

## CONTENTS



<i>Preface</i>	xi
----------------	----

### PART 1:

## DARK DATA: THEIR ORIGINS AND CONSEQUENCES

<b>CHAPTER 1: Dark Data: What We Don't See</b>	
<b>Shapes Our World</b>	3
<i>The Ghost of Data</i>	3
<i>So You Think You Have All the Data?</i>	12
<i>Nothing Happened, So We Ignored It</i>	17
<i>The Power of Dark Data</i>	22
<i>All around Us</i>	24
<b>CHAPTER 2: Discovering Dark Data: What We</b>	
<b>Collect and What We Don't</b>	28
<i>Dark Data on All Sides</i>	28
<i>Data Exhaust, Selection, and Self-Selection</i>	31
<i>From the Few to the Many</i>	43
<i>Experimental Data</i>	56
<i>Beware Human Frailties</i>	67
<b>CHAPTER 3: Definitions and Dark Data: What Do</b>	
<b>You Want to Know?</b>	72
<i>Different Definitions and Measuring the Wrong Thing</i>	72

<i>You Can't Measure Everything</i>	80
<i>Screening</i>	90
<i>Selection on the Basis of Past Performance</i>	94
<b>CHAPTER 4: Unintentional Dark Data: Saying One Thing, Doing Another</b>	98
<i>The Big Picture</i>	98
<i>Summarizing</i>	102
<i>Human Error</i>	103
<i>Instrument Limitations</i>	108
<i>Linking Data Sets</i>	111
<b>CHAPTER 5: Strategic Dark Data: Gaming, Feedback, and Information Asymmetry</b>	114
<i>Gaming</i>	114
<i>Feedback</i>	122
<i>Information Asymmetry</i>	128
<i>Adverse Selection and Algorithms</i>	130
<b>CHAPTER 6: Intentional Dark Data: Fraud and Deception</b>	140
<i>Fraud</i>	140
<i>Identity Theft and Internet Fraud</i>	144
<i>Personal Financial Fraud</i>	149
<i>Financial Market Fraud and Insider Trading</i>	153
<i>Insurance Fraud</i>	158
<i>And More</i>	163

<b>CHAPTER 7: Science and Dark Data: The Nature of Discovery</b>	167
<i>The Nature of Science</i>	167
<i>If Only I'd Known That</i>	172
<i>Tripping over Dark Data</i>	181
<i>Dark Data and the Big Picture</i>	184
<i>Hiding the Facts</i>	199
<i>Retraction</i>	215
<i>Provenance and Trustworthiness: Who Told You That?</i>	217

PART II:

ILLUMINATING AND USING DARK DATA

<b>CHAPTER 8: Dealing with Dark Data: Shining a Light</b>	223
<i>Hope!</i>	223
<i>Linking Observed and Missing Data</i>	224
<i>Identifying the Missing Data Mechanism</i>	233
<i>Working with the Data We Have</i>	236
<i>Going Beyond the Data: What If You Die First?</i>	241
<i>Going Beyond the Data: Imputation</i>	245
<i>Iteration</i>	252
<i>Wrong Number!</i>	256
<b>CHAPTER 9: Benefiting from Dark Data: Reframing the Question</b>	262
<i>Hiding Data</i>	262

<i>Hiding Data from Ourselves: Randomized Controlled Trials</i>	263
<i>What Might Have Been</i>	265
<i>Replicated Data</i>	269
<i>Imaginary Data: The Bayesian Prior</i>	276
<i>Privacy and Confidentiality Preservation</i>	278
<i>Collecting Data in the Dark</i>	287
<b>CHAPTER 10: Classifying Dark Data: A Route through the Maze</b>	291
<i>A Taxonomy of Dark Data</i>	291
<i>Illumination</i>	298
<i>Notes</i>	307
<i>Index</i>	319

## Chapter 1

# DARK DATA



## What We Don't See Shapes Our World

### The Ghost of Data

First, a joke.

Walking along the road the other day, I came across an elderly man putting small heaps of powder at intervals of about 50 feet down the center of the road. I asked him what he was doing. “It’s elephant powder,” he said. “They can’t stand it, so it keeps them away.”

“But there are no elephants here,” I said.

“Exactly!” he replied. “It’s wonderfully effective.”

Now, on to something much more serious.

Measles kills nearly a 100,000 people each year. One in 500 people who get the disease die from complications, and others suffer permanent hearing loss or brain damage. Fortunately, it’s rare in the United States; for example, only 99 cases were reported in 1999. But a measles outbreak led Washington to declare a statewide emergency in January 2019, and other states also reported dramatically increased numbers of cases.<sup>1</sup> A similar pattern was reported elsewhere. In Ukraine, an outbreak resulted in over 21,000 cases by mid-February 2019.<sup>2</sup> In Europe there were 25,863 cases in 2017, but in 2018 there were over 82,000.<sup>3</sup> From

1 January 2016 through the end of March 2017, Romania reported more than 4,000 cases and 18 deaths from measles.

Measles is a particularly pernicious disease, spreading undetected because the symptoms do not become apparent until some weeks after you contract it. It slips under the radar, and you have it before you even know that it's around.

But the disease is also preventable. A simple vaccination can immunize you against the risk of contracting measles. And, indeed, national immunization programs of the kind carried out in the United States have been immensely successful—so successful in fact that most parents in countries which carry out such programs have never seen or experienced the terrible consequences of such preventable diseases.

So, when parents are advised to vaccinate their children against a disease they have neither seen nor heard of any of their friends or neighbors having, a disease which the Centers for Disease Control and Prevention announced was no longer endemic in the United States, they naturally take the advice with a pinch of salt.

Vaccinate against something which is not there? It's like using the elephant powder.

Except that, unlike the elephants, the risks are still there, just as real as ever. It's merely that the information and data these parents need to make decisions are missing, so that the risks have become invisible.

My general term for the various kinds of missing data is *dark data*. Dark data are concealed from us, and that very fact means we are at risk of misunderstanding, of drawing incorrect conclusions, and of making poor decisions. In short, our ignorance means we get things wrong.

The term “dark data” arises by analogy with the dark matter of physics. About 27 percent of the universe consists of this

mysterious substance, which doesn't interact with light or other electromagnetic radiation and so can't be seen. Since dark matter can't be seen, astronomers were long unaware of its existence. But then observations of the rotations of galaxies revealed that the more distant stars were not moving more slowly than stars nearer the center, contradicting what we would have expected from our understanding of gravity. This rotational anomaly can be explained by supposing that galaxies have more mass than appears to be the case judging from the stars and other objects we can see through our telescopes. Since we can't see this extra mass, it has been called dark matter. And it can be significant (I almost said "it can matter"): our home galaxy, the Milky Way, is estimated to have some ten times as much dark matter as ordinary matter.

Dark data and dark matter behave in an analogous way: we don't see such data, they have not been recorded, and yet they can have a major effect on our conclusions, decisions, and actions. And as some of the later examples will show, unless we are aware of the possibility that there's something unknown lurking out there, the consequences can be disastrous, even fatal.

The aim of this book is to explore just how and why dark data arise. We shall look at the different kinds of dark data and see what leads to them. We shall see what steps we can take to avoid dark data's arising in the first place. We shall see what we can do when we realize that dark data are obscured from us. Ultimately, we shall also see that if we are clever enough, we can sometimes take advantage of dark data. Curious and paradoxical though that may seem, we can make use of ignorance and the dark data perspective to enable better decisions and take better actions. In practical terms, this means we can lead healthier lives, make more money, and take lower risks by judicious use of the unknown. This doesn't mean we should hide information from others

(though, as we shall also see, deliberately concealed data is one common kind of dark data). It is much more subtle than that, and it means that everyone can benefit.

Dark data arise in many different shapes and forms as well as for many different reasons, and this book introduces a taxonomy of such reasons, the *types* of dark data, labeled *DD-Type x*, for “Dark Data-Type x.” There are 15 *DD-Types* in all. My taxonomy is not exhaustive. Given the wealth of reasons for dark data, that would probably be impossible. Moreover, any particular example of dark data might well illustrate the effect of more than one *DD-Type* simultaneously—*DD-Types* can work together and can even combine in an unfortunate synergy. Nonetheless, an awareness of these *DD-Types*, and examination of examples showing how dark data can manifest, can equip you to identify when problems occur and protect you against their dangers. I list the *DD-Types* at the end of this chapter, ordered roughly according to similarity, and describe them in more detail in chapter 10. Throughout the book I have indicated some of the places when an example of a particular *Type* occurs. However, I have deliberately not tried to do this in an exhaustive way—that would be rather intrusive.

To get us going, let’s take a new example.

In medicine, trauma is serious injury with possible major long-term consequences. It’s one of the most serious causes of “life years lost” through premature death and disability, and is the commonest cause of death for those under age 40. The database of the Trauma Audit and Research Network (TARN) is the largest medical trauma database in Europe. It receives data on trauma events from more than 200 hospitals, including over 93 percent of the hospitals in England and Wales, as well as hospitals in Ireland, the Netherlands, and Switzerland. It’s clearly



a very rich seam of data for studying prognoses and the effectiveness of interventions in trauma cases.

Dr. Evgeny Mirkes and his colleagues from the University of Leicester in the UK looked at some of the data from this database.<sup>4</sup> Among the 165,559 trauma cases they examined, they found 19,289 with unknown outcomes. “Outcome” in trauma research means whether or not the patient survives at least 30 days after the injury. So the 30-day survival was unknown for over 11 percent of the patients. This example illustrates a common form of dark data—our *DD-Type 1: Data We Know Are Missing*. We know these patients had some outcome—we just don’t know what it was.

No problem, you might think—let’s just analyze the 146,270 patients for whom we do know the outcome and base our understanding and prognoses on those. After all, 146,270 is a big number—within the realm of medicine it’s “big data”—so surely we can be confident that any conclusions based on these data will be right.

But can we? Perhaps the missing 19,289 cases are very different from the others. After all, they were certainly different in that they had unknown outcomes, so it wouldn’t be unreasonable to suspect they might differ in other ways. Consequently, any analysis of the 146,270 patients with known outcomes might be misleading relative to the overall population of trauma patients. Thus, actions taken on the basis of such analysis might be the wrong actions, perhaps leading to mistaken prognoses, incorrect prescriptions, and inappropriate treatment regimes, with unfortunate, even fatal, consequences for patients.

To take a deliberately unrealistic and extreme illustration, suppose that all 146,270 of those with known outcomes survived and recovered without treatment, but the 19,289 with unknown

outcomes all died within two days of admission. If we ignored those with unknown outcomes, we would justifiably conclude there was nothing to worry about, and all patients with trauma recovered. On this basis, we wouldn't treat any incoming trauma cases, expecting them to recover naturally. And then we would be horrified and confused by the fact that more than 11 percent of our patients were dying.

Before I go any further with this story, I want to reassure the reader. My extreme illustration is very much a worst-case scenario—we might reasonably expect things not to be that bad in reality—and Dr. Mirkes and his colleagues are experts on missing data analysis. They are very aware of the dangers and have been developing statistical methods to cope with the problem; I describe similar methods later in this book. But the take-home message from this story is that *things may not be what they seem*. Indeed, if there were a single take-home message from this book, that would be a good approximation to it: while it helps to have lots of data—that is, “big data”—size is not everything. And what you don't know, the data you don't have, may be even more important in understanding what's going on than the data you do have. In any case, as we shall see, the problems of dark data are not merely big-data problems: they also arise with small data sets. They are ubiquitous.

My story about the TARN database may be exaggerated, but it serves as a warning. Perhaps the outcomes of the 19,289 patients were not recorded precisely *because* they'd all died within 30 days. After all, if the outcome was based on contacting the patients 30 days after admission to see how they were, none of those who died would respond to questions. Unless we were aware of this possibility, we'd never record that any patients had died.

This may sound a bit silly, but in fact such situations arise quite often. For example, a model built to determine the prognosis for

patients being given a particular treatment will be based on the outcomes of previous patients who received that treatment. But what if insufficient time had passed for all such previous patients to have reached an outcome? For those patients the eventual outcome would be unknown. A model built just on those for whom the outcome was known could be misleading.

A similar phenomenon happens with surveys, in which *non-response* is a source of difficulty. Researchers will typically have a complete list of people from whom they would ideally like answers, but, also typically, not everyone responds. If those who do respond differ in some way from those who do not, then the researchers might have cause to doubt whether the statistics are good summaries of the population. After all, if a magazine carried out a survey of its subscribers asking the single question, Do you reply to magazine surveys? then we could not interpret the fact that 100 percent of those who replied answered yes as meaning that all the subscribers replied to such surveys.

The preceding examples illustrate our first type of dark data. We know that the data for the TARN patients all exist, even if the values aren't all recorded. We know that the people on the survey list had answers, even if they did not give them. In general, we know that there are values for the data; we just don't know what those values are.

An illustration of a different kind of dark data (*DD-Type 2: Data We Don't Know Are Missing*) is the following.

Many cities have problems with potholes in road surfaces. Water gets into small cracks and freezes in the winter, expanding the cracks, which are then further damaged by car tires. This results in a vicious circle, ending with a tire- and axle-wrecking hole in the road. The city of Boston decided to tackle this problem using modern technology. It released a smartphone app which used the internal accelerometer of the phone to detect the

jolt of a car being driven over a pothole and then used GPS to automatically transmit the location of the hole to the city authorities.

Wonderful! Now the highway maintenance people would know exactly where to go to repair the potholes.

Again, this looks like an elegant and cheap solution to a real problem, built on modern data analytic technology—except for the fact that ownership of cars and expensive smartphones is more likely to be concentrated in wealthier areas. Thus, it's quite likely that potholes in poorer areas would not be detected, so that their location would not be transmitted, and some areas might never have their potholes fixed. Rather than solving the pothole problem in general, this approach might even aggravate social inequalities. The situation here is different from that in the TARN example, in which we knew that certain data were missing. Here we are unaware of them.

The following is another illustration of this kind of dark data. In late October 2012, Hurricane Sandy, also called “Superstorm Sandy,”<sup>5</sup> struck the Eastern Seaboard of the United States. At the time it was the second most costly hurricane in U.S. history and the largest Atlantic hurricane on record, causing damage estimated at \$75 billion, and killing more than 200 people in eight countries. Sandy affected 24 U.S. states, from Florida to Maine to Michigan to Wisconsin, and led to the closure of the financial markets owing to power cuts. And it resulted, indirectly, in a surge in the birth rate some nine months later.

It was also a triumph of modern media. The physical storm Hurricane Sandy was accompanied by a Twitter storm of messages describing what was going on. The point about Twitter is that it tells you what and where something is happening as it happens, as well as who it's happening to. The social media platform is a way to keep up in real time as events unfold. And that's exactly

what occurred with Hurricane Sandy. Between 27 October and 1 November 2012, there were more than 20 million tweets about it. Clearly, then, we might think, this is ideal material from which to get a continuously evolving picture of the storm as it develops, identifying which areas have been most seriously affected, and where emergency relief is needed.

But later analysis revealed that the largest number of tweets about Sandy came from Manhattan, with few tweets coming from areas like Rockaway and Coney Island. Did that mean that Rockaway and Coney Island were less severely affected? Now it's true that subways and streets of Manhattan were flooded, but it was hardly the worst-hit region, even of New York. The truth is, of course, that those regions transmitting fewer tweets may have been doing so not because the storm had less impact but simply because there were fewer Twitter users with fewer smartphones to tweet them.

In fact, we can again imagine an extreme of this situation. Had any community been completely obliterated by Sandy, then no tweets at all would have emerged. The superficial impression would be that everybody there was fine. Dark data indeed.

As with the first type of dark data, examples of this second kind, in which we don't know that something is missing, are ubiquitous. Think of undetected fraud, or the failure of a crime-victim survey to identify that any murders have been committed.

You might have a sense of *déjà vu* about those first two types of dark data. In a famous news briefing, former U.S. Secretary of Defense Donald Rumsfeld nicely characterized them in a punchy sound bite, saying “there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns—the ones we don't know we don't know.”<sup>6</sup> Rumsfeld attracted considerable media ridicule for that

convoluted statement, but the criticism was unfair. What he said made very good sense and was certainly true.

But those first two types are just the beginning. In the next section we introduce some of the other types of dark data. These, and others described later, are what this book is all about. As you will see, dark data have many forms. Unless we are aware that data might be incomplete, that observing something does not mean observing everything, that a measurement procedure might be inaccurate, and that what is measured might not really be what we want to measure, then we could get a very misleading impression of what's going on. Just because there's no one around to hear that tree fall in the forest doesn't mean that it didn't make a noise.

## So You Think You Have All the Data?

The customer arrives at the supermarket checkout with a full shopping cart. The laser scans the barcode of each item, and the till emits its electronic beep as it adds up the total cost. At the end of this exercise, the customer is presented with the overall bill and pays. Except that's not really the end. The data describing the items bought and the price of each are sent to a database and stored. Later, statisticians and data scientists will pore over the data, extracting a picture of customer behavior from details of what items were bought, which items were bought together, and indeed what sort of customer bought the items. Surely there's no opportunity for missing data here? Data of the transaction have to be captured if the supermarket is to work out how much to charge the customer—short of a power cut, register failure, or fraud, that is.

Now it seems pretty obvious that the data collected are all the data there are. It's not just *some* of the transactions or details of

just *some* of the items purchased. It's *all* the transactions made by *all* the customers on *all* the items in that supermarket. It is, as is sometimes simply said, "data = all."

But is it really? After all, these data describe what happened *last week* or *last month*. That's useful, but if we are running the supermarket, what we probably really want to know is what will happen tomorrow or next week or next month. We really want to know who will buy what when, and how much of it they will buy in the future. What's likely to run out if we don't put more on the shelves? What brands will people prefer to buy? We really want data that have not been measured. Dark data *DD-Type 7: Changes with Time* describes the obscuring nature of time on data.

Indeed, beyond that complication, we might want to know how people *would have behaved* had we stocked different items, or arranged them differently on the shelves, or changed the supermarket opening times. These are called *counterfactuals* because they are contrary to fact—they are about what would have happened if what actually happened hadn't. Counterfactuals are dark data *DD-Type 6: Data Which Might Have Been*.

Needless to say, counterfactuals are of concern not just to supermarket managers. You've taken medicines in the past. You trusted the doctor who prescribed them, and you assumed they'd been tested and found to be effective in alleviating a condition. But how would you feel if you discovered that they hadn't been tested? That no data had been collected on whether the medicines made things better? Indeed, that it was possible they made things worse? Or that even if they had been tested and found to help, the medicines hadn't been compared with simply leaving the condition alone, to see if they made it get better more quickly than natural healing processes? Or the medicines hadn't been compared with other ones, to see if they were more effective than

familiar alternatives? In the elephant powder example, a comparison with doing nothing would soon reveal that *doing nothing was just as effective at keeping the elephants away* as putting down the heaps of powder. (And that, in turn could lead to the observation that there were actually no elephants to be kept away.)

Returning to the notion of “data=all,” in other contexts the notion that we might have “all” the data is *clearly* nonsensical. Consider your weight. This is easy enough to measure—just hop on your bathroom scale. But if you repeat the measurement, even very soon afterward, you might find a slightly different result, especially if you try to measure it to the nearest ounce or gram. All physical measurements are subject to potential inaccuracies as a result of measurement error or random fluctuations arising from very slight changes in the circumstances (*DD-Type 10: Measurement Error and Uncertainty*). To get around this problem, scientists measuring the magnitude of some phenomenon—the speed of light, say, or the electric charge of the electron—will take multiple measurements and average them. They might take 10 measurements, or 100. But what they obviously cannot do is take “all” the measurements. There is no such thing as “all” in this context.

A different type of dark data is illustrated when you ride on London’s red buses: you will know that more often than not they are packed with passengers. And yet data show that the occupancy of the average bus is just 17 people. What can explain this apparent contradiction? Is someone manipulating the figures?

A little thought reveals that the answer is simply that more people are riding on the buses when they are full—that’s what “full” means. The consequence is that more people see a full bus. At the opposite extreme, an empty bus will have no one to report that it was empty. (I’m ignoring the driver in all this, of



course.) This example is an illustration of dark data *DD-Type 3: Choosing Just Some Cases*. Furthermore, that mode of dark data can even be a necessary consequence of collecting data, in which case it illustrates *DD-Type 4: Self-Selection*. The following are my two favorite examples of opposite extremes in terms of significance.

The first is the cartoon showing a man looking at one of those maps which are placed outside railway stations. In the middle of the map is a red dot with a label saying “You are here.” “How,” thinks the man, “did they know?” They knew because they recognized that *everyone* looking at the red dot had to be in front of the sign. It was a highly selected sample and *necessarily* missed everyone standing elsewhere.

The point is that data can be collected only if there is someone or something—a measuring instrument, for example—there to collect them. And the second extreme manifestation of this is described by the *anthropic principle*, which essentially says that the universe has to be like it is, or we would not be here to observe it. We cannot have data from very different universes because we could not exist in those and so could not collect data from them. This means any conclusions we draw are necessarily limited to our (type of) universe: as with the potholes, there might be all sorts of other things going on which we don’t know about.

There’s an important lesson for science here. Your theory might be perfectly sound for your data, but your data will have limits. They might not refer to very high temperatures or long times or vast distances. And if you extrapolate beyond the limits within which your data were collected, then perhaps your theory will break down. Economic theories built on data collected during benign conditions can fail dramatically in recessions, and Newton’s laws work fine unless tiny objects or high velocities or

other extremes are involved. This is the essence of *DD-Type 15: Extrapolating beyond Your Data*.

I have a T-shirt with an *xkcd* cartoon with two characters talking to each other. One character says “I used to think correlation implied causation.” In the next frame, he goes on to say, “Then I took a statistics class. Now I don’t.” Finally, the other character says, “Sounds like the class helped,” and the first character replies, “Well, maybe.”<sup>7</sup>

Correlation simply means that two things vary together: for example, positive correlation means that when one is big then the other is big, and when the first is small, the second is small. That’s different from causation. One thing is said to *cause* another if a change in the first induces a change in the second. And the trouble is that two things can vary together without changes in one being the cause of changes in the other. For example, observations over the early years of schooling show that children with a larger vocabulary tend, on average, to be taller. But you wouldn’t then believe that parents who wanted taller offspring should hire tutors to expand their vocabulary. It’s more likely that there are some unmeasured dark data, a third factor which accounts for the correlation—such as the ages of the children. When the *xkcd* character says, “Well, maybe,” he’s acknowledging that it’s possible that taking the statistics class caused his understanding to change, but maybe there was some other cause. We shall see some striking examples of this situation, characterized by *DD-Type 5: Missing What Matters*.

I’ve now mentioned several dark data types. But there are more. The aim of this book is to reveal them, to show how they can be identified, to observe their impact, and to show how to tackle the problems they cause—and even how to take advantage of them. They are listed at the end of this chapter, and their content is summarized in chapter 10.

## Nothing Happened, So We Ignored It

A final example illustrates that dark data can have disastrous consequences and that they are not especially a problem of large data sets.

Thirty years ago, on 28 January 1986, 73 seconds into its flight and at an altitude of 9 miles, the space shuttle *Challenger* experienced an enormous fireball caused by one of its two booster rockets and broke up. The crew compartment continued its trajectory, reaching an altitude of 12 miles, before falling into the Atlantic. All seven crew members, consisting of five astronauts and two payload specialists, were killed.

A later presidential commission found that NASA middle managers had violated safety rules requiring data to be passed up the chain of command. This was attributed to economic pressures, making it very important that the launch schedule should be maintained: the launch date had already slipped from January 22nd to the 23rd, then to the 25th, and then to the 26th. Since temperature forecasts for that day suggested an unacceptably low temperature, the launch was again rescheduled, for the 27th. Countdown proceeded normally until indicators suggested a hatch lock had not closed properly. By the time that was fixed the wind was too strong, and again the launch was postponed.

On the night of January 27th, a three-hour teleconference was held between Morton Thiokol, which was the company that made the booster rockets, NASA staff at the Marshall Space Flight Center, and people from the Kennedy Space Center. Larry Wear, of the Marshall Center, asked Morton Thiokol to check the possible impact of low temperatures on the solid rocket motors. In response, the Morton Thiokol team pointed out that the O-rings would harden in low temperatures.

The O-rings were rubber-like seals, with a cross-section diameter of about a quarter of an inch, which fitted in the joint around the circumference between each of the four rocket motor segments. The solid rocket boosters were 149 feet high and 38 feet in circumference. Under launch conditions, the 0.004 inch gap that the O-rings normally sealed typically opened to a maximum of 0.06 inch: just six one-hundredths of an inch. And during launch this larger gap remained open for just six-tenths of a second.

Robert Ebeling of Morton Thiokol had been concerned that at low temperatures the hardening of the O-rings meant they would lose their ability to create an effective seal between segments when the gaps expanded by that 0.056 inch for that 0.6 second. At the teleconference Robert Lund, vice president of Morton Thiokol, said that the O-ring operating temperature must not be less than the previous lowest launch temperature, 53°F. Extensive, sometimes heated, discussion ensued, both in the conference and off-line in private conversations. Eventually, Morton Thiokol reconsidered and recommended launch.

Precisely 58.79 seconds after the launch a flame burst from the right solid rocket motor near the last joint. This flame quickly grew into a jet which broke the struts joining the solid rocket motor to the external fuel tank. The motor pivoted, hitting first the Orbiter's wing and then the external fuel tank. The jet of flame then fell onto this external tank containing the liquid hydrogen and oxygen fuel. At 64.66 seconds the tank's surface was breached, and 9 seconds later *Challenger* was engulfed in a ball of flame and broke into several large sections.<sup>8</sup>

One thing we have to remember is that space travel is all about risk. No mission, even under the very best of circumstances, is a risk-free enterprise: the risk cannot be reduced to zero. And there are always competing demands.

Furthermore, as with any incident like this, the notion of “cause” is complicated. Was it due to violation of safety rules, undue pressure put on managers because of economic considerations, other consequences of budget tightening, or perhaps media pressure following the fact that the launch of the previous shuttle, *Columbia*, had been delayed seven times, each delay greeted with press ridicule? For example, here’s Dan Rather’s script for the evening news on Monday, January 27th, following the four delays to the *Challenger* launch: “Yet another costly, red-faces-all-around space-shuttle-launch delay. This time a bad bolt on a hatch and a bad-weather bolt from the blue are being blamed.” Or was it a consequence of political pressure. After all, there was significantly more interest in this launch than earlier launches because it carried an “ordinary person,” Christa McAuliffe, a teacher, and the president’s State of the Union address was scheduled for the evening of January 28th.

In such situations, multiple factors typically come together. Complex and obscure interactions can lead to unexpected consequences. But in this case there was another factor: dark data.

After the disaster, a commission headed by former secretary of state William Rogers drew attention to the fact that flights which had not had any O-rings showing distress had not been included in the diagram discussed at the teleconference (dark data *DD-Type 3: Choosing Just Some Cases* but also *DD-Type 2: Data We Don’t Know Are Missing*). The report said (p. 146): “The managers compared as a function of temperature the flights for which thermal distress of O-rings had been observed—not the frequency of occurrence based on all flights.”<sup>9</sup> And that’s the giveaway: *data from some flights were not included in the analysis*. My earlier examples have shown the sorts of problems leaving out some of the data can lead to.

The report went on: “In such a comparison [that is, using the limited set of data presented], there is nothing irregular in the distribution of O-ring ‘distress’ over the spectrum of joint temperatures at launch between 53 degrees Fahrenheit and 75 degrees Fahrenheit,” meaning: there is no apparent relationship between temperature and number of O-rings showing distress. However, “when the entire history of flight experience is considered, including ‘normal’ flights with no erosion or blow-by, the comparison is substantially different”; that is, if you include all the data, you get a different picture. In fact, flights which took place at higher temperatures were much more likely to show no problems, and these were the dark data not shown in the plot. But if the higher the temperature, the less the chance of a problem, then, conversely, the lower the temperature, the greater the chance of a problem. And the ambient temperature was predicted to be just 31°F.

This section of the report concluded: “Consideration of the entire launch temperature history indicates that the probability of O-ring distress is increased to *almost a certainty* if the temperature of the joint is less than 65[°F].” (my italics)

The situation is graphically illustrated in the two diagrams in Figure 1. Figure 1(a) shows the diagram discussed at the teleconference. This is a plot of the number of distressed O-rings on each launch plotted against launch temperature in degrees Fahrenheit. So, for example, at the lowest launch temperature in the past, 53°F, three of the O-rings experienced distress, and at the highest launch temperature in the past, 75°F, two of the O-rings experienced distress. There is no clear relationship between launch temperature and the number of distressed O-rings.

However, if we add the missing data—showing the launches which led to no O-ring distress, we obtain Figure 1(b). The pattern is now very clear. In fact, *all* the launches which occurred

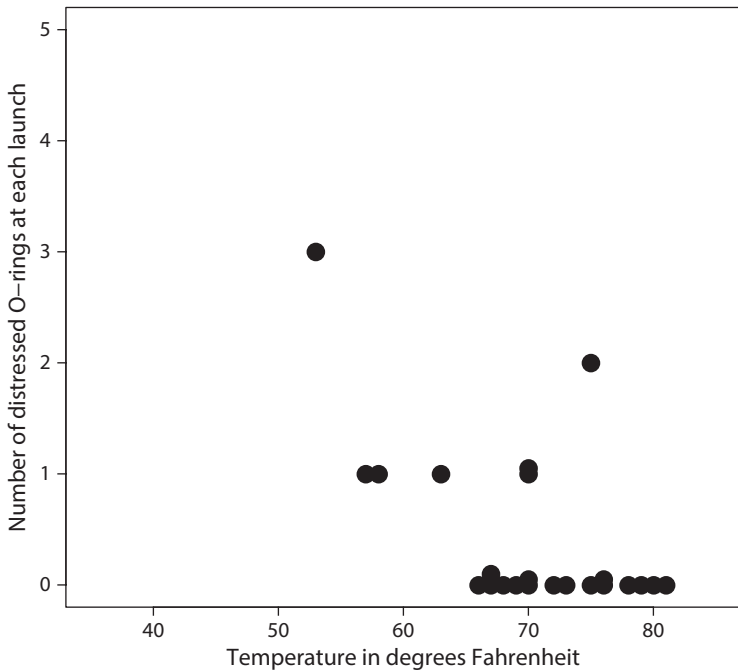
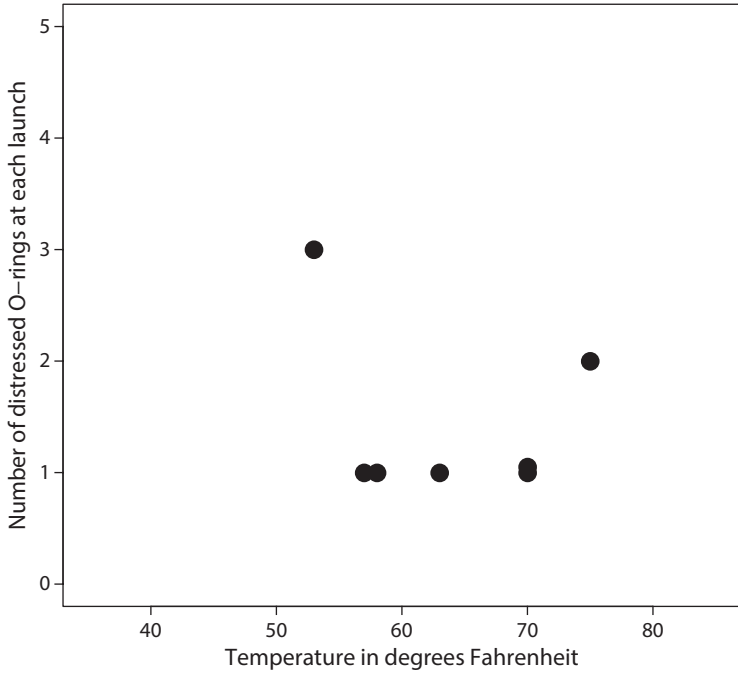


FIGURE 1. (a) Data examined in the *Challenger* prelaunch teleconference; (b) complete data.

when the temperature was less than 65°F experienced some O-ring distress, but only four of the 21 launches which took place at higher temperatures led to O-ring distress. The figure shows that the lower the temperature, the greater the risk. And worse, the predicted launch temperature was way below any previously experienced (*DD-Type 15: Extrapolating beyond Your Data*).

The missing data are crucial to understanding what is going on.

There is an interesting corollary to this story. Although it took months for the official report to arrive at its conclusion, Morton Thiokol's stock price crashed by 11.86 percent on the day of the disaster. Based on Morton Thiokol price movements prior to the incident, stock price changes of even as much as 4 percent would be rare. The stock prices of other companies involved in constructing the shuttle launch vehicle also slumped, but by far less. It was as if the market knew what was responsible for the crash. Dark data again?

## The Power of Dark Data

That last example shows just how catastrophically badly wrong things can go when we fail to allow for dark data. It seems that they represent a real threat. But in fact the picture is not all gloomy. It turns out that an understanding of dark data can be used to our advantage, in a sort of data science judo. We can do this in several ways, as I describe in the second part of the book. Here's one way.

I introduce so-called randomized controlled trials in chapter 2, and in chapter 9 I return to look at them from a different perspective. In a medical context, in the simplest such trial two treatments are compared by giving one treatment to one group of people and the other to another group. However, there is a



risk. If the researchers know which treatment has been given to which people, this knowledge could influence the study. The researchers might be tempted to treat one of the groups more carefully than the other. For example, if the study aimed to compare an untested new treatment with the standard treatment, the researchers might monitor the former group for side effects more closely (perhaps subconsciously) or take more care in their measurement of possible outcomes. To overcome this potential bias, in such studies the treatment allocation is concealed from the researchers (*DD-Type 13: Intentionally Darkened Data*). The term *blinded* is used, to indicate that the data are dark.

Another familiar way in which dark data are used to advantage is in sample surveys. We might want to know the opinions of the people in a town, or of those who buy a company's products, and perhaps it is prohibitively expensive to ask them all. Certainly it is likely to be a time-consuming exercise, and opinions might change over its course. An alternative to asking everyone is to ask just some of them. The opinions of the others, the ones you don't ask, will be dark data. This might look like a high-risk strategy—it clearly resembles the TARN example. But it turns out that by using judicious methods of selecting which people you approach, you can get accurate and reliable answers—more quickly and cheaply than by attempting to approach everyone.

Yet a third way to use dark data to advantage is through so-called smoothing of data. As we shall see in chapter 9, this is equivalent to revealing unobserved and unobservable kinds of dark data (*DD-Type 14: Fabricated and Synthetic Data*), and it enables more accurate estimates and better predictions.

We shall explore other uses of dark data in chapter 9, where we will see that they often have exotic names. Some of them are widely applied in areas such as machine-learning and artificial intelligence.

## All around Us

We've seen that dark data are ubiquitous. They can arise anywhere and everywhere, and one of their most dangerous aspects is that, by definition, we may not know that they are *not* there. It means we have to be constantly on the alert, asking ourselves, *what are we missing?*

Are we failing to notice large amounts of fraud because the police catch the inept criminals while the really good ones escape unnoticed? Bernie Madoff established his firm Bernard L. Madoff Investment Securities LLC in 1960 but wasn't arrested until 2008, and sentenced (to 150 years in prison) in 2009, when he was already 71—he almost got away with it.

Are we not noticing many potentially curable sick people simply because the more severe cases are obvious, but the less severe don't show so many symptoms?

Are the social networks established by modern social media dangerous simply because they reflect what we already know and believe, not challenging us because they don't show us facts or events outside our comfort zone?

Perhaps worse still, the descriptions people choose to post on social media may give us a false impression of how wonderful everyone else's life is, casting us into depression because in contrast our lives have so many obstacles.

We tend to think of data as numerical. But data don't have to be just numbers. And that means that dark data also don't have to be numerical. The following is an example in which the crucial missing information is a single letter.

The Arctic expeditions of 1852, 1857, and 1875 were stocked with a supply of Allsopp's Arctic Ale, an ale with an especially low freezing point prepared by brewer Samuel Allsopp. Alfred Barnard sampled the beer in 1889, describing it as "of a nice brown

colour, and of a vinous, and at the same time, nutty flavor, and as sound as on the day it was brewed. . . . Owing to the large amount of unfermented extract still remaining in it, it must be considered as an extremely valuable and nourishing food.”<sup>10</sup> Just the sort of thing you need to sustain you on Arctic expeditions.

In 2007 a bottle of the 1852 batch came up for sale on eBay, with a reserve price of \$299. Or at least that was the aim. In fact the vendor, who had had the bottle for 50 years, misspelled the beer’s name—he missed one of the *p*’s in Allsopp. As a consequence, the item didn’t show up in the searches carried out by most vintage beer enthusiasts, so that there were just two bids. The winning bid, for \$304, was from 25-year-old Daniel P. Woodul. Aiming to appraise the value of the bottle, Woodul immediately relisted it on eBay, but this time with the correct spelling. This time there were 157 bids, with the winning one being for \$503,300.

That missing *p* clearly mattered, to the tune of some half a million dollars.\* This shows that missing information can have significant consequences. In fact, as we shall see, a mere half-million-dollar loss is nothing compared with the losses that other missing data situations have led to. Indeed, missing data can wreck lives, destroy companies, and (as with the *Challenger* disaster) can even lead to death. In short, missing data matter.

In the case of Allsopp’s Arctic Ale, a little care would have avoided the problem. But while carelessness is certainly a common cause of dark data, there are many others. The painful fact

\*In fact it turned out that the winning bid was a practical joke, and the bidder had no intention of paying. But Woodul is nevertheless doubtless still sitting on a tidy profit: a private collector from Scotland recently auctioned a bottle from the 1875 expedition for £3,300 (~\$4,300).

is that data can be dark for a tremendously wide variety of reasons, as we shall see in this book.

It is tempting to regard dark data as simply synonymous with data which could have been observed but which for some reason were not. That is certainly the most obvious kind of dark data. The missing salary levels in a survey in which some people refused to divulge how much they were paid are certainly dark data, but so also are the salary levels for those who do not work and hence do not have a salary level to divulge. Measurement error obscures true values, data summaries (such as averages) hide the details, and incorrect definitions misrepresent what you want to know. More generally still, any unknown characteristic of a population can be thought of as dark data (statisticians often refer to such characteristics as *parameters*).

Since the number of possible causes of dark data is essentially unlimited, knowing what *sort* of thing to keep an eye open for can be immensely useful in helping avoid mistakes and missteps. And that is the function of the *DD-Types* described in this book. These are not basic causes (like failure to include the final outcome for patients who have been in a study for only a short time) but provide a more general taxonomy (like the distinction between data we know are missing and data we don't know are missing). An awareness of these *DD-Types* can help in protecting against mistakes, errors, and disasters arising from ignorance about what you do not know. The *DD-Types*, which are introduced in this book, and which are summarized in chapter 10, are as follows:

*DD-Type 1: Data We Know Are Missing*

*DD-Type 2: Data We Don't Know Are Missing*

*DD-Type 3: Choosing Just Some Cases*

*DD-Type 4: Self-Selection*

- DD-Type 5: Missing What Matters*
- DD-Type 6: Data Which Might Have Been*
- DD-Type 7: Changes with Time*
- DD-Type 8: Definitions of Data*
- DD-Type 9: Summaries of Data*
- DD-Type 10: Measurement Error and Uncertainty*
- DD-Type 11: Feedback and Gaming*
- DD-Type 12: Information Asymmetry*
- DD-Type 13: Intentionally Darkened Data*
- DD-Type 14: Fabricated and Synthetic Data*
- DD-Type 15: Extrapolating beyond Your Data*

## INDEX



- 419 scam, 148
- A/B trial, 58
- abandoned calls, 37, 38
- acquiescence bias. *See* biases:  
    subconscious
- adaptive allocation, 65
- administrative data, 31–45, 73, 74, 112,  
    284
- Administrative Data Research  
    Network, 112, 284
- adverse selection, 130, 138
- advertisements, 67
- AIDS, 92
- Air Canada, 106
- Akerlof, George, 129, 137, 197
- Alitalia Airlines, 105
- Allsopp's Arctic Ale*, 24, 25
- Alzheimer's disease, 76, 77, 93
- Amazon, 66, 300
- American Academy of Pediatrics, 77
- American Psychological Association,  
    201
- Amsterdam City Council, 104
- Analysis and Detection Center, 155
- anonymizing data, 283, 284, 286, 287
- anthropic principle, 15
- AOL, 41, 286
- Apple Inc., 281
- Argentinian inflation rate, 300
- arithmetic mean. *See* mean: arithmetic
- arthroscopic surgery, 180, 181
- artificial intelligence, 23, 80, 155, 269, 302
- asthma, 58, 60, 61, 130, 131
- Atlantic salmon, 194
- Atlas Venture*, 184
- Atomic Energy Establishment, 175
- autism, 68, 77, 217
- automated transactions, 80
- availability bias. *See* biases:  
    subconscious
- Aviva Insurance, 161
- axiom system, 115
- B-2 Spirit stealth bomber, 109
- Babbage, Charles, 202, 211–215
- Bacon, Francis, 168
- Baesens, Bart, 142, 143
- balanced designs, 58, 65, 246
- bandwagon effect. *See* biases:  
    subconscious
- Bankia, 157
- banknotes, paper, 123, 142
- Banque Générale*, 123
- Barclay Hedge Fund Index, 233
- Barings, 154
- Barrett, Amy, 39
- Barth-Jones, Daniel, 286
- base rate fallacy, 68, 69
- basket of goods, 78–80
- bathroom scales, 110
- Bayesian statistics, 268, 269, 276–278
- BBC, 37
- Bebo, 41
- Begley, C. Glenn, 185
- Behavioural Insights Team, 291

- Behavioral Risk Factor Surveillance Survey, 52
- Belgian Constitutional Court, 135
- belief bias. *See* biases: subconscious
- Benford distribution, 258, 259, 278
- Benjamini, Yoav, 197
- Bennett, Craig, 195
- Beringer, Johann Bartholomeus  
Adam, 203, 204
- Betamax, 304
- biases, subconscious: acquiescence bias, 70; availability bias, 67–69, 77; bandwagon effect, 70; belief bias, 70; bizarreness effect, 70; confirmation bias, 69, 70, 127, 168, 181, 203, 207; conjunction fallacy, 69; negativity bias, 70
- bicycle speed record, 260
- big bang, 178, 182
- Billion Prices Project, 300, 301
- bipolar disorder, 180
- bizarreness effect. *See* biases: subconscious
- blinding, 264
- blood pressure, 58, 99, 100, 103, 242, 272; diastolic, 100, 272; systolic, 100, 272
- bloodletting, 59, 179
- Blunt, Gordon, 143
- BMI. *See* body mass index
- body mass index, 58, 136, 225–229
- Boesky, Ivan, 154, 155
- Boghossian, Peter, 206
- Bohannon, John, 205
- boiler room scams, 166
- bomb: hydrogen, 173; fission, 174
- bomber. *See* B-2 Spirit stealth bomber
- Bonferroni correction, 197. *See also* false discovery rate
- boosting, 273, 298
- Booth, Bruce, 184, 190
- bootstrapping, 273, 274
- bottomcoding, 101, 110
- BP Deepwater Horizon, 161
- Bradley, James, 278
- brain damage, 3, 196
- Brewer, Devon, 101
- British Crime Survey, 74
- British Hypertension Guidelines, 100
- British Transport Police, 75
- Briton, Deborah, 160
- Broad, William, 214
- Brodeski, Brent, 39
- Brown, Herbert, 183
- Brownian motion, 215
- Bryson, Maurice, 49
- bubbles, financial, 122–127, 297
- budesonide, 60
- Bureau of Labor Statistics, 79, 121
- Burt, Cyril, 201, 202
- Byrd, Malcolm, 147
- Byrne, John, 181
- caffeine, 103, 104
- calorie consumption, 291–293
- camouflage, 142
- Campbell's law, 116
- Canadian Supreme Court, 114
- cancer, 81, 92–94, 185, 209, 225, 242–244, 265
- card not present transactions, 151
- cardiovascular outcomes, 181
- car-parts, 65, 246
- Caruana, Rich, 130
- causality, 41, 81, 190, 264, 295
- causation, 16
- Cavallo, Alberto, 300
- Cavendish Laboratory, 176
- ceiling effect, 110, 297
- cellphones, 79

- censored data, 243, 244
- census, 28, 29, 43, 44, 45, 56, 107, 206
- Census Bureau, 107, 154
- Centers for Disease Control, 4
- Challenger*. *See* space shuttle
- Chalmers, Iain, 59
- champion/challenger trial, 58
- Chancellor of the Exchequer, 105
- charge of the electron. *See* electron
- check digit, 257
- chemotherapy, 92
- Chinese Ministry of Science and Technology, 209
- chip and PIN, 151
- Churchill, Randolph, 105
- Churchill, Winston, 105
- churning. *See* fee churning
- Clever Hans, 302
- climate change, 181
- Climate Orbiter, 106
- clinical trial, 59, 60, 65, 66, 84, 190, 197, 257, 283, 296
- Coase, Ronald, 192
- Cobange, Ocorrafoo, 205
- cocaine, 147, 234
- coin: fair, 193, 267; French, 124; toss, 193, 266, 267, 275, 276, 280, 288
- cold fusion, 174, 175, 185
- Colonial Bank, 157
- Columbia*. *See* space shuttle
- Commission for Health Improvement, 120
- complete case analysis, 236
- confidence trick, 141
- confirmation bias. *See* biases: subconscious
- Congressional Budget Office, 138
- Congressional report on credit scoring, 133
- conjunction fallacy. *See* biases: subconscious
- consumer banking, 32, 240
- Consumer Price Index, 79
- contingent nature of science, 169
- cooking data, 203, 213–215
- Corey, Robert, 176
- correlation, 81, 89, 190, 191, 196, 202, 234, 238; causation, 16; coefficients, 102, 201; definition, 16
- cosmic microwave background radiation, 182
- cosmological constant, 178, 179
- counterfactuals, 13, 30, 58, 244, 265, 296
- counterfeit, 142; drugs, 142; machine, 141; crash-for-cash, 160
- credit card, 29, 32, 33, 42, 142, 145, 148–152, 259
- Credit Sesame, 135
- Crews, Frederick, 171
- Crick, Francis, 176
- cricket bat, 207
- crime maps, 127
- crime prevention, 62
- Crime Survey for England and Wales, 74, 75, 143
- Criterion Capital Partners, 41
- cryptography, 280, 281
- CSE&W. *See* Crime Survey for England and Wales
- Cubelli, Roberto, 196
- curbstoning, 56, 107, 206
- customer loyalty, 66
- Dalton, John, 200
- dark energy, 179
- dark matter, 4
- Darsee, John, 208, 209
- Darwin, Charles, 172, 173, 204, 208



- Dasatinib, 209
- data errors, 127; correction of, 260;  
detection of, 257; prevention of, 256
- data exhaust, 31, 42
- data mining, 80, 107, 159, 166
- data shadow, 42
- data=all, 14, 56
- Datashield, 146
- Dawkins, Richard, 206
- Dawson, Charles, 206–208
- Day of the Jackal scam, 145
- DD-Type 1: Data We Know Are Missing*,  
7, 26, 33, 46, 48, 50, 60, 110, 223, 293,  
294
- DD-Type 2: Data We Don't Know Are  
Missing*, 9, 19, 26, 38, 39, 50, 119, 194,  
223, 293, 294
- DD-Type 3: Choosing Just Some Cases*,  
15, 19, 26, 70, 95, 187, 214, 223, 295
- DD-Type 4: Self-Selection*, 15, 26, 33, 37,  
46, 48, 50, 189, 223, 293, 295
- DD-Type 5: Missing What Matters*, 16,  
27, 81, 82, 295
- DD-Type 6: Data Which Might Have  
Been*, 13, 27, 296
- DD-Type 7: Changes with Time*, 13, 27,  
39, 60, 107, 223, 256, 296, 301
- DD-Type 8: Definitions of Data*, 27, 74,  
77, 78, 121, 224, 296
- DD-Type 9: Summaries of Data*, 27, 102,  
137, 223, 256, 297
- DD-Type 10: Measurement Error and  
Uncertainty*, 14, 27, 110, 223, 256, 293,  
297
- DD-Type 11: Feedback and Gaming*, 27,  
75, 115, 189, 292, 293, 297
- DD-Type 12: Information Asymmetry*,  
27, 128, 154, 224, 297
- DD-Type 13: Intentionally Darkened  
Data*, 23, 27, 141, 298
- DD-Type 14: Fabricated and Synthetic  
Data*, 23, 27, 56, 203, 217, 269, 298
- DD-Type 15: Extrapolating beyond Your  
Data*, 16, 22, 27, 167, 224, 298
- de Lusignan, Simon, 99
- death-penalty sentences, 87, 88
- death spiral, 137, 139
- Deepwater Horizon. *See* BP Deepwater  
Horizon
- Della Sala, Sergio, 196
- dementia, 76–78, 249
- Dempster, Arthur, 254
- Department of Education, 77
- depression, 225
- deuterium, 173, 174
- diabetes, 225
- diagnosis, 68, 76–78, 93, 270, 296
- diastolic blood pressure. *See* blood  
pressure: diastolic
- Dicke radiometer, 182
- Dicke, Robert, 182
- Dickens, Charles, 164
- digit preference, 101
- digit transposition, 106
- Direct Line Insurance Group, 127
- Dirksen, Everett, 154
- disclosure control, computational,  
285
- disclosure requirements, 130
- discretizing, 101
- discrimination, 132, 135. *See also*  
disparate impact; disparate  
treatment; Equality Act; protected  
characteristics
- disguise, 40, 143, 163
- disparate impact, 132
- disparate treatment, 132
- Dixon, Jeane, 69, 70
- DJIA. *See* Dow Jones Industrial  
Average

- DNA, 175, 176
- Donohue, Jerry, 176
- dot-com bubble, 126
- doubly labelled water studies, 292
- Dow Jones Industrial Average, 40, 232, 233
- dropouts, 59, 60, 61, 225–228, 248
- drugs: counterfeit, 142; possession of, 75, 147; regulatory approval of, 190; trials for, 84, 85, 264; use of, 101
- dumpster diving, 145
- duplicate records, 113, 302
- Dvorsky, Georg, 243
- echo chamber, 71, 128
- ecological fallacy, 89
- Eddington, Arthur, 168, 169
- Edward Lee Thorndike Award, 201
- Efron, Brad, 273, 275
- Eiffel Tower, 140
- Einstein, Albert, 168, 169, 177–179
- Einstein's biggest blunder, 179
- electoral roll, 45, 285, 286
- electron, 14, 175, 213
- electron microscope, 175
- elephant powder, 3, 4, 14
- Ellis, Lee M., 185
- EM algorithm, 254, 255
- emergency calls, 37
- Enron, 156, 157
- Equality Act, 131
- European Court of Justice, 135
- exercise and weight, 291
- exercise, effect of caffeine on, 103
- experimental data, 30, 56, 299
- external examiner system, 118
- extreme gradient boosting, 273
- Facebook, 41, 42, 66, 67, 147
- factor loadings, 102
- fake: accidents, 160; data, 202, 206, 211; deaths, 159, 160; news, 128, 217, 218; tooth, 207; website, 152
- false discovery rate, 198
- false facts, 128, 217
- false negatives, 94
- false positives, 94, 194
- falsifiability, 167, 170, 172
- Fang, Ferric, 190
- fat finger errors, 105, 113
- FBI, 161, 281
- Federal Trade Commission, 143
- fee churning, 158
- feedback, 75, 77, 114, 116, 122–128, 129, 189, 291, 293, 297
- Fenofibrate, 181
- file-drawer effect, 188
- Fisher, Ronald, 65, 81
- fission bomb. *See* bomb: fission
- flat Earth, 170
- Fleischmann, Martin, 174, 175
- Fleming, Alexander, 183
- flight simulators, 266
- floor effect, 110, 297
- Foreign and Commonwealth Office, 160
- forging, 203, 206, 208–210, 215
- fossil, 203; fuel, 173, 174, 299; radiation, 182
- Foucault, Léon, 278
- fractional factorial designs, 65
- Franklin, Rosalind, 176
- fraud: advance fee fraud, 148; credit card fraud, 149, 151; e-commerce fraud, 143; financial fraud, 149–158, 166; hard fraud, 158; insurance fraud, 158–163; internet fraud, 144–149; mortgage fraud, 159; scientific fraud, 202, 216, 217; soft fraud, 158; tax fraud, 165;

- fraud (*continued*)  
undetected fraud, 11. *See also*  
confidence trick; Benford's  
distribution; insider trading;  
curbstoning
- Freedman, Leonard, 185
- Freud, Sigmund, 171, 172
- Galileo, 110, 170, 200, 278
- Galton, Francis, 97
- gambling, 163, 164
- gaming, 114–122
- Gamow, George, 179
- GCSEs, 120
- GDPR. *See* General Data Protection  
Regulation
- Gender Directive, 114, 135
- General Data Protection Regulation,  
35, 219
- General Election, UK, 49
- Genesis probe, 106
- Gentry, Craig, 289
- geometric mean. *See* mean: geometric
- Gigerenzer, Gerd, 92
- Global Synthetic Equities, 153
- Gödel, Kurt, 115
- gold sample, 239, 241
- Goodhart's law, 116
- Goodstein, David, 214
- Google, 301
- grade inflation, 117, 118, 297
- gravity, 5, 171, 173, 178, 204, 213, 214
- Greco-Latin square designs, 65
- "greed is good speech," 155
- Greenwood, Elizabeth, 159, 160
- grievance studies, 206
- group residences, 75
- HARKing, 198, 199
- harmonic mean. *See* mean: harmonic
- Harvard, 208, 234; medical students  
at, 68
- Harwell, 175
- Hawking, Stephen, 172
- Hawthorne effect, 65, 122
- Health and Social Care Information  
Centre, 35
- heaping, 101, 297
- Hearst, William Randolph, 218
- Heckit, 232
- Heckman, James, 232
- helix, 176
- Herschel, William, 183
- hoaxing, 203–206, 215
- Hochberg, Yosef, 197
- homeopathy, 179
- homomorphic computation, 289
- "horse" in machine learning, 302
- hot deck imputation, 251
- housing bubble, 126
- Houston Natural Gas, 156
- Hoyle, Fred, 177, 178
- HSCIC. *See* Health and Social Care  
Information Centre
- Hugh-Jones, David, 288
- humbug, 204, 208
- hurricane, 268; Katrina, 161; Sandy, 10, 11
- Hurwitz, Michael, 117
- hydroboration, 183
- hydrogen bomb. *See* bomb: hydrogen
- hypotensive drug, 242
- IBM, 108, 289
- identity theft, 144–149
- Ig Nobel Prize, 195
- immigration, 73–74, 76
- immunization, 4
- impact factor, 189, 190
- The Improbability Principle*, 107
- Impulse Airlines, 155

- imputation, 112, 240, 245–252, 255;  
    multiple, 252
- incontinence, 180
- individual mandate, 138
- inflation, 78, 79, 171, 300, 301. *See also*  
    grade inflation
- information asymmetry, 114, 128–130,  
    139, 297
- informatively missing. *See* missing
- injection-moulded, 65, 246
- insider trading, 142, 153–158, 165, 297
- Insurance Information Institute, 146
- insurance premium, 135, 136, 159
- integration, in money laundering, 163
- intelligence, 123, 130, 171, 201. *See also*  
    artificial intelligence
- International Passenger Survey, 73
- internet, 54, 55, 66, 126, 142, 144–149,  
    151, 205, 287
- Internet Movie Database, 287
- InterNorth, 156
- investment funds, 39
- IPS. *See* International Passenger Survey
- irrigation, 65
- Javelin Strategy and Research, 146
- Jehovah, 203
- John, Ioannidis, 184, 186, 189
- Johnston, Mike, 38
- Kahneman, Daniel, 141
- Kamin, Leon, 201
- Kaplan, Edward, 243
- Kelvin. *See* Lord Kelvin
- Keynes, John Maynard, 123
- Kho, Michelle, 37
- Kim, Joonseok, 181
- Kohler, Eddie, 205
- Kruskal, William, 107
- Kuhn, Thomas, 183
- Labour Force Survey, 51, 52
- Laird, Nan, 254
- land speed bicycle record. *See* bicycle  
    speed record
- Landon, Alfred, 48–51, 55
- last observation carried forward,  
    248–249
- latent variable models, 255
- Law, John, 123–125
- law of large numbers, 43–45, 267
- law of likelihood, 253
- Lay, Kenneth, 156, 157
- layering, in money laundering, 163
- league tables, 118, 119
- Leaning Tower of Pisa, 170
- Lee, Jason, 117
- Leeson, Nick, 153
- Lehman Brothers, 41, 104
- Leigh, Andrew, 62, 64, 66
- Lemaitre, Georges, 178
- lemons, 129, 137, 240, 298
- length-time bias, 93
- lethargy, 180
- light rays, 178
- likelihood, 253, 254
- Lincoln Emergency Communications  
    Center, 38
- linking data, 109, 111–113, 284, 301
- literacy, 89, 90
- Literary Digest*, 49, 50
- lithium, 173, 174
- Little, Roderick, 235
- Little Dorrit*, 164
- Living Costs and Food Survey, 292
- Livio, Mario, 179
- loans, 32, 241
- Local Government Transparency  
    Code, 219
- LOCF. *See* last observation carried  
    forward

- log cabins, 38  
logarithms, 258  
London buses, 14  
London Zoo, 28  
Long Term Capital Management, 304  
Long-Term International Migration, 73  
Lord Kelvin, 172, 173, 174, 177  
Los Angeles County, 112  
lottery, 41, 42, 59  
LTCM. *See* Long Term Capital Management  
LTIM. *See* Long-Term International Migration  
lung: cancer, 81; function, 60  
Lustig, Victor, 140, 141, 144  
  
machine-learning, 23, 130, 131, 155, 166, 269, 272, 273, 302  
Madoff, Bernie, 24, 165, 210  
Madoff Investment Securities LLC, 24  
magnetic stripe, 150, 152  
Malmquist bias, 111  
Malmquist, Gunnar, 111  
market for lemons, 137, 298  
Manning, Bradley/Chelsea, 130  
Mars, 106  
Marsh, Cathie, 229, 232  
Marshall Space Flight Center, 17  
*Martin Chuzzlewit*, 164  
Massachusetts Group Insurance Commission, 285  
maximum likelihood, 253, 254  
May, Theresa, 127  
Mazières, David, 205  
McAuliffe, Christa, 19  
mean: arithmetic mean, 78, 103, 233; geometric mean, 78; harmonic mean, 78  
measles, 3, 4  
Medical Research Council, 59  
  
Meier, Paul, 243  
Mendel, Gregor, 200  
Meng, Xiao-Li, 234  
Meteorological Office, 109  
miasma theory of disease, 168  
Michelson, Albert, 177  
Michelson-Morley experiment, 177  
micron, 98  
Millikan, Robert, 200, 213–215  
minimum income, 62  
Mirkes, Evgeny, 7, 8  
missing: informatively missing, 226; missing at random, 226, 227; missing completely at random, 227; non-ignorably missing, 226, 227  
Mississippi Company, 124–126  
Mizuho Securities Co, 105  
money laundering, 158, 163, 259  
money printing machine, 141  
Moniz, António Egas, 180  
Morgan Stanley, 153, 154  
Morningstar database, 39  
Morton Thiokol, 17, 18, 22  
Moseley, Bruce, 181  
MRI scan, 195  
Mueller-Korenek, Denise, 260  
multiple comparisons, 194, 195  
multiple imputation. *See* imputation: multiple  
multiple regression, 90  
multivitamins, 181  
murder, 11, 75, 87, 146  
  
NA, 46, 82, 236, 237, 239, 251  
Narayanan, Arvind, 287  
NASA, 17, 106  
NASDAQ, 126  
Nashef, Samer, 266  
National Crime Victimization Survey, 74

- National Dementia Strategy, 78  
National Diet and Nutrition Survey  
and Health Survey, 292  
National Health Interview Survey, 52, 53  
National Health Service, 35, 266  
National Insurance Numbers, 73  
National Statistical Institute, 112, 281  
Natural History Museum, 206  
NDD. *See* not data dependent  
near-field communication, 150  
negativity bias. *See* biases: subconscious  
neighborhood policing, 63  
Netflix Prize, 286, 287  
new oil, 299  
New York Police Department, 127  
Newcomb, Simon, 258, 278  
Newton, Isaac, 125, 126, 172, 200  
Newton's Laws, 15, 169, 171  
NHS. *See* National Health Service  
Nigrini, Mark, 259  
NINos. *See* National Insurance  
Numbers  
Nobel Prize, 129, 141, 175, 180, 182, 183,  
213, 232  
non-ignorably missing. *See* missing  
nonresponse, 9, 46, 48, 53, 54, 55, 263,  
295  
Nosek, Brian, 185  
not data dependent, 227, 228- 237, 245,  
247, 253, 255, 263  
notifiable offences, 75  
nuclear fission bomb. *See* bomb:  
fission  
nuclear fusion, 173, 174  
Nudge Unit, 291–293  
Nuremberg Code, 60, 65, 67  
Obamacare health plan, 138  
obesity, 225, 292  
observational data, 59, 299  
Office for National Statistics, 73–75, 277  
Office for Standards in Education,  
Children's Services, and Skills, 120  
Office of Research Integrity, 209  
Ofsted. *See* Office for Standards in  
Education, Children's Services, and  
Skills  
oil drop experiment, 213, 214  
oil spill. *See* BP Deepwater Horizon  
omitted variable bias, 90  
online shopping, 79  
ONS. *See* Office for National Statistics  
opt in, 37, 295  
opt out, 36, 135, 136, 295  
O-rings, 17–21  
osteoarthritis, 180, 225  
osteoporosis, 90, 91, 92  
Pacioli, Luca, 162  
parameters, 26, 252, 273  
Pareto principle, 260  
Paris Exposition, 140  
passwords, 144, 145, 147, 152, 279, 280  
Pasteur, Louis, 200  
patterns of missing values, 238–239  
Pauling, Linus, 175, 176, 190  
peaking, 101  
penicillin, 183  
Penzias, Arno, 182  
people trafficking, 163  
periodic table, 169  
Personal Identification Number, 150  
personnel selection, 94  
Pfungst, Oskar, 302  
p-hacking, 191, 192, 196–198, 213, 295  
Pharmasset, 156  
Philips, Christopher, 39  
phishing, 152  
piling of numbers, 101  
Piltown Man, 206–208

- PIN. *See* Personal Identification Number
- Pinker, Stephen, 206
- placebo, 60, 61, 181, 264
- placement, in money laundering, 163, 164
- plagiarism, 142, 215
- plane crashes, 67
- plasma, 174
- plastic cards, 32, 148, 150
- platinobarium screen, 183
- Plewes, Thomas, 51, 55
- pneumonia, 130, 131
- Poisson, André, 140, 141
- police forces, 38, 75, 121
- Police Recorded Crime, 74
- poll, 48–51
- pollen count, 57
- pollution, 174, 266, 268
- Pons, Stanley, 174, 175
- Ponzi scheme, 164, 165
- Popperian, 167
- potholes, 9, 10, 15
- Potti, Anil, 209
- PRC. *See* Police Recorded Crime
- preclinical medical research, 185
- predictive models, 270
- predictor variable, 90
- prefrontal: lobe, 180; lobotomy, 179
- presidential commission, 17
- presidential election, 48
- principal agent problem, 116
- prior distribution, 276–278
- prison, 24, 63, 155, 160, 165
- protected characteristics, 132, 133
- provenance of data, 217–220
- pseudonymizing data, 36, 283, 286
- public policy, 32, 62, 112, 116, 284, 296
- publication bias, 187
- public-key cryptography, 281
- PubMed, 216
- Pulitzer, Joseph, 218
- p-value, 192–194, 197
- PwC, 163
- Quantas, 155
- question wording, 71
- radio telescope, 182
- radio-frequency identification, 150
- random assignment, 263–265
- randomized controlled trials, 22, 58–59, 263–265
- randomized response, 287, 288
- Rao, B. C. Subba, 182
- RCTs. *See* randomized controlled trials
- Referendum, UK, 49
- regression: coefficients, 102; multiple, 90; to the mean, 40, 95, 97, 188, 189, 195
- regulatory arbitrage, 116
- reidentification, 36, 283, 286
- reject inference, 240
- relativity, 168, 169, 177, 178
- reliability analysis, 242
- replication, 80, 175, 185, 188, 190, 195, 198, 201, 215, 269–276, 298
- reproducibility, 184, 186
- Reproducibility Project*, 185
- response rate, 51–54, 82, 293. *See also* nonresponse
- response variable, 90
- retail sales, 79
- retraction, 190, 210, 215–217
- retractionwatch.com, 216
- Reurink, Arjan, 130, 149, 150
- Rigobon, Roberto, 300
- risk factor, 52, 91
- road accidents, 145, 159, 161
- Roberts: John, 101; Paul, 160
- Robinson, W. S., 89

- rocket, 17, 18, 187  
Rogers, William, 19  
Röntgen, Wilhelm, 177, 183  
Roosevelt, Franklin D., 48–51, 55  
rounding, 98–101, 103, 297  
Rubin, Donald, 225–229, 254  
Rumsfeld, Donald, 11, 294
- S&P 500, 40, 233  
Saki, 230  
salmon, 194, 195  
Samsung Securities Co, 105, 106  
SAT examinations, 117  
Satyam, 157  
schizophrenia, 180  
Schlanger, Todd, 39  
Schön, Jan Hendrik, 209  
school grades, 117–119  
scientific revolution, 167  
scorecard, 33, 132–134  
screening, 52, 90–94, 239  
SDD. *See* seen data dependent  
SEC. *See* Securities and Exchange Commission  
secure multiparty computation, 288, 289  
Securities and Exchange Commission, 154  
seen data dependent, 227  
sensitive survey questions, 56  
series events, 51  
*Sesame Street*, 62  
sexual assaults, 38  
Shipman, Harold, 266  
Shmatikov, Vitaly, 287  
short-term migration, 74  
significance test, 192, 193, 194  
Simpson's paradox, 84, 89, 296  
simulation, 265–269, 276, 283, 298  
skewed distribution, 103, 212, 244, 297  
skimmers, 151  
Slutsky, Bob, 209  
smoking, 81, 137  
smurfing, 164  
social engineering, 151  
Sokal, Alan, 204, 205  
South Sea bubble, 125  
space program, 187  
space shuttle; *Challenger*, 17–22, 25, 295, 298; *Columbia*, 19  
spousal abuse, 51  
spouse's income, 239, 265  
Stapel, Diederik, 210  
statistical significance. *See* significance test  
steady state theory, 178  
Steen, R. Grant, 216  
Steiner, Peter, 144  
streptomycin, 59  
Subdø, Jon, 209  
subprime: applicants, 140; lending, 126; mortgages, 153  
Sumitomo, 153  
Survey of Consumer Attitudes, 51  
surveys, 43–56  
survival analysis, 242, 245, 254  
survival rates, 82–84  
survivor bias, 40, 41, 60, 233  
Swanson, Randel, 196  
Sweeney, Latanya, 285  
Swissair, 304  
Sylvester, Rachel, 119, 120  
systolic blood pressure. *See* blood pressure
- TARN. *See* Trauma Audit and Research Network  
tax: avoidance, 115; credits, 43; evasion, 166; filed, 146; fraud, 163; French, 124; laws, 115; office, 284; payments, 43, 44, 73; regime, 115, 116



- testability, 167, 171, 178  
*The Naked Surgeon*, 266  
*The Right Stuff*, 187  
*The Structure of Scientific Revolutions*, 183  
thermometer, 110  
Thomson, William, 172  
time series, 39, 80, 296  
*Titanic*, 82–84, 87–89  
Tokyo Stock Exchange, 106  
topcoding, 101, 110  
torturing data, 192, 194  
Tourangeau, Roger, 51, 55  
Toyota, 106  
transparency, 130, 162, 165, 218, 219, 220  
trauma, 6–8, 196  
Trauma Audit and Research Network,  
    6–10, 23  
trimming, 203, 211–213, 215  
Trump, Donald, 121  
truncation, 101  
trusted third party, 284  
tuberculosis, 59  
Tulip bubble, 125  
Twitter, 10, 11, 41  
Twyman's Law, 107  
types of dark data, 6, 11, 12, 163.  
    See also *DD-Type*
- UBS, 153  
UDD. See unseen data dependent  
Universal Declaration of Human  
    Rights, 279  
unseen data dependent, 227  
Uranus, 183
- vaccination, 4  
Van Helmont, Jean-Baptiste, 59  
Van Vlasselaer, Véronique, 142
- Vanguard, 39  
VHS, 304  
Visa card transactions, 32  
Vitamin C, 181  
von Helmholtz, Hermann, 173  
von Osten, Wilhelm, 302
- Wade, Nicholas, 214  
waterbed effect, 151  
Watkins, Sherron, 157  
Watson, James, 176  
Watson Research Center, IBM, 289  
web survey, 50, 55, 294  
web-scraping, 80, 301  
weight of American men, 102  
Weld, William, 285, 286  
West Yorkshire Metropolitan  
    Ambulance Service NHS Trust,  
        120  
wheat yields, 95  
whiplash injury, 160, 161  
Wiener, Norbert, 180  
Wikileaks, 130  
Willetts, David  
Wilson, Robert Woodrow, 182  
Winsorizing, 211  
Wolfe, Tom, 187  
Woodward, Arthur Smith, 206,  
    207  
WorldCom, 156, 157
- xkcd*, 16, 41  
X-ray, 90, 176, 177, 183
- yellow journalism, 217  
YouTube, 42  
Zoosk, 147