# 16

# Artificial Intelligence, Economics, and Industrial Organization

Hal Varian

## 16.1 Introduction

Machine learning (ML) and artificial intelligence (AI) have been around for many years. However, in the last five years, remarkable progress has been made using multilayered neural networks in diverse areas such as image recognition, speech recognition, and machine translation. Artificial intelligence is a general purpose technology that is likely to impact many industries. In this chapter I consider how machine learning availability might affect the industrial organization of both firms that *provide* AI services and industries that *adopt* AI technology. My intent is not to provide an extensive overview of this rapidly evolving area, but instead to provide a short summary of some of the forces at work and to describe some possible areas for future research.

## 16.2 Machine-Learning Overview

Imagine we have a set of digital images along with a set of labels that describe what is depicted in those images—things like cats, dogs, beaches, mountains, cars, or people. Our goal is to use this data to train a computer

Hal Varian is an emeritus professor at the University of California, Berkeley, and chief economist at Google.

to learn how to predict labels for some new set of digital images. (For a nice demonstration, see cloud.google.com/vision where you can upload a photo and retrieve a list of labels appropriate for that photo.)

The classical approach to machine vision involved creating a set of rules that identified pixels in the images with human-recognizable features such as color, brightness, and edges and then use these features to predict labels. This "featurization" approach had limited success. The modern approach is to work directly with the raw pixels using layered neural networks. This has been remarkably successful, not only with image recognition but also with voice recognition, language translation, and other traditionally difficult machine-learning tasks. Nowadays computers can outperform humans in many of these tasks.

This approach, called *deep learning*, requires (a) labeled data for training, (b) algorithms for the neural nets, and (c) special-purpose hardware to run the algorithms. Academics and tech companies have provided training data and algorithms for free, and compute time in cloud-computing facilities is available for a nominal charge.

1. Training data. Examples are OpenImages, a 9.5 million data set of labeled images and the Stanford Dog Data set, 20,580 images of 120 breeds of dogs.

2. Algorithms. Popular open-source packages include TensorFlow, Caffe, MXNet, and Theano.

3. Hardware. CPUs (central processing units), GPUs (graphical processing units), and TPUs (Tensor processing units), are available via cloud-computing providers. These facilities allow the user to organize vast amounts of data, which can be used to train machine-learning models.

Of course, it is also important to have experts who can manage the data, tune the algorithms, and nurture the entire process. These skills are, in fact, the main bottleneck at the moment, but universities are rapidly rising to the challenge of providing the education and training necessary to create and utilize machine learning.

In addition to machine vision, the deep learning research community has made dramatic advances in speech recognition and language translation. These areas also have been able to make this progress without the sorts of feature identification that had been required for previous ML systems.

Other types of machine learning are described in the Wikipedia entry on this topic. One important form of machine learning is *reinforcement learning*. This is a type of learning where a machine optimizes some task such as winning at chess or video games. One example of reinforcement learning is a multiarmed bandit, but there are many other tools used, some of which involve deep neural nets.

Reinforcement learning is a type of sequential experimentation and is therefore fundamentally about causality: moving a particular chess piece

from one position to another *causes* the probability of a win to increase. This is unlike passive machine-learning algorithms that use only observational data.

Reinforcement learning can also be implemented in an adversarial context. For example, in October 2017 DeepMind announced a machine-learning system, Alpha Go 0, that developed a highly effective strategy by playing Go games against itself!

The model of "self-taught machine learning" is an interesting model for game theory. Can deep networks learn to compete and/or learn to cooperate with other players entirely their own? Will the learned behavior look anything like the equilibria for game-theoretic models we have built? So far these techniques have been applied primarily to full information games. Will they work in games with incomplete or asymmetric information?

There is a whole subarea of AI known as *adversarial AI* (or *adversarial ML*) that combines themes from AI, game theory, and computer security that examines ways to attack and defend AI systems. Suppose, for example, that we have a trained image recognition system that performs well, on average. What about its worst-case performance? It turns out that there are ways to create images that appear innocuous to humans that will consistently fool the ML system. Just as "optical illusions" can fool humans, these "ML illusions" can fool machines. Interestingly, the optimal illusions for humans and machines are very different. For some examples, see Goodfellow et al. (2017) for illustrative examples and Kurakin, Goodfellow, and Bengio (2016) for a technical report. Computer science researchers have recognized the connections with game theory; in my opinion, this area offers many interesting opportunities for collaboration. (See, e.g., Sreevallabh and Liu 2017).

### 16.2.1    What Can Machine Learning Do?

The example of machine learning presented in the popular press emphasizes novel applications, such as winning at games such as chess, Go, and Pong. However, there are also many practical applications that use machine learning to solve real-world business problems. A good place to see what kinds of problem ML can solve is Kaggle. This company sets up machine-learning competitions. A business or other organization provides some data, a problem statement, and some prize money. Data scientists then use the data to solve the problem posed. The winners get to take home the prize money. There are well over 200 competitions on the site. Here are a few of the most recent.

- Passenger Threats. Improve accuracy of Homeland Security threat recognition: $1,500,000.
- Home Prices. Improve accuracy of Zillow's home-price prediction: $1,200,000.

- Traffic to Wikipedia Pages. Forecast future traffic to Wikipedia pages: $25,000.
- Personalized Medicine. Predict effect of genetic variants to enable personalized medicine: $15,000.
- Taxi Trip Duration. Predict total ride duration of taxi trips in New York: $30,000.
- Product Search Relevance. Predict relevance of search results on homedepot.com: $40,000.
- Clustering Questions. Can you identify question pairs that have the same intent?: $25,000.
- Cervical cancer screening. Which cancer treatments will be most effective?: $100,000.
- Click Prediction. Can you predict which recommended content each user will click?: $25,000.
- Inventory Demand. Maximize sales and minimize returns of bakery goods: $25,000.

What is nice is that these are real questions and real money from organizations that want real answers for real problems. Kaggle gives concrete examples of how machine learning can be applied for practical business questions.[1]

16.2.2    What Factors Are Scarce?

Suppose you want to deploy a machine-learning system in your organization. The first requirement is to have a data infrastructure that collects and organizes the data of interest—a *data pipeline*. For example, a retailer would need a system that can collect data at point of sale, and then upload it to a computer that can then organize the data into a database. This data would then be combined with other data, such as inventory data, logistics data, and perhaps information about the customer. Constructing this data pipeline is often the most labor intensive and expensive part of building a data infrastructure, since different businesses often have idiosyncratic legacy systems that are difficult to interconnect.

Once the data has been organized, it can be collected together to in a data warehouse. The data warehouse allows easy access to systems that can manipulate, visualize, and analyze the data.

Traditionally, companies ran their own data warehouses that required not only purchase of costly computers, but also required human system administrators to keep everything functioning properly. Nowadays, it is more and more common to store and analyze the data in a cloud-computing facility

---

1. Disclosure: I was an angel investor in Kaggle up till mid-2017 when it was acquired by Google. Since then, I have had no financial interest in the company.

such as Amazon Web Services, Google Cloud Platform, or Microsoft Azure Cloud.

The cloud provider takes care of managing and updating the hardware and software necessary to host the databases and tools for data analysis. From an economic point of view, what is interesting is that what was previously a fixed cost to the users (the data center) has now turned into a variable cost (renting time on the data center). An organization can purchase virtually any amount of cloud services, so even small companies can start at a minimal level and be charged based on usage. Cloud computing is much more cost effective than owning your own data center, since compute and data resources can be purchased on an as-needed basis. Needless to say, most tech start-ups today use a cloud provider for their hardware, software, and networking needs.

Cloud providers also offer various machine-learning services such as voice recognition, image recognition, translation, and so on. These systems are already trained by the vendor and can be put to immediate use by the customer. It is no longer necessary for each company to develop its own software for these tasks.

Competition among the cloud providers is intense. Highly detailed and specific image recognition capabilities are offered at a cost of a tenth-of-a-cent per image or less, with volume discounts on top of that price.

A user may also have idiosyncratic data relevant to its own business like the point-of-sale data mentioned above. The cloud provider also provides up-to-date, highly optimized hardware and software than implements popular machine-learning algorithms. This allows the use immediate access to high-powered tools . . . providing that they have the expertise to use them.

If the hardware, software, and expertise are available, all that is needed is the labeled data. There are a variety of ways to acquire such data.

- As By-Product of Operations. Think of a chain of restaurants where some perform better than others, and management may be interested in factors that are associated with performance. Much of the data in the Kaggle competitions mentioned above are generated as a byproduct of day-to-day operations.
- Web Scraping. This is a commonly used way to extract data from websites. There is a legal debate about what exactly is permitted with respect to both the collection of data and how it is used. The debate is too complex to discuss here, but the Wikipedia entry on Web scraping is good. An alternative is to use data that others have scraped. For example, the Common Crawl database contains petabytes of data compiled over eight years of Web crawling.
- Offering a Service. When Google started its work on voice recognition, it had no expertise and no data. It hired the expertise and they came up

with the idea of a voice-input telephone directory as a way to acquire data. Users would say "Joe's Pizza, University Avenue, Palo Alto" and the system would respond with a phone number. The digitized question and the resulting user choices were uploaded to the cloud and machine learning was used to evaluate the relationship between Google's answer and the user action—for example, to call the suggested number. The ML training used data from millions of individual number requests and learned rapidly. ReCAPTCHA applies a similar model where humans label images to prove they are human and not a simple bot.

- Hiring Humans to Label Data. Mechanical Turk and other systems can be used to pay people to label data (see Hutson 2017).
- Buying Data from Provider. There are many providers of various sorts of data such as mail lists, credit scores, and so on.
- Sharing Data. It may be mutually advantageous to parties to share data. This is common among academic researchers. The Open Images Data set contains about nine million labeled images contributed by universities and research labs. Sharing may be mandated for a variety reasons, such as concerns for public safety. Examples are black boxes from airplanes or medical data on epidemics.
- Data from Governments. There are vast amounts of data available from governments, universities, research labs, and nongovernmental agencies.
- Data from Cloud Providers. Many cloud providers also provide public data repositories. See, for example, Google Public Data sets, Google Patents Public Data set, or AWS Public Data sets.
- Computer-Generated Data. The Alpha Go 0 system mentioned earlier generated its own data by playing Go games against itself. Machine-vision algorithms can be trained using "synthetic images," which are actual images that have been shifted, rotated, and scaled in various ways.

### 16.2.3    Important Characteristics of Data

Information science uses the concept of a "data pyramid" to depict the relationship between data, information, and knowledge. Some system has to collect the raw data, and subsequently organize and analyze that data in order to turn it into information—something such as a textual document image that can be understood by humans. Think of the pixels in an image being turned into human-readable labels. In the past this was done by humans; in the future more and more of this will be done by machines. (See figure 16.1.)

This insights from the information can then turned into knowledge, which generally is embodied in humans. We can think of data being stored in bits, information stored in documents, and knowledge stored in humans. There are well-developed markets and regulatory environments for information
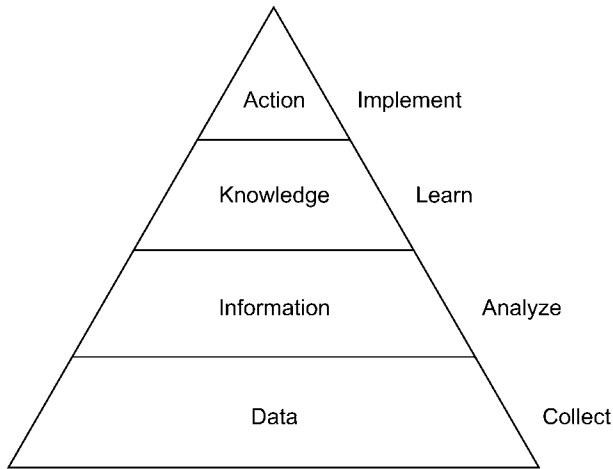
**Fig. 16.1    The information pyramid**

(books, articles, web pages, music, videos) and for knowledge (labor markets, consultants). Markets for data—in the sense of unorganized collections of bits—are not as developed. Perhaps this is because raw data is often heavily context dependent and is not very useful until it is turned into information.

*Data Ownership and Data Access*

It is said that "data is the new oil." Certainly, they are alike in one respect: both need to be refined in order to be useful. But there is an important distinction: oil is a *private good* and consumption of oil is *rival*: if one person consumes oil, there is less available for someone else to consume. But data is *nonrival:* one person's use of data does not reduce or diminish another person's use.

So instead of focusing on data "ownership"—a concept appropriate for private goods—we really should think about data access. Data is rarely "sold" in the same way private goods are sold, rather it is licensed for specific uses. Currently there is a policy debate in Europe about "who should own autonomous vehicle data?" A better question is to ask "who should have access to autonomous vehicle data and what can they do with it?" This formulation emphasizes that many parties can simultaneously access autonomous vehicle data. In fact, from the viewpoint of safety it seems very likely that multiple parties should be allowed to access autonomous vehicle data. There could easily be several data collection points in a car: the engine, the navigation system, mobile phones in rider's pockets, and so on. Requiring exclusivity without a good reason for doing so would unnecessarily limit what can be done with the data.

Ross Anderson's description of what happens when there is an aircraft

crash makes an important point illustrating why it may be important to allow several parties to access data.

> When an aircraft crashes, it is front page news. Teams of investigators rush to the scene, and the subsequent enquiries are conducted by experts from organisations with a wide range of interests—the carrier, the insurer, the manufacturer, the airline pilots' union, and the local aviation authority. Their findings are examined by journalists and politicians, discussed in pilots' messes, and passed on by flying instructors. In short, the flying community has a strong and institutionalised learning mechanism. (Anderson 1993)

Should we not want the same sort of learning mechanism for autonomous vehicles? Some sorts of information can be protected by copyright. But in the United States, raw data such as a telephone directory is not protected by copyright. (See Wikipedia entry on the legal case *Feist Publications, Inc v. Rural Telephone Service Co*.)

Despite this, data providers may compile some data and offer to license on certain terms to other parties. For example, there are several data companies that merge US census data with other sorts of geographic data and offer to license this data. These transactions may prohibit resale or relicensing. Even though there is no protectable intellectual property, the terms of the contract form a private contract that can be enforced by courts, as with any other private contract.

*Decreasing Marginal Returns*

Finally, it is important to understand that data typically exhibits decreasing returns to scale like any other factor of production. The same general principle applies for machine learning. Figure 16.2 shows how the accuracy of the Stanford dog breed classification behaves as the amount of training data increases. As one would expect, accuracy improves as the number of training images increases, but it does so at a decreasing rate.

Figure 16.3 shows how the error rate in the ImageNet competition has declined over the last several years. An important fact about this competition is that the number of training and test observations has been fixed during this period. This means that the improved performance of the winning systems cannot depend on sample size since it has been constant. Other factors such as improved algorithms, improved hardware, and improved expertise have been much more important than the number of observations in the training data.

## 16.3   Structure of ML-Using Industries

As with any new technology, the advent of machine learning raises several economic questions.
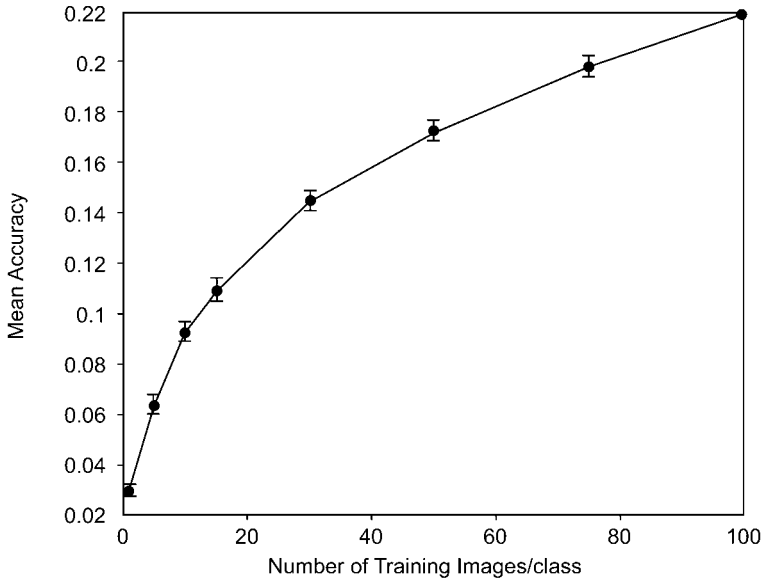
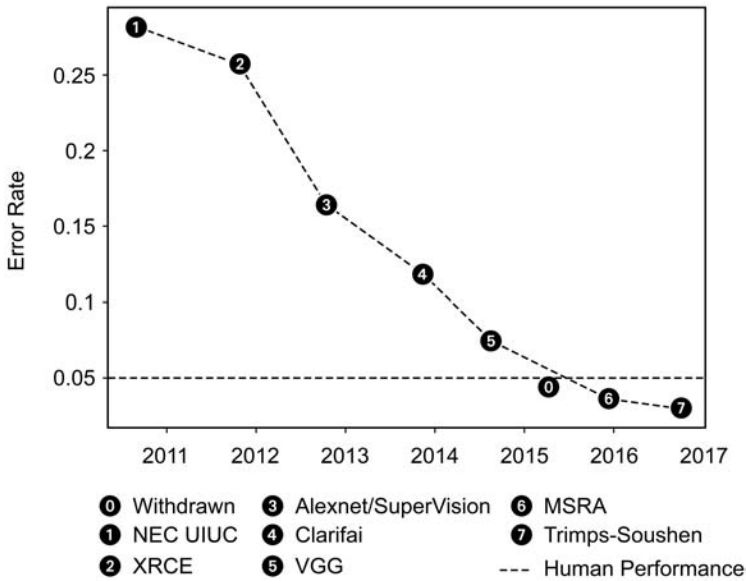**Fig. 16.2  Machine-learning adoption by economic sector**
*Source:* http://vision.stanford.edu/aditya86/ImageNetDogs/.



**Fig. 16.3  Imagenet image recognition**
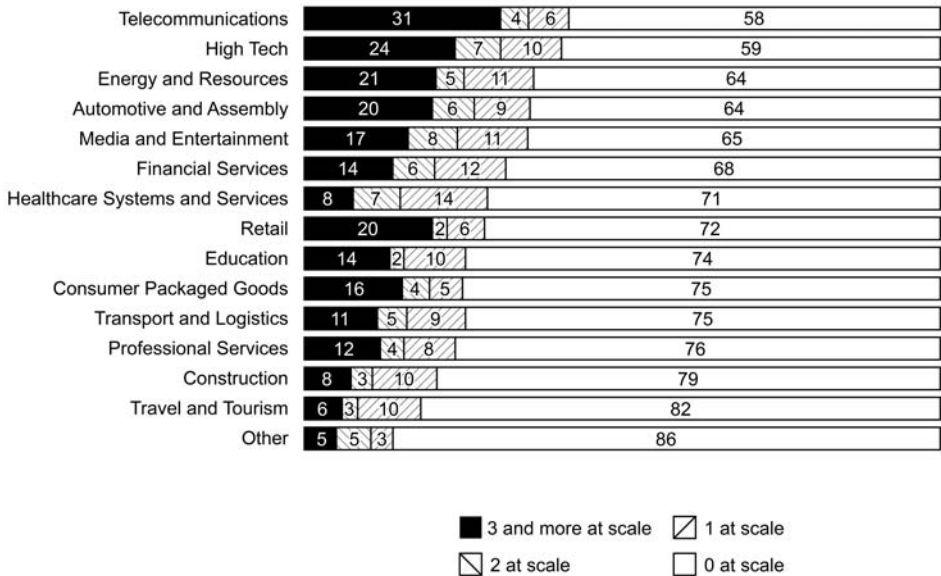*Source:* Eckersley and Nasser (2017).

Fig. 16.4    **Number of AI-related technologies adopted at scale or in a core part of the business**
*Source:* McKinsey (2017).

- Which firms and industries will successfully adopt machine learning?
- Will we see heterogeneity in the timing of adoption and the ability to use ML effectively?
- Can later adopters imitate early adopters?
- What is the role of patents, copyright, and trade secrets?
- What is the role of geography in adoption patterns?
- Is there a large competitive advantage for early, successful adopters?

Bughin and Hazan (2017) recently conducted a survey of 3,000 "AI Aware" C-level executives about adoption readiness. Of these executives, 20 percent are "serious adopters," 40 percent are "experimenting," and 28 percent feel their firms "lack the technical capabilities" to implement ML. McKinsey identifies key enablers of adoption to be leadership, technical ability, and data access. Figure 16.4 breaks down how ML adoption varies across economic sectors. Not surprisingly, sectors such as telecom, tech, and energy are ahead of less tech-savvy sectors such as construction and travel.

### 16.3.1    Machine Learning and Vertical Integration

A key question for industrial organization is how machine-learning tools and data can be combined to create value. Will this happen within or across corporate boundaries? Will ML users develop their own ML capabilities or purchase ML solutions from vendors? This is the classic make versus buy

question that is the key to understanding much of real-world industrial organization.

As mentioned earlier, cloud vendors provide integrated hardware and software environments for data manipulation and analysis. They also offer access to public and private databases, provide labeling services, consulting, and other related services that enable one-stop shopping for data manipulation and analysis. Special-purpose hardware provided by cloud providers such as GPUs and TPUs have become key technologies for differentiating provider services.

As usual there is a tension between standardization and differentiation. Cloud providers are competing intensely to provide standardized environments that can be easily maintained. At the same time, they want to provide services that differentiate their offerings from competitors.

Data manipulation and machine learning are natural areas to compete with respect to product speed and performance.

### 16.3.2    Firm Size and Boundaries

Will ML increase or decrease minimum efficient scale? The answer depends on the relationship between fixed costs and variable costs. If firms have to spend significant amounts to develop customized solutions to their problems, we might expect that fixed costs are significant and firm size must be large to amortize those costs. On the other hand, if firms can buy off-the-shelf services from cloud vendors, we would expect that fixed costs and minimum efficient scale to be small.

Suppose, for example, that an oil change service would like to greet returning customers by name. They can accomplish this using a database that joins license plate numbers with customer names and service history. It would be prohibitively expensive for a small provider to write the software to enable this, so only the large chains could provide such services. On the other hand, a third party might develop a smartphone app that could provide this service for a nominal cost. This service might allow minimum efficient scale to decrease. The same considerations apply for other small service providers such as restaurants, dry cleaners, or convenience stores.

Nowadays new start-ups are able to outsource a variety of business processes since there are a several providers of business services. Just as fast-food providers could perfect a model with a single establishment and then go national, business service companies can build systems once and replicate them globally.

Here is a list of how a start-up might outsource a dozen business processes.

- Fund your project on Kickstarter.
- Cloud cloud computing and network from Google, Amazon, or Micro-Soft.
- Use open-source software like Linux, Python, Tensorflow, and so forth.

- Manage your software using GitHub.
- Become a micromultinational and hire programmers from abroad.
- Set up a Kaggle competition for machine learning.
- Use Skype, Hangouts, Google Docs, and so forth for team communication.
- Use Nolo for legal documents (company, patents, NDAs).
- Use QuickBooks for accounting.
- Use AdWords, Bing, or Facebook for marketing.
- Use ZenDesk for user support.

This is only a partial list. Most start-ups in Silicon Valley and SOMA avail themselves of several of these business-process services. By choosing standardizing business processes, the start-ups can focus on their core competency and purchases services as necessary as they scale. One would expect to see more entry and more innovation as a result of the availability of these business-process services.

### 16.3.3   Pricing

The availability of cloud computing and machine learning offers lots of opportunities to adjust prices based on customer characteristics. Auctions and other novel pricing mechanisms can be implemented easily. The fact that prices can be so easily adjusted implies that various forms of differential pricing can be implemented. However, it must be remembered that customers are not helpless; they can also avail themselves of enhanced search capabilities. For example, airlines can adopt strategies that tie purchase price to departure date. But services can be created that reverse-engineer the airline algorithms and advise consumers about when to purchase (see, e.g., Etzioni et al. (2003). See Acquisti and Varian (2005) for a theoretical model of how consumers might respond to attempts to base prices on consumer history and how the consumers can respond to such attempts.

### 16.3.4 Price Differentiation

Traditionally, price differentiation has been classified into three categories:

1. First degree (personalized),
2. second degree (versioning: same price menu for all consumers, but prices vary with respect to quantity or quality), and
3. third degree (group pricing based on membership).

Fully personalized pricing is unrealistic, but prices based on fine-grained features of consumers may well be feasible, so the line between third degree and first degree is becoming somewhat blurred. Shiller (2013) and Dubé and Misra (2017) have investigated how much consumer surplus can be extracted using ML models.

Second-degree price discrimination can also be viewed as pricing by

group membership, but recognizing the endogeneity of group membership and behavior. Machine learning using observational data will be of limited help in designing such pricing schemes. However, reinforcement learning techniques such as multiarmed bandits may also be helpful.

According to most noneconomics, the only thing worse than price differentiation is price discrimination! However, most economists recognize that price differentiation is often beneficial from both an efficiency and an equity point of view. Price differentiation allows markets to be served that would otherwise not be served and often those unserved markets involve low-income consumers.

DellaVigna and Gentzkow (2017) suggest that "the uniform pricing we document significantly increases the prices paid by poorer households relative to the rich." This effect can be substantial. The authors show that "consumers of [food] stores in the lowest income decile pay about 0.7 percent higher prices than they would pay under flexible pricing, but consumers of stores in the top income decile pay about 9.0 percent lower prices than under flexible pricing."

### 16.3.5   Returns to Scale

There are at least three types of returns to scale that could be relevant for machine learning.

1. Classical supply-side returns to scale (decreasing average cost).
2. Demand-side returns to scale (network effects).
3. Learning by doing (improvement in quality or decrease in cost due to experience).

*Supply-Side Marginal Returns*

It might seem like software is the paradigm case of supply-side returns to scale: there is a large fixed cost of developing the software, and a small variable cost of distributing it. But if we compare this admittedly simple model to the real world, there is an immediate problem.

Software development is not a one-time operation; almost all software is updated and improved over time. Mobile phone operating systems are a case in point: there are often monthly releases of bug fixes and security improvements coupled with yearly releases of major upgrades.

Note how different this is from physical goods—true, there are bug fixes for mechanical problems in a car, but the capabilities of the car remain more or less constant over time. A notable exception is the Tesla brand, where new updated operating systems are released periodically.

As more and more products become network enabled we can expect to see this happen more often. Your TV, which used to be a static device, will be able to learn new tricks. Many TVs now have voice interaction, and we can expect that machine learning will continue to advance in this area. This

means that your TV will become more and more adept at communication and likely will become better at discerning your preferences for various sorts of content. The same goes for other appliances—their capabilities will no longer be fixed at time of sale, but will evolve over time.

This raises interesting economic questions about the distinction between goods and services. When someone buys a mobile phone, a TV, or a car, they are not just buying a static good, but rather a device that allows them to access a whole panoply of services. This, in turn, raises a whole range of questions about pricing and product design.

*Demand-Side Returns to Scale*

Demand-side economies of scale, or network effects, come in different varieties. There are *direct network effects*, where the value of a product or service to an incremental adopter depends on the total number of other adopters, and there are *indirect network effects* where there are two or more types of complementary adopters. Users prefer an operating system with lots of applications and developers prefer operating systems with lots of users.

Direct network effects could be relevant to choices of programming languages used in machine-learning systems, but the major languages are open source. Similarly, it is possible that prospective users might prefer cloud providers that have a lot of other users. However, it seems to me that this is no different than many other industries. Automobile purchasers may well have a preference for popular brands since dealers, repair shops, parts, and mechanics are readily available.

There is a concept that is circulating among lawyers and regulators called "data network effects." The model is that a firm with more customers can collect more data and use this data to improve its product. This is often true—the prospect of improving operations is what makes ML attractive—but it is hardly novel. And it is certainly not a network effect! This is essentially a supply-side effect known as "learning by doing" (also known as the "experience curve" or "learning curve"). The classical exposition is Arrow (1962); Spiegel and Hendel (2014) contain some up-to-date citations and a compelling example.

*Learning by Doing*

Learning by doing is generally modeled as a process where unit costs decline (or quality increases) as cumulative production or investment increases. The rough rule of thumb is that a doubling of output leads to a unit cost decline of 10 to 25 percent. Though the reasons for this efficiency increase are not firmly established, the important point is that learning by doing requires intention and investment by the firm and described in Stiglitz and Greenwald (2014).

This distinguishes learning by doing from demand-side or supply-side network effects that are typically thought to be more or less automatic.

This is not really true either; entire books have been written about strategic behavior in the presence of network effects. But there is an important difference between learning by doing and so-called "data network effects." A company can have huge amounts of data, but if it does nothing with the data it produces no value.

In my experience the problem is not lack of resources but lack of skills. A company that has data but no one to analyze it is in a poor position to take advantage of that data. If there is no existing expertise internally, it is hard to make intelligent choices about what skills are needed and how to find and hire people with those skills. Hiring good people has always been a critical issue for competitive advantage. But since the widespread availability of data is comparatively recent, this problem is particularly acute. Automobile companies can hire people who know how to build automobiles, since that is part of their core competency. They may or may not have sufficient internal expertise to hire good data scientists, which is why we can expect to see heterogeneity in productivity as this new skill percolates through the labor markets. Bessen (2016, 2017) has written perceptively about this issue.

### 16.3.6    Algorithmic Collusion

It has been known for decades that there are many equilibrium in repeated games. The central result in this area is the so-called "folk theorem," which says that virtually any outcome can be achieved as an equilibrium in a repeated game. For various formulations of this result, see the surveys by Fudenberg (1992) and Pierce (1992).

Interaction of oligopolists can be viewed as a repeated game, and in this case particular attention is focused on collusive outcomes. There are very simple strategies that can be used to facilitate collusion.

*Rapid Response Equilibrium.* For example, consider the classic example of two gas stations across the street from each other who can change prices quickly and serve a fixed population of consumers. Initially, they are both pricing above marginal cost. If one drops its price by a penny, the other quickly matches the price. In this case, both gas stations do worse off because they are selling at a lower price. Hence, there is no reward to price cutting and high prices prevail. Strategies of this sort may have been used in online competition, as described in Varian (2000). Borenstein (1997) documents related behavior in the context of airfare pricing.

*Repeated Prisoner's Dilemma.* In the early 1980s, Robert Axelrod (1984) conducted a prisoner's dilemma tournament. Researches submitted algorithmic strategies that were played against each other repeatedly. The winner by a large margin was a simple strategy submitted by Anatol Rapoport called "tit for tat." In this strategy, each side starts out cooperating (charging high prices). If either player defects (cuts its price), the other player matches. Axelrod then constructed a tournament where strategies reproduced according to their payoffs in the competition. He found that the best-performing

strategies were very similar to tit for tat. This suggests that artificial agents might learn to play cooperative strategies in a classic duopoly game.

*NASDAQ Price Quotes*. In the early 1990s, price quotes in the NASDAQ were made in eighths of a dollar rather than cents. So if a bid was three-eighths and an ask was two-eighths, a transaction would occur with the buyer paying three-eighths and the seller receiving two-eighths. The difference between the bid and the ask was the "inside spread," which compensated the traders for risk bearing and maintaining the capital necessary to participate in the market. Note that the bigger the inside spread, the larger the compensation to the market makers doing the trading.

In the mid-1990s two economists, William Christie and Paul Schultz, examined trades for the top seventy listed companies in NASDAQ and found to their surprise that there were virtually no transactions made at odd-eighth prices. The authors concluded that "our results most likely reflected an understanding or implicit agreement among the market makers to avoid the use of odd-eighth price fractions when quoting these stocks" (Christie and Schultz 1995, 203).

A subsequent investigation was launched by the Department of Justice (DOJ), which was eventually settled by a $1.01 billion fine that, at the time, was the largest fine ever paid in an antitrust case.

As these examples illustrate, it appears to be possible for implicit (or perhaps explicit) cooperation to occur in the context of repeated interaction—what Axelrod refers to as the "evolution of cooperation."

Recently, issues of these sort have reemerged in the context of "algorithmic collusion." In June 2017, the Organisation for Economic Co-operation and Development (OECD) held a roundtable on algorithms and collusion as a part of their work on competition in the digital economy. See OECD (2017) for a background paper and Ezrachi and Stucke (2017) for a representative contribution to the roundtable.

There are a number of interesting research questions that arise in this context. The folk theorem shows that collusive outcomes can be an equilibrium of a repeated game, but does not describe a specific algorithm that leads to such an outcome. It is known that very simplistic algorithms, such as finite automata with a small number of states cannot discover all equilibria (see Rubinstein 1986).

There are auction-like mechanisms that can be used to approximate monopoly outcomes; see Segal (2003) for an example. However, I have not seen similar mechanisms in an oligopoly context.

## 16.4   Structure of ML-Provision Industries

So far we have looked at industries that *use* machine learning, but it is also of interest to look at companies that *provide* machine learning.

As noted above, it is likely that ML vendors will offer several related ser-

vices. One question that immediately rises is how easy it will be to switch among providers. Technologies such as containers have been developed specifically make it easy to port applications from one cloud provider to another. Open-source implementation such as dockers and kubernetes are readily available. Lock in will not be a problem for small- and medium-size applications, but of course, there could be issues involving large and complex applications that involve customized applications.

Computer hardware also exhibits at least constant returns to scale due to the ease of replicating hardware installations at the level of the chip, motherboard, racks, or data centers themselves. The classic replication argument for constant returns applies here since the basic way to increase capacity is to just replicate what has been done before: add more core to the processors, add more boards to racks, add more racks to the data center, and build more data centers.

I have suggested earlier that cloud computing is more cost effective for most users than building a data center from scratch. What is interesting is that companies that require lots of data processing power have been able to replicate their existing infrastructure and sell the additional capacity to other, smaller entities. The result is an industry structure somewhat different than an economist might have imagined. Would an auto company build excess capacity that it could then sell off to other companies? This is not unheard of, but it is rare. Again it is the general purpose nature of computing that enables this model.

### 16.4.1   Pricing of ML Services

As with any other information-based industry, software is costly to produce and cheap to reproduce. As noted above, computer hardware also exhibits at least constant returns to scale due to the ease of replicating hardware installations at the level of the chip, motherboard, racks, or data centers themselves.

If services become highly standardized, then it is easy to fall into Bertrand-like price cutting. Even in these early days, machine pricing appears to be intensely competitive. For example, image recognition services cost about a tenth-of-a-cent per image at all major cloud providers. Presumably, we will see vendors try to differentiate themselves along dimensions of speed and capabilities. Those firms that can provide better services may be able to provide premium prices, to the extent that users are willing to pay for premium service. However, current speeds and accuracy are very high and it is unclear how users value further improvement in these dimensions.

## 16.5   Policy Questions

We have already discussed issues involving data ownership, data access, differential pricing, returns to scale, and algorithmic collusion, all of which

have significant policy aspects. The major policy areas remaining are security and privacy. I start with a few remarks about security.

### 16.5.1   Security

One important question that arises with respect to security is whether firms have appropriate incentives in this regard. In a classic article, Anderson (1993) compares US and UK policy with respect to automatic teller machines (ATMs). In the United States, the user was right unless the bank could prove them wrong, while in the United Kingdom, the bank was right unless the user could prove them wrong. The result of this liability assignment was that US banks invested in security practices such as security cameras, while the UK banks didn't bother with such elementary precautions.

This industry indicates how important liability assignment is in creating appropriate incentives for investment in security. The law and economics analysis of tort law is helpful in understanding the implications of different liability assignments and what optimal assignments might look like.

One principle that emerges is that of the "due care" standard. If a firm follows certain standard procedures such as installing security fixes within a few days of their being released, implementing two-factor authentication, educating their workforce about security practices, and so on, they have a safe harbor with respect to liability for costs associated with security incidents.

But where does the due care standard come from? One possibility is from the government, particularly from military or law enforcement practices. The Orange Book and its successor, the Common Criteria standard, are good examples. Another possibility is that insurance agencies offer insurance to parties that implement good practices security. Just as an insurer may require a sprinkler system to offer fire insurance, cyber insurance may only be offered to those companies that engage in best practices (see Varian 2000 for more discussion).

This model is an appealing approach to the problem. However, we know that there are many issues involving insurance such as adverse selection and moral hazard that need to be addressed. See the archives of the Workshop on the Economics of Information Security for more work in this area, and Anderson (2017) for an overview.

### 16.5.2   Privacy

Privacy policy is a large and sprawling area. Acquisti, Taylor, and Wagman (2016) provide a comprehensive review of the economic literature.

There are several policy questions that arise in the machine-learning area. For example, do firms have appropriate incentives to provide appropriate levels of privacy? What is the trade-off between privacy and economic performance? It is widely recognized that privacy regulations may limit ability of ML vendors to combine data from multiple sources and there may be

limits on transfer of data across corporate boundaries and/or sale of data. There is a tendency to promulgate regulation in this area that leads to unintended consequences. An example is the Health Insurance Portability and Accountability Act of 1996, commonly known as HIPAA. The original intent of the legislation was to stimulate competition among insurers by establishing standards for medical record keeping. However, many researchers argue that it has had a significant negative impact on the quantity and quality of medical research.

### 16.5.3   Explanations

European regulators are examining the idea of a "right to an explanation." Suppose information about a consumer is fed into a model to predict whether or not he or she will default on a loan. If the consumer is refused the loan, are they owed an "explanation" of why? If so, what would count as an explanation? Can an organization keep a predictive model secret because if it were revealed it could be manipulated? A notable example is the Discriminant Inventory Function. better known as the DIF function that the IRS uses to trigger audits. Is it legitimate to reverse engineer the DIF function? See CAvQM (2011) for a collection of links on the DIF function.

Can we demand more of an ML model than we can of a person? Suppose we show you a photo and that you correctly identify it as a picture of your spouse. Now we ask, "how do you know?" The best answer might be "because I've seen a lot of pictures that I know are pictures of my spouse, and that photo looks a lot like those pictures!" Would this explanation be satisfactory coming from a computer?

### 16.6   Summary

This chapter has only scratched the surface of how AI and ML might impact industrial structure. The technology is advancing rapidly, with the main bottleneck now being analysts who can implement these machine-learning systems. Given the huge popularity of college classes in this area and the wealth of online tutorials, we expect this bottleneck will be alleviated in the next few years.

## References

Acquisti, Alessandro, Curtis R. Taylor, and Liad Wagman. 2016. "The Economics of Privacy." *Journal of Economic Literature* 52 (2).

Acquisti, Alessandro, and Hal Varian. 2004. "Conditioning Prices on Purchase History." *Marketing Science* 24 (4): 367–81.

Anderson, Ross. 1993. "Why Cryptosystems Fail." *Proceedings of the 1st ACM Conference on Computer and Communications Security*. https://dl.acm.org/citation.cfm?id=168615.

———. 2017. "Economics and Security Resource Page." Working paper, Cambridge University. http://www.cl.cam.ac.uk/~rja14/econsec.html.

Arrow, Kenneth J. 1962. "The Economic Implications of Learning by Doing." *Review of Economic Studies* 29 (3): 155–73.

Axelrod, Robert. 1984. *The Evolution of Cooperation*. New York: Basic Books.

Bessen, James. 2016. *Learning by Doing: The Real Connection between Innovation, Wages, and Wealth*. New Haven, CT: Yale University Press.

———. 2017. "Information Technology and Industry." Law and Economics Research Paper no. 17-41, Boston University School of Law.

Borenstein, Severin. 1997. "Rapid Communication and Price Fixing: The Airline Tariff Publishing Company Case." Working paper. http://faculty.haas.berkeley.edu/borenste/download/atpcase1.pdf.

Bughin, Jacques, and Erik Hazan. 2017. "The New Spring of Artificial Intelligence." Vox CEPR Policy Portal. https://voxeu.org/article/new-spring-artificial-intelligence-few-early-economics.

CavQM. 2011. "Reverse Engineering The IRS DIF-Score." Comparative Advantage via Quantitative Methods blog, July 10. http://cavqm.blogspot.com/2011/07/reverse-engineering-irs-dif-score.html.

Christie, William G., and Paul H. Schultz. 1995. "Did Nasdaq Market Makers Implicitly Collude?" Journal of Economic Perspectives 9 (3): 199–208.

DellaVigna, Stefano, and Matthew Gentzkow. 2017. "Uniform Pricing in US Retail Chains." NBER Working Paper no. 23996, Cambridge, MA.

Dubé, Jean-Pierre, and Sanjog Misra. 2017. "Scalable Price Targeting." NBER Working Paper no. 23775, Cambridge, MA.

Eckersley, Peter, and Yomna Nassar. 2017. "Measuring the Progress of AI Research." Electronic Frontier Foundation. https://eff.org/ai/metrics.

Etzioni, Oren, Rattapoom Tuchinda, Craig Knoblock, and Alexander Yates. 2003. "To Buy or Not to Buy: Mining Airfare Data to Minimize Ticket Purchase Price." *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. www.doi.org/10.1145/956750.956767.

Ezrachi, A., and M. E. Stucke. 2017. "Algorithmic Collusion: Problems and Counter-Measures—Note." OECD Roundtable on Algorithms and Collusion. https://www.oecd.org/officialdocuments/publicdisplaydocumentpdf/?cote=DAF/COMP/WD%282017%2925&docLanguage=En.

Fudenberg, Drew. 1992. "Explaining Cooperation and Commitment in Repeated Games." In *Advances in Economic Theory: Sixth World Congress*, Econometric Society Monographs, edited by Jean-Jacques Laffont. Cambridge, MA: Cambridge University Press.

Goodfellow, Ian, Nicolas Papernot, Sandy Huang, Yan Duan, Pieter Abbeel, and Jack Clark. 2017. "Attacking Machine Learning with Adversarial Examples." OpenAI blog, Feb. 26. https://blog.openai.com/adversarial-example-research/.

Hutson, Matthew. 2017. "Will Make AI Smarter for Cash." *Bloomberg Business Week*, Sept. 11.

Kurakin, Alexy, Ian Goodfellow, and Samy Bengio. 2016. "Adversarial Examples in the Physical World." Cornell University Library, ArXiv 1607.02533. https://arxiv.org/abs/1607.02533.

Organisation for Economic Co-operation and Development (OECD). 2017. "Algorithms and Collusion: Competition Policy in the Digital Age." www.oecd.org/competition/algorithms-collusion-competition-policy-in-the-digital-age.htm.

Pierce, David G. 1992 "Repeated Games: Cooperation and Rationality." In *Advances in Economic Theory: Sixth World Congress*, Econometric Society Monographs, edited by Jean-Jacques Laffont. Cambridge, MA: Cambridge University Press.

Rubinstein, Arial. 1986. "Finite Automata Play the Repeated Prisoner's Dilemma." *Journal of Economic Theory* 39:83–96.

Segal, Ilya. 2003. "Optimal Pricing Mechanisms with Unknown Demand." *American Economic Review* 93 (3): 509–29.

Shiller, Benjamin Reed. 2013. "First Degree Price Discrimination Using Big Data." Working Paper no. 58, Department of Economics and International Business School, Brandeis University.

Spiegel, Yossi, and Igal Hendel. 2014. "Small Steps for Workers, A Giant Leap for Productivity." *American Economic Journal: Applied Economics* 6 (1): 73–90.

Sreevallabh, Chivukula, and Wei Liu. 2017. "Adversarial Learning Games with Deep Learning Models." International Joint Conference in Neural Networks. www.doi.org/10.1109/IJCNN.2017.7966196.

Stiglitz, Joseph E., and Bruce C. Greenwald. 2014. *Creating a Learning Society*. New York: Columbia University Press.

Varian, H. 2000. "Managing Online Security Risks." *New York Times*, June 1.

## Comment    Judith Chevalier

Varian provides an excellent overview of industrial organization issues arising out of the adoption of machine learning and artificial intelligence. A number of these issues have potential competition policy implications. For example, exploitation of AI technologies may either increase or decrease economies of scale, leading potentially to situations of market power. Ownership of data, if crucial to competition in a specific industry, may create barriers to entry. The potential for algorithmic collusion clearly leads to antitrust enforcement concerns. Here, I briefly address one of these issues, data ownership, and highlight some potential antitrust policy responses. While I focus here on data ownership as a barrier to entry, some of the policy trade-offs I discuss are germane to the other potential market structure changes highlighted in Varian.

Artificial intelligence and machine-learning processes often use raw data as an input. As Varian points out, it is not at all clear that data defies our usual expectation that a scarce asset or resource will eventually face decreasing returns to scale. Nonetheless, one can certainly imagine circumstances where exclusive ownership of a body of data will create a nearly insurmountable advantage to a market incumbent. While the concern that access to a

Judith Chevalier is the William S. Beinecke Professor of Finance and Economics at the Yale School of Management and a research associate of the National Bureau of Economic Research.