training set to "learn" how to map input data to the correct output or label.** The algorithm analyzes the features of each example in the dataset, such as the words used in the email or the pixels in an image, and tries to identify patterns or relationships between the features and the labels.

**Once the algorithm has been trained on the labeled dataset, it can be used to make predictions on new, unseen data.** For example, in the case of the spam email filter, the trained algorithm can be used to predict whether a new email is likely to be spam or not, based on its features.

## ▼ What are the two most important tasks that are supervised?

The two most important tasks that are supervised are:

1. **Classification:** **Classification is the task of assigning a label or category to an input based on a set of predefined labels.** In supervised learning, the algorithm is trained on a labeled

dataset where the input examples are labeled with the correct output category. Once the algorithm has been trained, it can be used to classify new, unseen examples into one of the predefined categories. Examples of classification tasks include image classification, sentiment analysis, and spam detection.

2. **Regression:** *Regression is the task of predicting a continuous numerical value based on a set of input features.* In supervised learning, the algorithm is trained on a labeled dataset where the input examples are paired with the correct output values. Once the algorithm has been trained, it can be used to predict the output value for new, unseen examples based on their input features. Examples of regression tasks include stock price prediction, weather forecasting, and housing price prediction.

### ▼ Can you think of four examples of unsupervised tasks?

1. **Clustering:** *Clustering is the task of grouping similar data points together in a dataset.* The algorithm identifies patterns and similarities in the data, and then groups the data points into clusters based on those similarities. Clustering is often used in customer segmentation, image segmentation, and anomaly detection.
2. **Dimensionality reduction:** *Dimensionality reduction is the task of reducing the number of input features in a dataset while retaining as much useful information as possible.* This is often done to make it easier to visualize and analyze high-dimensional data. Principal Component Analysis (PCA) is a common technique for dimensionality reduction.
3. **Association rule learning:** *Association rule learning is the task of discovering relationships or associations between variables in a dataset.* This is often used in market basket analysis to identify which products are frequently purchased together, or in recommendation systems to identify which items are frequently viewed together.
4. **Anomaly detection:** *Anomaly detection is the task of identifying unusual or abnormal data points in a dataset.* The algorithm identifies patterns in the data and then flags data points that do not fit those patterns as potential anomalies. Anomaly detection is often used in fraud detection, network intrusion detection, and equipment failure detection.

### ▼ State the machine learning model that would be best to make a robot walk through various unfamiliar terrains?

**Reinforcement learning** is a type of machine learning where an **agent learns to interact with an environment in order to maximize a reward signal**. In this case, the robot would be the agent and the environment would be the various unfamiliar terrains.

**The reinforcement learning algorithm would receive feedback in the form of a reward signal** that indicates how well the robot is performing the task of walking through the terrain. **The algorithm would use this feedback to update its policy or set of actions to take in order to maximize the reward.**

Through trial and error, the reinforcement learning algorithm would learn which actions are most effective for navigating different terrains and would improve its performance over time.

### ▼ Which algorithm will you use to divide your customers into different groups?

To divide customers into different groups, **a clustering algorithm would be a good choice.**

**Clustering is an unsupervised learning technique used to group similar data points together based on their characteristics.** In this case, we would use clustering to group customers with similar

characteristics or behaviors together.

One popular clustering algorithm is **K-means clustering**. This algorithm starts by randomly assigning each data point to one of K clusters. It then iteratively refines the cluster assignments by calculating the distance between each data point and the centroid (center) of its assigned cluster, and then reassigning data points to the nearest centroid. The algorithm continues to refine the cluster assignments until the cluster assignments stabilize.

Another popular clustering algorithm is **hierarchical clustering**. This algorithm starts by treating each data point as its own cluster and then iteratively merges the closest clusters together based on a distance metric. The algorithm continues to merge clusters until all data points belong to a single cluster, or until a predefined stopping criterion is met.

Both K-means clustering and hierarchical clustering are commonly used for customer segmentation in marketing and e-commerce applications, as they can help businesses identify distinct groups of customers with similar characteristics or behaviors.

### ▼ Will you consider the problem of spam detection to be a supervised or unsupervised learning problem?

***Spam detection is typically considered a supervised learning problem.***

Supervised learning involves training a machine learning algorithm on a labeled dataset, where the input examples are paired with the correct output values. In the case of spam detection, the input examples are emails or messages, and the output values are binary labels indicating whether the message is spam or not.

**To train a spam detection model, a dataset of emails or messages would be labeled as either spam or not spam.** The algorithm would then learn to identify patterns in the data that distinguish between spam and non-spam messages. Once the algorithm has been trained, it can be used to classify new, unseen messages as spam or not spam based on their content.

Unsupervised learning, on the other hand, involves identifying patterns or relationships in a dataset without explicit labels. While unsupervised learning techniques such as clustering or anomaly detection could potentially be used to detect spam messages, these approaches are not as commonly used as supervised learning for this specific task.

### ▼ What is the concept of an online learning system?

***An online learning system is a type of machine learning system that is designed to continuously learn from new data as it arrives, in real-time or near real-time, without requiring a complete retraining of the model.*** This is also known as online machine learning or incremental machine learning.

In an online learning system, the model is trained on a stream of incoming data and the model's parameters are updated incrementally as new data arrives. ***This approach allows the model to adapt to changing data and to make predictions or classifications in real-time.***

Online learning systems are commonly used in applications where data is constantly changing or evolving, such as in **recommender systems, fraud detection, and dynamic pricing**. In these applications, it is important to be able to update the model quickly and efficiently in order to provide accurate predictions or recommendations.

One of the key benefits of online learning is its ability to handle large amounts of data in real-time, without requiring significant computing resources or storage. Additionally, ***online learning allows for***

*faster adaptation to changes in the data, as the model can immediately incorporate new data into its training.*

### ▼ What is out-of-core learning, and how does it differ from core learning?

**Out-of-core learning** is a type of machine learning that is designed *to handle very large datasets that cannot fit into memory* or that exceed the capacity of a single computer's processing power.

In out-of-core learning, **the data is stored in external storage**, such as hard disk drives or network-attached storage, and the learning algorithm processes the data in small, manageable chunks, or batches. The algorithm processes each batch of data sequentially and updates the model parameters incrementally as each batch is processed.

Out-of-core learning differs from **in-core learning**, also known as batch learning, where the **entire dataset is loaded into memory and processed at once**. In-core learning is typically faster than out-of-core learning, but it requires significant amounts of memory and computing resources, and may not be able to handle extremely large datasets.

Out-of-core learning **can be implemented using** a variety of algorithms, including **stochastic gradient descent**, which updates the model parameters after processing each individual data point or small batch of data. Other algorithms that are commonly used in out-of-core learning include incremental PCA (Principal Component Analysis), k-means clustering, and Naïve Bayes.

Out-of-core learning is particularly useful in applications where large amounts of data need to be processed and analyzed, such as in image or video processing, natural language processing, and recommendation systems. By processing data in small batches, out-of-core learning enables machine learning algorithms to learn from extremely large datasets while efficiently using available computing resources.

### ▼ What kind of learning algorithm makes predictions using a similarity measure?

A type of learning algorithm that makes predictions using a similarity measure is called a **"nearest neighbor"** algorithm.

Nearest neighbor algorithms are a type of instance-based learning, where the algorithm learns from specific examples **in the training data and uses those examples to make predictions on new, unseen data**. The algorithm uses a similarity measure, such as **Euclidean distance or cosine similarity**, to determine the closest or most similar instances in the training data to the new data point.

The algorithm can be used for both classification and regression tasks. **In classification tasks, the nearest neighbor algorithm assigns the class label of the closest training instance to the new data point. In regression tasks, the algorithm predicts the value of the dependent variable based on the values of the closest training instances.**

**One advantage of nearest neighbor algorithms is that they can be used for non-linear relationships and can handle high-dimensional data.** However, they can be computationally expensive, especially when dealing with large datasets, and may not perform well in high-dimensional spaces or when dealing with noisy or irrelevant features.

Overall, nearest neighbor algorithms are a useful tool in machine learning for tasks where similarity or distance is a meaningful measure for predicting outcomes.

### ▼ What's the difference between a model parameter and a hyperparameter in a learning algorithm?

In machine learning, **a model parameter is a configuration variable that is internal to the model and is learned from training data. Model parameters are adjusted during training using an optimization algorithm**, such as **gradient descent**, to minimize a loss function or maximize a likelihood function. Model parameters are specific to the model architecture and are learned during training, **so they are not set by the user**.

On the other hand, **a hyperparameter is a configuration variable that is external to the model and is set by the user before training begins. Hyperparameters control the behavior of the learning algorithm and affect how the model parameters are learned. Examples** of hyperparameters include **learning rate, regularization strength, number of hidden layers in a neural network, or the type of kernel used in a support vector machine**.

**The values of hyperparameters are typically set using a trial-and-error process**, where different values are tried and the performance of the model is evaluated on a validation set. The goal is to find the combination of hyperparameters that gives the best performance on the validation set.

**The key difference between model parameters and hyperparameters is that model parameters are learned during training, while hyperparameters are set before training begins and control the learning process itself.** Model parameters are specific to the model architecture and are learned from data, while hyperparameters are external to the model and affect how the model is learned.

### ▼ What are the criteria that model-based learning algorithms look for? What is the most popular method they use to achieve success? What method do they use to make predictions?

Model-based learning algorithms, **also known as parametric learning algorithms**, aim to learn a model that can be used to make predictions on new, unseen data. **They do this by estimating a set of parameters that define the model and allow it to make accurate predictions.**

The criteria that model-based learning algorithms look for are typically related to the goodness of fit between the model and the training data, and the ability of the model to generalize to new, unseen data. Some common criteria include:

1. Minimizing the difference between the predictions of the model and the actual values in the training data (e.g., minimizing the mean squared error).
2. Maximizing the likelihood of the training data given the model.
3. Minimizing the complexity of the model, to prevent overfitting and improve generalization to new data.

The most popular method that model-based learning algorithms use to achieve success is maximum likelihood estimation (MLE). MLE is a statistical method for estimating the parameters of a model that maximize the likelihood of the observed data. MLE assumes that the data is generated by a probability distribution that depends on the model parameters, and uses optimization algorithms, such as gradient descent, to find the parameter values that maximize the likelihood of the observed data.

To make predictions, model-based learning algorithms use the learned model to compute a prediction based on the input features of new, unseen data. The exact method used to make predictions depends on the specific model being used. For example, linear regression models use a linear combination of the input features to compute a prediction, while logistic regression models use a sigmoid function to transform the linear combination of the input features into a probability of belonging to a particular class.

### ▼ Can you name four of the most important Machine Learning challenges?

1. **Data quality:** Machine learning algorithms rely on large amounts of high-quality data to learn accurate models. However, in *many real-world scenarios, the available data is incomplete, noisy, or biased*. Ensuring data quality is crucial for the success of machine learning models.
2. **Overfitting:** *Overfitting occurs when a machine learning model is too complex and fits the training data too closely, leading to poor generalization to new, unseen data.* Avoiding overfitting is a major challenge in machine learning, and various techniques, such as regularization and early stopping, have been developed to address it.
3. **Interpretability:** As machine learning models become more complex, *it can become difficult to interpret their decisions and understand how they are making predictions.* Interpretability is an important challenge in machine learning, particularly in domains where decisions have significant real-world consequences, such as healthcare or finance.
4. **Scalability:** *Machine learning models can require significant computational resources, particularly for large datasets and complex models.* Scalability is a major challenge in machine learning, and techniques such as distributed computing, parallel processing, and hardware acceleration are used to address it.

▼ **What happens if the model performs well on the training data but fails to generalize the results to new situations? Can you think of three different options?**

If a machine learning model performs well on the training data but fails to generalize to new situations, *it is said to have overfit to the training data*. This is a common problem in machine learning, and there are several options for addressing it:

1. **Regularization:** *Regularization is a technique that adds a penalty term to the model's objective function to encourage simpler models that are less likely to overfit.* There are several types of regularization, including L1 regularization (which encourages sparse models), L2 regularization (which discourages large parameter values), and dropout regularization (which randomly drops out units in the model during training).
2. **Cross-validation:** *Cross-validation is a technique for evaluating a model's performance on a separate validation set that is not used for training.* By testing the model on new data that it has not seen before, *cross-validation can help identify whether a model is overfitting to the training data.*
3. **Feature selection:** *Feature selection is a technique for selecting a subset of the input features that are most relevant for predicting the target variable.* By removing irrelevant or redundant features, feature selection *can help reduce the complexity* of the model and improve its generalization performance.

Overall, *the goal of these techniques is to strike a balance between model complexity and generalization performance*, so that the model can perform well on both the training data and new, unseen data.

▼ **What exactly is a test set, and why would you need one?**

In machine learning, *a test set is a dataset that is used to evaluate the performance of a trained model on new, unseen data. The test set is separate from the training set, which is used to train the model, and the validation set, which is used to tune the model's hyperparameters.*

**The purpose of the test set is to provide an unbiased estimate of the model's performance on new data. By evaluating the model on a test set that is separate from the training set, we can get a more accurate estimate of how well the model will perform on new, unseen data.** This is important because the **ultimate goal of a machine learning model is to make accurate predictions on new data, not just the data that it was trained on.**

**Without a test set, it is difficult to know how well a model will perform on new data.** If we evaluate the model on the training set, we may get overly optimistic estimates of its performance because the model has already seen the training data. On the other hand, if we evaluate the model on the validation set, we may inadvertently tune the model's hyperparameters to perform well on the validation set, rather than generalizing well to new data.

By reserving a portion of the data for a test set, we can get a more accurate estimate of the model's performance on new, unseen data, and make more informed decisions about how to improve the model.

### ▼ What is a validation set's purpose?

In machine learning, **a validation set is a dataset that is used to evaluate the performance of a trained model during the training process.** The validation set is separate from the training set, which is used to train the model, and the test set, which is used to evaluate the model's performance on new, unseen data.

**The purpose of the validation set is to provide a way to tune the model's hyperparameters and prevent overfitting.** Hyperparameters are settings that are not learned during training, but instead are set by the user or the machine learning engineer. Examples of hyperparameters include the learning rate of the optimization algorithm, the number of hidden layers in a neural network, and the strength of regularization.

**During training, the model is updated based on the training data, but the validation set is used to evaluate the performance of the model on new data that it has not seen before.** By monitoring the model's performance on the validation set during training, we can adjust the hyperparameters to improve the model's performance and prevent overfitting.

If we were to use the test set to tune the hyperparameters, we would risk overfitting to the test set and getting overly optimistic estimates of the model's performance on new data. By using a separate validation set, we can avoid this problem and get a more accurate estimate of the model's performance on new, unseen data.

### ▼ What precisely is the train-dev kit, when will you need it, how do you put it to use?

The train-dev set, **also known as the development set, is a subset of the training set that is used to monitor the performance of a machine learning model during training. The purpose of the train-dev set is to detect overfitting and identify potential problems with the model before it is tested on the test set.**

The train-dev set is created by splitting the training set into two parts: the train set and the dev set. The train set is used to train the model, while the dev set is used to evaluate the model's performance during training.

The train-dev set is typically used in situations where the training set is large and diverse, and it is difficult to evaluate the model's performance on new data. In these situations, the dev set can be used



to identify potential problems with the model, such as overfitting, and make adjustments to the model to improve its performance.

To use the train-dev set, the model is trained on the train set, and its performance is evaluated on the dev set after each epoch of training. The performance on the dev set is used to make decisions about the model's hyperparameters and architecture, such as the learning rate, regularization strength, and the number of hidden layers in a neural network.

It is important to note that the train-dev set should be kept separate from the test set, which is used to evaluate the final performance of the model. Using the same data for both training and testing can lead to overfitting and biased estimates of the model's performance.

## ▼ What could go wrong if you use the test set to tune hyperparameters?

*If you use the test set to tune hyperparameters, you risk overfitting the model to the test set.*

*The purpose of the test set is to evaluate the final performance of the model on new, unseen data. If you use the test set to select hyperparameters or make other decisions during model development, you risk biasing your estimates of the model's performance on new data.*

Here are some specific problems that can arise if you use the test set to tune hyperparameters:

1. **Overfitting:** If you use the test set to select hyperparameters, you risk overfitting the model to the test set. This can lead to overly optimistic estimates of the model's performance on new data, which may not generalize well.
2. **Leakage:** *If you use the test set to make decisions during model development, you risk leaking information from the test set into the training process.* This can lead to biased estimates of the model's performance on new data, as the model has "seen" some of the test data during training.
3. **Limited sample size:** If you use the test set to tune hyperparameters, you are essentially reducing the size of your dataset, as you are using some of the data for both training and testing. This can lead to limited sample sizes for both training and testing, which can make it difficult to estimate the model's performance accurately.

To avoid these problems, it is important to keep the test set separate from the training and validation sets, and use the validation set to tune hyperparameters and make decisions during model development. Once the model is fully trained and tuned, it can be evaluated on the test set to get a final estimate of its performance on new, unseen data.