# What is supervised learning?

[Supervised learning](#) is a machine learning approach that's defined by its use of labeled datasets. These datasets are designed to train or "supervise" algorithms into classifying data or predicting outcomes accurately. Using labeled inputs and outputs, the model can measure its accuracy and learn over time.

Supervised learning can be separated into two types of problems when [data mining](#): classification and regression:

- **Classification** problems use an algorithm to accurately assign test data into specific categories, such as separating apples from oranges. Or, in the real world, supervised learning algorithms can be used to classify spam in a separate folder from your inbox. Linear classifiers, support vector machines, decision trees and [random forest](#) are all common types of classification algorithms.
- **Regression** is another type of supervised learning method that uses an algorithm to understand the relationship between dependent and independent variables. Regression models are helpful for predicting numerical values based on different data points, such as sales revenue projections for a given business. Some popular regression algorithms are linear regression, logistic regression and polynomial regression.

# What is unsupervised learning?

[Unsupervised learning](#) uses machine learning algorithms to analyze and cluster unlabeled data sets. These algorithms discover hidden patterns in data without the need for human intervention (hence, they are "unsupervised").

Unsupervised learning models are used for three main tasks: clustering, association and dimensionality reduction:

- **Clustering** is a data mining technique for grouping unlabeled data based on their similarities or differences. For example, K-means clustering algorithms assign similar data points into groups, where the K value represents the size of the grouping and granularity. This technique is helpful for market segmentation, image compression, etc.
- **Association** is another type of unsupervised learning method that uses different rules to find relationships between variables in a given dataset. These methods are frequently used for market basket analysis and recommendation engines, along the lines of "Customers Who Bought This Item Also Bought" recommendations.
- **Dimensionality reduction** is a learning technique used when the number of features  (or dimensions) in a given dataset is too high. It reduces the number of data inputs to a manageable size while also preserving the data integrity. Often,

this technique is used in the preprocessing data stage, such as when autoencoders remove noise from visual data to improve picture quality.

## The main difference between supervised and unsupervised learning: Labeled data

The main distinction between the two approaches is the use of labeled datasets. To put it simply, supervised learning uses labeled input and output data, while an unsupervised learning algorithm does not.

In supervised learning, the algorithm "learns" from the training dataset by iteratively making predictions on the data and adjusting for the correct answer. While supervised learning models tend to be more accurate than unsupervised learning models, they require upfront human intervention to label the data appropriately. For example, a supervised learning model can predict how long your commute will be based on the time of day, weather conditions and so on. But first, you'll have to train it to know that rainy weather extends the driving time.

Unsupervised learning models, in contrast, work on their own to discover the inherent structure of unlabeled data. Note that they still require some human intervention for validating output variables. For example, an unsupervised learning model can identify that online shoppers often purchase groups of products at the same time. However, a data analyst would need to validate that it makes sense for a recommendation engine to group baby clothes with an order of diapers, applesauce and sippy cups.

## Other key differences between supervised and unsupervised learning

- **Goals:** In supervised learning, the goal is to predict outcomes for new data. You know up front the type of results to expect. With an unsupervised learning algorithm, the goal is to get insights from large volumes of new data. The machine learning itself determines what is different or interesting from the dataset.
- **Applications**: Supervised learning models are ideal for spam detection, sentiment analysis, weather forecasting and pricing predictions, among other things. In contrast, unsupervised learning is a great fit for anomaly detection, recommendation engines, customer personas and medical imaging.
- **Complexity:** Supervised learning is a simple method for machine learning, typically calculated through the use of programs like R or Python. In unsupervised learning, you need powerful tools for working with large amounts of unclassified data. Unsupervised learning models are computationally complex because they need a large training set to produce intended outcomes.
- **Drawbacks**: Supervised learning models can be time-consuming to train, and the labels for input and output variables require expertise. Meanwhile,

unsupervised learning methods can have wildly inaccurate results unless you have human intervention to validate the output variables.

## Supervised vs. unsupervised learning: Which is best for you?

Choosing the right approach for your situation depends on how your data scientists assess the structure and volume of your data, as well as the use case. To make your decision, be sure to do the following:

- **Evaluate your input data:** Is it labeled or unlabeled data? Do you have experts that can support additional labeling?
- **Define your goals:** Do you have a recurring, well-defined problem to solve? Or will the algorithm need to predict new problems?
- **Review your options for algorithms:** Are there algorithms with the same dimensionality you need (number of features, attributes or characteristics)? Can they support your data volume and structure?

Classifying big data can be a real challenge in supervised learning, but the results are highly accurate and trustworthy. In contrast, unsupervised learning can handle large volumes of data in real time. But, there's a lack of transparency into how data is clustered and a higher risk of inaccurate results. This is where semi-supervised learning comes in.