

44 Everyday Chaos

Models We Can Understand

We have not always insisted on understanding our predictions. For example, some of the Founding Fathers of the United States made daily records of the weather and the factors they thought were related to it: when plants start blooming, the first frost, and so forth. They hoped this aggregated data would reveal reliable correlations, such as the daffodils' blooming early signifying that there's a good chance it will be a wet summer. Until the early 1900s, that sort of weather forecasting worked better than not predicting at all.

As Nate Silver explains in *The Signal and the Noise*, this is statistical forecasting: we gather data and use it to make an informed guess about what will happen, based on the assumption that the data is expressing a regularity.¹⁰ Silver says that is how hurricanes were predicted until about thirty years ago. It works pretty well, at least as long as the natural system is fairly consistent.

Statistical forecasting doesn't need a model of the sort proposed in 1900 by Vilhelm Bjerknes, which we looked at in chapter 1. Bjerknes's model explained the dynamics of global weather using seven variables and Newtonian physics: relevant factors connected by rules governing their interactions.¹¹ But there was a problem: even using only seven variables, the computations were so complex that in 1922 a mathematician named Lewis Fry Richardson spent six full weeks doing the work required to predict the weather on a day years earlier, based on data gathered on the days before it. He wasn't even close. After all that grueling work, Richardson's calculated air pressure was 150 times too high.¹²

These days we track hundreds of variables to forecast the weather, as well as to predict the longer-term changes in our climate. We do so with computers that chortle at the 1940s computer—the ENIAC (Electronic Numerical Integrator and Computer)—that took twenty-four hours to predict the next day's weather.¹³ Nevertheless, until machine learning, we relied on the model-based technique that harks back to Pierre-Simon Laplace's demon: if we know the rules governing the behavior of the seven factors that determine the weather,

and if we have the data about them for any one moment in the life of the Earth, we should be able to predict what the next moment's weather will be.

The problem is that so very many factors can affect the weather. In fact, Silver says "the entire discipline of chaos theory developed out of what were essentially frustrated attempts to make weather forecasts."¹⁴ Literally everything on the surface of the planet affects the weather to one degree or another. It is not a coincidence that the example forever associated with Chaos Theory involves a butterfly that creates a catastrophic weather event thousands of miles away.

So if we want to make a prediction about a system like the weather—a third level of predictive complexity, in the terms discussed in the previous chapter—we seem to be left with bad choices. We can rely on statistics and hope that we've been gathering the relevant ones, and that the future will repeat the patterns of the past as surely as the Nile overflows after the Dog Star returns. Or we can figure out the laws governing change and hope that the system is as simple as the model we're using . . . and that it is not disrupted by, say, the Krakatoa volcano that erupted in 1883, spewing forth enough ash to cool the seasons for a year and to chill the oceans for a full century afterward.¹⁵

Bjerknes's seven-factor weather model has the advantage of providing a working model that at least crudely reflects its conceptual model. But we don't always insist on that. The following four examples show different ways successful working models may or may not coincide with our conceptual models. They'll also let us see how deeply machine learning models break with our traditional practices and age-old assumptions about how things happen . . . and our assumptions about how suited we humans are to understanding what happens.

Spreadsheets

Although computerized spreadsheets date back to the early 1960s,¹⁶ they only became widely used after 1978 when Dan Bricklin, a student working on his MBA at Harvard Business School, was annoyed

46 Everyday Chaos

by a class assignment that required calculating the financial implications of a merger. Having to recalculate all the dependent numbers when any single variable changed was more than irksome.¹⁷

So, in the spring of 1978, Bricklin prototyped a spreadsheet on an Apple II personal computer, using a gaming controller in place of a mouse.¹⁸ With the rise of PCs and with the decision by Bricklin and his partner, Bob Frankston, not to patent the software,¹⁹ spreadsheets became a crucial way businesses understood themselves and made decisions: a company's conceptual model of itself now could be expressed in a working model that let the business see the effects of the forces affecting it and of decisions the company was contemplating.

In a remarkably prescient article in 1984, Stephen Levy wrote, "It is not far-fetched to imagine that the introduction of the electronic spreadsheet will have an effect like that brought about by the development during the Renaissance of double-entry bookkeeping."²⁰ He was right. "The spreadsheet is a tool, and it is also a world view—reality by the numbers," Levy wrote.

A spreadsheet is what a business looks like to a traditional computer: quantitative information connected by rules. The rules—formulas—and some of the data, such as fixed costs, are relatively stable. But some of the data changes frequently or even constantly: sales, expenses, headcount, and so on. Personal computers running spreadsheets made keeping the working model up to date so easy and fast that a new decision-making process was made feasible: a spreadsheet is a temptation to tiddle, to try out new futures by plugging in different numbers or by tweaking a relationship. This makes them very different from most traditional models, which focus on representing unchanging relationships, whether they're Newtonian laws or the effect that raising taxes has on savings. Spreadsheets are models that encourage play: you "run the numbers," but then you poke at them to try out "what if this" or "what if that." This was a model meant to be played with.

A spreadsheet thus is a simple example of a working model based on a fully understandable conceptual model. It lets you plug in data or play with the rules to see what the future might or could look like.

Of course, they are inexact, they can't capture all of the relationships among all of the pieces, and the predictions made from their models may be thrown off by events that no one predicted. Because spreadsheets are tools and not perfect encapsulations of every possible eventuality, we accept some distance between the working model and the conceptual model, and between the conceptual model and the real world. We continue to use them because, as George E. P. Box said, "[a]ll models are wrong but some are useful."²¹

Armillary

In the Galileo Museum in Florence sits a beautiful set of nested geared rings, 6.5 feet tall.²² If we today had to guess the point of this intricate mechanism just by looking at it, we might suppose that it's some type of clock. If we were contemporaries of it, we'd be far more likely to recognize that it shows the positions of the major heavenly bodies in Earth's skies for any night.

Antonio Santucci finished this object, called an armillary, in 1593, after five years of work. Although forty-six years earlier Nicolaus Copernicus had shown that the Earth revolves around the sun, Santucci still put the Earth at the center, circled by seven rings that display the positions of the seven known planets. An eighth ring has the fixed stars on it, as well as the markings of the zodiac. Adjust the rings on this wood-and-metal machine, and the planets and fixed stars will align relative to one another and to the Earth. Now gild it and paint in the four winds, the shield of your patron's Medici-related family, and an image of God Himself, and you have a beautiful room-size model of the universe.²³

According to no less an authority than Galileo, even Santucci eventually came around to Copernicus's idea.²⁴ But the armillary's model of the universe is odd beyond its Earth-centric view. It simulates the movement of the heavenly bodies using only circles as components of the mechanism because, from the early Greeks on, it was commonly assumed that because the heavens were the realm of perfection, and circles were the perfect shape, the heavenly bodies must move in perfect circles. That makes the planets a problem, for they

48 Everyday Chaos

1 wander through Earth's sky in distinctly noncircular ways; *planet*
2 comes from the Greek word for *wanderer*. Therefore, if the armillary
3 were to be truthful to its conceptual model, not only did it have to
4 get the planets in the right places relative to Earth, it also had to
5 do it the way the universe does: by using circles. So Santucci set
6 smaller gears turning as they revolved around larger gears that
7 were themselves turning, adding in as many as necessary to model
8 the paths of the planets accurately.²⁵

9 The result is a successful working model that uses a convoluted
10 mechanism dictated by a conceptual model that has been shown to
11 be wildly wrong.

12 The error in its conceptual model also happens to make the work-
13 ing model quite beautiful.

Tides

14
15
16
17 “Unlike the human brain, this one cannot make a mistake.”²⁶

18 That's how a 1914 article in *Scientific American* described a tide-
19 predicting machine made of brass and wood that made mistakes all
20 the time. And its creators knew it.

21 Newton had shown that the gravitational pull of the sun and
22 moon accounted for the rise and fall of the tides around Earth. But
23 his formulas only worked approximately, for, as the *Scientific Amer-*
24 *ican* article pointed out,

25
26 *the earth is not a perfect sphere, it isn't covered with water*
27 *to a uniform depth, it has many continents and islands and*
28 *sea passages of peculiar shapes and depths, the earth does*
29 *not travel about the Sun in a circular path, and Earth, Sun*
30 *and Moon are not always in line. The result is that two tides*
31 *are rarely the same for the same place twice running, and*
32 *that tides differ from each other enormously in both times*
33 *and in amplitude.*²⁷

34
35 In his book *Tides: The Science and Spirit of the Ocean*, Jonathan
36 White notes, “There are hundreds of these eccentricities, each call-

ing out to the oceans—some loudly, some faintly, some repeating every four hours and others every twenty thousand years.” Newton knew he was ignoring these complications, but they were too complicated to account for. (It’s quite possible that he never saw an ocean himself.)²⁸

It was Laplace who again got Newton righter than Newton did, creating formulas that included the moon’s eight-year cycle of distances from the Earth, its varying distance north and south of the equator, the effect of the shape and depth of the ocean’s basin, the texture of the ocean floor, the water’s fluctuating temperatures, and other conditions.²⁹

This added nuance to Newton’s model, but a vast number of additional factors also affect the tides. It took about another hundred years for Lord Kelvin, in 1867, to come up with a way of predicting tides that takes all the factors into account without having to know what all of them are.³⁰

As the 1914 *Scientific American* article explains it, imagine a pencil floating up and down in an ocean, creating a curve as it draws on a piece of paper scrolling past it. Imagine lots of pencils placed at uniform distances from one another. Now imagine the ocean lying still, without any bodies exerting gravitational forces on it. Finally, imagine a series of fictitious suns and moons above Earth in exactly the right spots for their gravity to pull that pencil to create exactly those curves. Wherever you have a curve that needs explaining, add another imaginary sun or moon in the right position to get the expected result. Lord Kelvin ended up with a “very respectable number” of imaginary suns and moons circling the Earth, as the article puts it. If adding sea serpents would have helped, presumably Lord Kelvin would have added them as well.³¹

With the assistance of George Darwin—brother of Charles—Lord Kelvin computed formulas that expressed the pull of these imaginary bodies, then designed a machine that used chains and pulleys to add up all of those forces and to draw the tidal curves. By 1914, this had evolved into the beast feted in the *Scientific American* article: fifteen thousand parts that, combined, could draw a line showing the tides at any hour.

50 Everyday Chaos

1 Lord Kelvin was in fact not the first to imagine a science-fiction
2 Earth circled by multiple suns and moons that create the wrinkled
3 swells and ebbs of tides caused by the vagaries of the Earth's geog-
4 raphy, topology, weather, and hundreds of other factors. Laplace
5 himself "imagined a stationary Earth with these tide components
6 circling as satellites."³² Lord Kelvin's machine and its iterations took
7 this to further levels of detail, while accepting that the actual tides
8 are subject to still more factors that simply could not be captured
9 in the machine's model—the influx of melted snow from a particu-
10 larly long winter, the effect of storms, and all the other influences
11 Earth is heir to. The *Scientific American* article could claim the ma-
12 chine never makes a mistake because Kelvin's machine was as
13 accurate as the tools and data of the time allowed, so it became the
14 accuracy we counted as acceptable . . . all while relying on a ficti-
15 tious model.

16 It set this level of accuracy by building a working model that is
17 knowingly, even wildly, divorced from its conceptual model.

The River

20 In 1943, the US Army Corps of Engineers set Italian and German
21 prisoners of war to work building the largest scale model in history:
22 two hundred acres representing the 41 percent of the United States
23 that drains into the Mississippi River. By 1949 the model was being
24 used to run simulations to determine what would happen to cities
25 and towns along the way if water flooded in. It's credited with pre-
26 venting \$65 million in damage from a flood in Omaha in 1952.³³ In
27 fact, some claim its simulations are more accurate than the exist-
28 ing digital models.³⁴

30 Water was at the heart of another type of physical model: the MO-
31 NIAC (Monetary National Income Analogue Computer) economic
32 simulator built in 1949 by the New Zealand economist Alban William
33 Housego Phillips.³⁵ The MONIAC used colored water in transparent
34 pipes to simulate the effects of Keynesian economic policies. Tanks
35 of water represented "households, business, government, exporting
36

and importing sectors of the economy,” measuring income, spending, and GDP.³⁶

It worked, given its limitations. The number of variables it could include was constrained by the number of valves, tubes, and tanks that could fit in a device about the size of a refrigerator.³⁷ But because it only took account of a relative handful of the variables that influence the state of a national economy, it was far less accurate than the Mississippi River simulator. Yet the flow of water through a river the size of the Mississippi is also affected by more variables than humans can list. So how could the Mississippi model get predictions so right?

The Mississippi had the advantage of not requiring its creators to have a complete conceptual model of how a river works. For example, if you want to predict what will happen if you place a boulder in a rapids, you don’t have to have a complete model of fluid dynamics; you can just build a working scale model that puts a small rock into a small flow. So long as scale doesn’t matter, your model will give you your answer. As Stanford Gibson, a senior hydraulic engineer in the Army Corps of Engineers, said about the Mississippi basin project, “The physical model will simulate the processes on its own.”³⁸

So this working model can deal with more complexity *because* it doesn’t have a conceptual model: it puts the actual forces to use in a controlled and adjustable way. Because the model is not merely a symbolic one—real water is rolling past a real, scaled-down boulder—the results aren’t limited by what we know to factor in. That’s the problem with the MONIAC: it sticks with factors that we know about. It’s like reducing weather to seven known factors.

Still, the Mississippi River basin model may seem to make no assumptions about what affects floods, but of course it does. It assumes that what happens at full scale also happens at 1/2000 scale, which is not completely accurate for the dynamics of water; for example, the creators of a model of San Francisco Bay purposefully distorted the horizontal and vertical scales by a factor of ten in order to get the right flow over the tidal flats.³⁹ Likewise, the Mississippi model does

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

52 Everyday Chaos

1 not simulate the gravitational pull of the sun and the moon. Nor
2 does it grow miniature crops in the fields. The model assumes those
3 factors are not relevant to the predictions it was designed to en-
4 able. Using the Mississippi model to simulate the effects of climate
5 change or the effect of paddle wheelers on algae growth probably
6 wouldn't give reliable results, for those phenomena are affected by
7 factors not in the model and are sensitive to scale.

8 The Mississippi model wasn't constructed based on an explicit
9 conceptual model of the Mississippi River basin, and it works without
10 yielding one. Indeed, it works because it doesn't require us to under-
11 stand the Mississippi River: it lets the physics of the simulation do its
12 job without imposing the limitations of human reason on it. The re-
13 sult is a model that is more accurate than one like the MONIAC that
14 was constructed based on human theory and understanding. So the
15 advent of machine learning is not the first time we have been pre-
16 sented with working models for which we have no conceptual model.

17 But, as we'll see, machine learning is making clear a problem with
18 the very idea of conceptual models. Suppose our concepts and the
19 world they model aren't nearly as alike as we've thought? After
20 all, when it comes to the Mighty Mississippi, the most accurate
21 working model lets physical water flow deeper than our conceptual
22 understanding.

23
24
25 Despite the important differences among all these models—from
26 spreadsheets to the Mississippi—it's the similarities that tell us
27 the most about how we have made our way in a wildly unpredict-
28 able world.

29 In all these cases, models *stand in* for the real thing: the armil-
30 lary is not the heavenly domain, the spreadsheet is not the business,
31 the tubes filled with colored water are not the economy. They do so
32 by *simplifying* the real-world version. A complete tidal model would
33 have to include a complete weather model, which would have to in-
34 clude a complete model of industrial effects on the climate, until the
35 entire world and heavens have been included. Models simplify sys-
36 tems until they yield acceptably accurate predictions.

Models thereby assume that we humans can identify the elements that are relevant to the thing we are modeling: the factors, rules, and principles that determine how it behaves. Even the model of the Mississippi, which does not need to understand the physics of fluid dynamics, assumes that floods are affected by the curves and depths of the river and not by whether the blue vervain growing along the sides of the river are in flower. This also implies that models assume some degree of *regularity*. The armillary assumes that the heavenly bodies will continue to move across the skies in their accustomed paths; the tidal machine assumes the gravitational mass of the sun and moon will remain constant; the spreadsheet assumes that sales are always going to be added to revenues.

Because the simplification process is *done by human beings*, models reflect our strengths and our weaknesses. The strengths include our ability to see the order beneath the apparent flux of change. But we are also inevitably prone to using unexamined assumptions, have limited memories and inherent biases, and are willing to simplify our world to the point where we can understand it.

Despite models' inescapable weaknesses due to our own flawed natures, they have been essential to how we understand and control our world. They have become the *stable frameworks that enable us to predict and explain the ever-changing and overwhelming world in process all around us*.

Beyond Explanation

We are transitioning to a new type of working model, one that does not require knowing how a system works and that does not require simplifying it, at least not to the degree we have in the past. This makes the rise of machine learning one of the most significant disruptions in our history.⁴⁰

In the introduction, we talked about Deep Patient, a machine learning system that researchers at Mount Sinai Hospital in New York fed hundreds of pieces of medical data about seven hundred thousand patients. As a result, it was able to predict the onset of

54 Everyday Chaos

1 diseases that have defied human diagnostic abilities. Likewise, a
2 Google research project analyzed the hospital health records of
3 216,221 adults. From the forty-six billion data points, it was able to
4 predict the length of a patient's stay in the hospital, the probability
5 that the patient would exit alive, and more.⁴¹

6 These systems work: they produce probabilistically accurate out-
7 comes. But why?

8 Both of these examples use *deep learning*, a type of machine learn-
9 ing that looks for relationships among the data points without
10 being instructed what to look for. The system connects the nodes
11 into a web of probabilistic dependencies, and then uses that web—
12 an “artificial neural network”—to refine the relationships again
13 and again. The result is a network of data nodes, each with a “weight”
14 that is used to determine whether the nodes it is connected to will
15 activate; in this way, artificial neural networks are like the brain's
16 very real neural network.

17 These networks can be insanely complicated. For example, Deep
18 Patient looked at five hundred factors for each of the hundreds of
19 thousands patients whose records it analyzed, creating a final data
20 set of two hundred million pieces of data. To check on a particular
21 patient's health, you run her data through that network and get back
22 probabilistic predictions about the medical risks she faces. For ex-
23 ample, Deep Patient is unusually good at telling which patients are
24 at risk of developing schizophrenia, a condition that is extremely
25 hard for human doctors to predict.⁴²

26 But the clues the system uses to make these predictions are not
27 necessarily like the signs doctors typically use, the way tingling and
28 numbness can be an early sign of multiple sclerosis, or sudden thirst
29 sometimes indicates diabetes. In fact, if you asked Deep Patient how
30 it came to classify people as likely to develop schizophrenia, there
31 could be so many variables arranged in such a complex constella-
32 tion that we humans might not be able see the patterns in the data
33 even if they were pointed out to us. Some factor might increase the
34 probability of a patient becoming schizophrenic but only in con-
35 junction with other factors, and the set of relevant factors may it-
36 self vary widely, just as your spouse dressing more formally might

mean nothing alone but, in conjunction with one set of “tells,” might be a sign that she is feeling more confident about herself and, with other sets, might mean that she is aiming for a promotion at work or is cheating on you. The number and complexity of contextual variables mean that Deep Patient simply cannot always explain its diagnoses as a conceptual model that its human keepers can understand.

Getting explanations from a machine learning system is much easier when humans have programmed in the features the system should be looking for. For example, the Irvine, California–based company Bitvore analyzes news feeds and public filings to provide real-time notifications to clients about developments relevant to them. To do this, its dozens of algorithms have been trained to look for over three hundred different sorts of events, including CEO resignations, bankruptcies, lawsuits, and criminal behavior, all of which might have financial impacts. Jeff Curie, Bitvore’s president, says that it’s like having several hundred subject experts each scouring a vast stream of data.⁴³ When one of these robotic experts finds something relevant to its area of expertise, it flags it, tags it, and passes it on to the rest, who add what they know and connect it to other events. This provides clients—including intelligence agencies and financial houses—not just an early warning system that sounds an alarm but also contextualized information about the alarm.

Bitvore’s system is designed so that its conclusions will always be explicable to clients. The company’s chief technology officer, Greg Bolcer, told me about a time when the system flagged news about cash reserves as relevant to its municipal government clients. It seemed off, so Bolcer investigated. The system reported that the event concerned not cash reserves but a vineyard’s “special reserve” wines and was of no relevance to Bitvore’s clients. To avoid that sort of machine-based confusion, Bitvore’s system is architected so that humans can always demand an explanation.⁴⁴

Bitvore is far from the only system that keeps its results explicable. Andrew Jennings, the senior vice president of scores and analytics at FICO, the credit-scoring company, told me, “There are a number of long standing rules and regulations around credit scoring in the

56 Everyday Chaos

1 US and elsewhere as a result of legislation that require[s] people
2 who build credit scores to manage the tradeoff between things that
3 are predictively useful and legally permitted.”⁴⁵ Machine learning
4 algorithms might discover—to use a made-up example—that the
5 Amish generally are good credit risks but, say, Episcopalians are
6 not. Even if this example were true, that knowledge could not be
7 used in computing a credit score because US law prevents discrimi-
8 nation on the basis of religion or other protected classes. Credit
9 score companies are also prohibited from using data that is a sur-
10rogate for these attributes, such as an applicant’s subscribing to
11 *Amish Week* magazine or, possibly, the size of someone’s monthly
12 electricity bills.

13 There are additional constraints on the model that credit score
14 companies can use to calculate credit risk. If a lender declines a loan
15 application, the lender has to provide the reasons why the applicant’s
16 score was not higher. Those reasons have to be addressable by the
17 consumer. For example, Jennings explained, an applicant might be
18 told, “Your score was low because you’ve been late paying off your
19 credit cards eight times in the past year,” a factor that the applicant
20 can improve in the future.

21 But suppose FICO’s manually created models turn out to be less
22 predictive of credit risk than a machine learning system would be?
23 Jennings says that they have tested this and found the differences
24 between the manual and machine learning models to be insignifi-
25 cant. But the promise of machine learning is that there are times
26 when the machine’s inscrutable models may be more accurately pre-
27 dictive than manually constructed, human-intelligible ones.

28 As such systems become more common, the demand for keeping
29 their results understandable is growing. It’s easy to imagine a pa-
30tient wanting to know why some future version of Deep Patient has
31 recommended that she stop eating high-fat foods, or that she preemptively
32 get a hysterectomy. Or a job applicant might want to know
33 whether her race had anything to do with her being ruled out of the
34 pool of people to interview. Or a property owner might want to know
35 why a network of autonomous automobiles sent one of its cars
36 through her fence as part of what that network thought was the op-

timal response to a power line falling onto a highway. Sometimes these systems will be able to report on what factors weighed the heaviest in a decision, but sometimes the answer will consist of the weightings of thousands of factors, with no one factor being dominant. These systems are likely to become more inexplicable as the models become more complex and as the models incorporate outputs from other machine learning systems.

But it's controversial. As it stands, in most fields developers generally implement these systems aiming at predictive accuracy, free of the requirement to keep them explicable. While there is a strong contingent of computer scientists who think that we will always be able to wring explanations out of machine learning systems, what counts as an explanation, and what counts as understanding, is itself debatable.⁴⁶ For example, the *counterfactual* approach proposed by Sandra Wachter, Brent Mittelstadt, and Chris Russell at Oxford could discover whether, say, race was involved in why someone was put into the "do not insure" bin by a machine learning application: in the simplest case, resubmit the same application with only the race changed, and if the outcome changes, then you've shown race affected the outcome.⁴⁷ It does not at all take away from the usefulness of the counterfactual approach to point out that it produces a very focused and minimal sense of "explanation," and even less so of "understanding."

In any case, in many instances, we'll accept the suggestions of these systems if their performance records are good, just as we'll accept our physician's advice if she can back it up with a study we can't understand that shows that a treatment is effective in a high percentage of cases—and just as many of us already accept navigation advice from the machine learning-based apps on our phones without knowing how those apps come up with their routes. The riskier or more inconvenient the medical treatment, the higher the probability of success we'll demand, but the justification will be roughly the same: a good percentage of people who follow this advice do well. That's why we took aspirin—initially in the form of willow bark—for thousands of years before we understood why it works.

58 Everyday Chaos

1 As machine learning surpasses the predictive accuracy of old-
2 style models, and especially as we butt our heads against the wall
3 of inexplicability, we are coming to accept a new model of models,
4 one that reflects a new sense of how things happen.

Four New Ways of Happening

5
6
7
8
9 Suppose someday in the near future your physician tells you to cut
10 down on your potassium intake; no more banana smoothies for you.
11 When you ask why, she replies that Deep Asclepius—a deep learn-
12 ing system I’ve made up—says you fit the profile of people who are
13 40 percent more likely to develop Parkinson’s disease at some point
14 in their lives if they take in too much potassium (which I’m also mak-
15 ing up).

16 “What’s that profile?” you may ask.

17 Your physician explains: “Deep Asclepius looks at over one thou-
18 sand pieces of data for each person, and Parkinson’s is a complex
19 disease. We just don’t know why those variables combine to suggest
20 that you are at risk.”

21 Perhaps you’ll accept your physician’s advice without asking
22 about her reasons, just as you tend to accept it when your physician
23 cites studies you’re never going to look up and couldn’t understand
24 if you did. In fact, Deep Asclepius’s marketers will probably forestall
25 the previous conversation by turning the inexplicability of its results
26 into a positive point: “Medical treatment that’s as unique as you
27 are . . . and just as surprising!”

28 Casual interactions such as these will challenge the basic as-
29 sumptions of our past few thousand years of creating models.

30 First, we used to assume that we humans made the models: in
31 many cases (but not all, as we’ve seen) we came up with the simpli-
32 fied conceptual model first, and then we made a working model. But
33 *deep learning’s models are not created by humans*, at least not di-
34 rectly.⁴⁸ Humans choose the data and feed it in, humans head the
35 system toward a goal, and humans can intercede to tune the weights
36 and the outcomes. But humans do not necessarily tell the machine

what features to look for. For example, Google fed photos that included dumbbells into a machine learning system to see if it could pick out the dumbbells from everything else in the scene. The researchers didn't give the system any characteristics of dumbbells to look for, such as two disks connected by a rod. Yet without being told, the system correctly abstracted an image of two disks connected by a bar. On the other hand, the image also included a muscular arm holding the dumbbell, reflecting the content of the photos in the training set.⁴⁹ (We'll talk in the final chapter about whether that was actually a mistake.)

Because the models deep learning may come up with are not based on the models we have constructed for ourselves, they can be opaque to us. This does not mean, however, that deep learning systems escape human biases. As has become well known, they can reflect and even amplify the biases in the data itself. If women are not getting hired for jobs in tech, a deep learning system trained on existing data is likely to "learn" that women are not good at tech. If black men in America are receiving stiffer jail sentences than white men in similar circumstances, the training based on that data is very likely to perpetuate that bias.⁵⁰

This is not a small problem easily solved. Crucially, it is now the subject of much attention, research, and development.

The second assumption about models now being challenged comes from the fact that our conceptual models cover more than one case; that's what makes them models. We have therefore tended to construct them out of general principles or rules: Newton's laws determine the paths of comets, lowering prices tends to increase sales, and all heavenly bodies move in circles, at least according to the ancient Greeks. Principles find simpler regularities that explain more complex particulars. But *deep learning models are not generated premised on simplified principles, and there's no reason to think they are always going to produce them*, just as A/B testing may not come up with any generalizable rules for how to make ads effective.

Sometimes a principle or at least a rule of thumb does emerge from a deep learning system. For example, in a famous go match between Lee Sedol, a world-class master, and Google's AlphaGo, the computer

60 Everyday Chaos

1 initially played aggressively. But once AlphaGo had taken over the
2 left side of the board, it started to play far more cautiously. This
3 turned out to be part of a pattern: when AlphaGo is 70 percent con-
4 fident it's going to win, it plays less aggressively. Perhaps this is a
5 generalizable heuristic for human go players as well.⁵¹ Indeed, in
6 2017, Google launched a program that brings together human play-
7 ers and AlphaGo so that the humans can learn from the machine.⁵²

8 A later version of AlphaGo took the next step. Rather than train-
9 ing AlphaGo on human games of go, the programmers fed in nothing
10 but the rules of the game and then had the machine play itself. After
11 just three days, the system so mastered the game that it was able to
12 beat the prior version of AlphaGo a hundred games out of a hun-
13 dred.⁵³ When experts studied the machine-vs.-machine games that
14 Google published, some referred to the style of play as "alien."⁵⁴

15 Isn't that the literal truth?

16 If so, it's because of the third difference: *deep learning systems do*
17 *not have to simplify the world to what humans can understand.*

18 When we humans build models for ourselves, we like to find the
19 general principles that govern the domain we're modeling. Then we
20 can plug in the specifics of some instance and read out the date and
21 time of an eclipse or whether the patient has type 2 diabetes. Deep
22 learning systems typically put their data through artificial neural
23 networks to identify the factors (or "dimensions") that matter and
24 to discern their interrelationships. They typically do this several
25 times, sometimes making the relationships among the pieces under-
26 standable only by understanding the prior pass, which may have
27 surpassed our understanding on its own.

28 The same holds for the data we input in order to get, say, a diag-
29 nosis from my hypothetical Deep Asclepius system. Deep Asclepius
30 doesn't have to confine itself to the handful of factors a patient is
31 typically asked to list on a three-page form while sitting in the wait-
32 ing room. It can run the patient's lifetime medical record against
33 its model, eventually even pulling in, perhaps, environmental data,
34 travel history, and education records, noting relationships that
35 might otherwise have been missed (and assuming privacy has been
36

waived). Simplification is no longer required to create a useful working model.

The success of deep learning suggests to us that the world does not separate into neatly divided events that can be predicted by consulting a relative handful of eternal laws. The comet crossing paths with Jupiter, Saturn, and the sun is not a three-body or four-body problem but rather an all-body problem, for, as Newton well knew, every gravitational mass exerts some pull on every other. Calculating a comet's path by computing the gravitational effect of three massive bodies is a convenient approximation that hides the alien complexity of the truth.

As we gasp at what our machines can now do, we are also gasping at the clear proof of what we have long known but often suppressed: our old, oversimplified models were nothing more than the rough guess of a couple of pounds of brains trying to understand a realm in which everything is connected to, and influenced by, everything.

Fourth, where we used to assume that our conceptual models were stable if not immutable, everything being connected to everything means that *machine learning's model can constantly change*. Because most of our old models were based on stable principles or laws, they were slower to change. The classic paradigm for this was put forward by Thomas Kuhn in his 1962 book *The Structure of Scientific Revolutions*. Historically, Kuhn says, a science's overarching model (which he calls a paradigm) maintains itself as data piles up that doesn't fit very well.⁵⁵ At some point—it's a nonlinear system—a new paradigm emerges that fits the anomalous data, as when germ theory replaced the long-held idea that diseases such as malaria were caused by bad air. But changes in machine learning models can occur simply by retraining them on new data. Indeed, some systems learn continuously. For example, our car navigation systems base our routes on real-time information about traffic and can learn from that data that Route 128 tends to get backed up around four o'clock in the afternoon. This can create a feedback loop as the navigation system directs people away from Route 128 at that time, perhaps

62 Everyday Chaos

1 reducing the backups. These feedback loops let the model constantly
2 adjust itself to changing conditions and optimize itself further.

3 As we'll see, this reveals a weakness in our traditional basic strat-
4 egy for managing what will happen, for the elements of a machine
5 learning model may not have the sort of one-to-one relationship
6 that we envision when we search for the right "lever" to pull. When
7 everything affects everything else, and when some of those rela-
8 tionships are complex and nonlinear—that is, tiny changes can
9 dramatically change the course of events—butterflies can be as
10 important as levers.

11 Overall, these changes mean that while models have been the sta-
12 ble frameworks that enable explanation, *now we often explain some-*
13 *thing by trying to figure out the model our machines have created.*

14 The only real continuity between our old types of models and our
15 new ones is that both are representations of the world. But one is a
16 representation that we have created based on our understanding, a
17 process that works by reducing the complexity of what it encounters.
18 The other is generated by a machine we have created, and into which
19 we have streamed oceans of data about everything we have thought
20 might possibly be worth noticing. The source, content, structure,
21 and scale of these two types of representations are vastly, discon-
22 certingly different.

Explanation Games

23
24
25
26
27 "JAL 123 was twelve minutes into its flight when a bang was heard
28 on the flight deck."

29 On August 12, 1985, thirty-two minutes after that, the pilots lost
30 their struggle to keep the plane aloft as its right wingtip clipped a
31 mountain. The Boeing 747 came down with such force that three
32 thousand trees in its path were destroyed. Of its 509 passengers, 505
33 were killed. It is to this day the plane crash that claimed the most
34 victims.⁵⁶

35 The task facing the investigators who arrived from multiple
36 organizations and countries was made more difficult by the impend-