

Machine Learning

BY DAVID A. PRICE

Customers of online music services have long been able to explore new music, or revisit old music, through the services' playlists. Whether you like '80s pop, '90s rap, or new country, your online music service has had a playlist for you, handmade by music experts. But in 2015, Spotify added something different: individually personalized playlists that each of its millions of users received every Monday. The feature, known as Discover Weekly, gained devotees. One wrote, "It felt like an intimate gift from someone who knew my tastes inside and out."

Of course, Spotify didn't scale up its staff of human music experts to create weekly playlists for what are now reportedly 87 million subscribers. Discover Weekly relies instead on a user's past listening habits and those of others with apparently similar tastes — and on machine learning software that converts this data into predictions of what a user would like.

Music is just one of a range of industries being affected by machine learning technology. Machine learning is likely to improve high-tech products in applications from spam filtering to face recognition. In medicine, machine learning may improve the interpretation of X-rays and other scans, as well as suggest diagnoses based on detailed patient information. Within the financial sector, some applications include detecting fraud, estimating insurance risks, and analyzing investments. In some industries, the adoption of machine learning may change the profile of skills sought by employers and even reduce employment numbers outright.

But what is it, exactly? Historically, it has a number of fields in its family tree: computer science, cognitive science, and statistics, among others. It's sometimes said to be a branch of artificial intelligence, or AI, but not the general, human-like AI seen in the fictional computers of *2001: A Space Odyssey* and *Star Trek*. Rather, it's a type of software that learns from examples — that is, it autonomously constructs models based on data fed into it. The data may represent transactions, images, or anything else in digital form.

Machine learning systems fall into one of two broad categories: supervised or unsupervised. In supervised machine learning, the system receives training data: a set of examples and information about the correct classification of each example. The latter is the "supervision." For instance, the training data could be images of furniture with information about whether each item is, say, a chair, a desk, or a sofa. With sufficient training data, the system would be able

to predict the correct category of an image of an item of furniture it hasn't seen before. Alternatively, the training data could be individuals' financial information, together with an indicator for each individual of whether he or she has a home mortgage default on record. The system would use that data to build a model for predicting whether a loan applicant is likely to default on a loan. (The person creating the system may hold back some of the data he or she has on hand to test the reliability of the model.)

In unsupervised machine learning, the system receives records, such as images or financial information, but no information on how to classify them. The task for the system is to discover categories within the data on its own.

In both supervised and unsupervised machine learning, the potential performance of the system improves as the system receives more data. Commonly, what goes into a machine learning system is an enormous dataset, so-called "big data," comprising millions of observations. Indeed, part of what has fueled the growth of machine learning is the availability of such datasets within technology companies as a byproduct of their operations as they capture data on transactions and other user behavior.

One important difference between machine learning and conventional techniques is that conventional statistical techniques produce models that can be interpreted by humans. Someone can look at the coefficients of a multiple regression analysis and see how it works — which variables count positively, which count negatively, and by how much. In contrast, complex machine learning models are like black boxes and cannot be translated into a form that lets humans understand the model's workings.

Within the discipline of economics, some researchers, such as Susan Athey of Stanford University, foresee that machine learning may become an increasingly important tool, transforming economic research. But for the time being, at least, switching from conventional statistical methods to machine learning comes at a price: Compared to machine learning, econometrics is better suited to asking about causation. Machine learning is about classification and prediction. Econometrics is too, but it also lets a researcher make inferences about whether and how one variable among many has been influencing the phenomenon that the researcher is studying. That distinction could erode, however, as researchers are seeking to combine machine learning with analysis of causation. **EF**

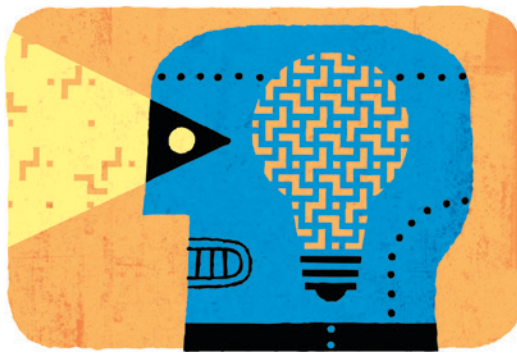


ILLUSTRATION: TIMOTHY COOK