

WEAPONS OF MATH DESTRUCTION



HOW BIG DATA INCREASES INEQUALITY
AND THREATENS DEMOCRACY

CATHY O'NEIL

INTRODUCTION

When I was a little girl, I used to gaze at the traffic out the car window and study the numbers on license plates. I would reduce each one to its basic elements—the prime numbers that made it up. $45 = 3 \times 3 \times 5$. That's called factoring, and it was my favorite investigative pastime. As a budding math nerd, I was especially intrigued by the primes.

My love for math eventually became a passion. I went to math camp when I was fourteen and came home clutching a Rubik's Cube to my chest. Math provided a neat refuge from the messiness of the real world. It marched forward, its field of knowledge expanding relentlessly, proof by proof. And I could add to it. I majored in math in college and went on to get my PhD. My thesis was on algebraic number theory, a field with roots in all that factoring I did as a child. Eventually, I became a tenure-track professor at Barnard, which had a combined math department with Columbia University.

And then I made a big change. I quit my job and went to work as a quant for D. E. Shaw, a leading hedge fund. In leaving academia for finance, I carried mathematics from abstract theory into practice. The operations we performed on numbers translated into trillions of dollars sloshing from one account to another. At first I was excited and amazed by working in this new laboratory, the global economy. But in the autumn

of 2008, after I'd been there for a bit more than a year, it came crashing down.

The crash made it all too clear that mathematics, once my refuge, was not only deeply entangled in the world's problems but also fueling many of them. The housing crisis, the collapse of major financial institutions, the rise of unemployment—all had been aided and abetted by mathematicians wielding magic formulas. What's more, thanks to the extraordinary powers that I loved so much, math was able to combine with technology to multiply the chaos and misfortune, adding efficiency and scale to systems that I now recognized as flawed.

If we had been clear-headed, we all would have taken a step back at this point to figure out how math had been misused and how we could prevent a similar catastrophe in the future. But instead, in the wake of the crisis, new mathematical techniques were hotter than ever, and expanding into still more domains. They churned 24/7 through petabytes of information, much of it scraped from social media or e-commerce websites. And increasingly they focused not on the movements of global financial markets but on human beings, on us. Mathematicians and statisticians were studying our desires, movements, and spending power. They were predicting our trustworthiness and calculating our potential as students, workers, lovers, criminals.

This was the Big Data economy, and it promised spectacular gains. A computer program could speed through thousands of résumés or loan applications in a second or two and sort them into neat lists, with the most promising candidates on top. This not only saved time but also was marketed as fair and objective. After all, it didn't involve prejudiced humans digging through reams of paper, just machines processing cold numbers. By 2010 or so, mathematics was asserting itself as never before in human affairs, and the public largely welcomed it.

Yet I saw trouble. The math-powered applications powering the data economy were based on choices made by fallible human beings. Some of these choices were no doubt made with the best intentions. Nevertheless, many of these models encoded human prejudice, misunderstanding, and bias into the software systems that increasingly managed our lives. Like gods, these mathematical models were opaque, their workings invisible to all but the highest priests in their domain: mathematicians and computer

scientists. Their verdicts, even when wrong or harmful, were beyond dispute or appeal. And they tended to punish the poor and the oppressed in our society, while making the rich richer.

I came up with a name for these harmful kinds of models: Weapons of Math Destruction, or WMDs for short. I'll walk you through an example, pointing out its destructive characteristics along the way.

As often happens, this case started with a laudable goal. In 2007, Washington, D.C.'s new mayor, Adrian Fenty, was determined to turn around the city's underperforming schools. He had his work cut out for him: at the time, barely one out of every two high school students was surviving to graduation after ninth grade, and only 8 percent of eighth graders were performing at grade level in math. Fenty hired an education reformer named Michelle Rhee to fill a powerful new post, chancellor of Washington's schools.

The going theory was that the students weren't learning enough because their teachers weren't doing a good job. So in 2009, Rhee implemented a plan to weed out the low-performing teachers. This is the trend in troubled school districts around the country, and from a systems engineering perspective the thinking makes perfect sense: Evaluate the teachers. Get rid of the worst ones, and place the best ones where they can do the most good. In the language of data scientists, this "optimizes" the school system, presumably ensuring better results for the kids. Except for "bad" teachers, who could argue with that? Rhee developed a teacher assessment tool called IMPACT, and at the end of the 2009–10 school year the district fired all the teachers whose scores put them in the bottom 2 percent. At the end of the following year, another 5 percent, or 206 teachers, were booted out.

Sarah Wysocki, a fifth-grade teacher, didn't seem to have any reason to worry. She had been at MacFarland Middle School for only two years but was already getting excellent reviews from her principal and her students' parents. One evaluation praised her attentiveness to the children; another called her "one of the best teachers I've ever come into contact with."

Yet at the end of the 2010–11 school year, Wysocki received a miserable score on her IMPACT evaluation. Her problem was a new scoring system known as value-added modeling, which purported to measure her effectiveness in teaching math and language skills. That score, generated

by an algorithm, represented half of her overall evaluation, and it outweighed the positive reviews from school administrators and the community. This left the district with no choice but to fire her, along with 205 other teachers who had IMPACT scores below the minimal threshold.

This didn't seem to be a witch hunt or a settling of scores. Indeed, there's a logic to the school district's approach. Administrators, after all, could be friends with terrible teachers. They could admire their style or their apparent dedication. Bad teachers can *seem* good. So Washington, like many other school systems, would minimize this human bias and pay more attention to scores based on hard results: achievement scores in math and reading. The numbers would speak clearly, district officials promised. They would be more fair.

Wysocki, of course, felt the numbers were horribly unfair, and she wanted to know where they came from. "I don't think anyone understood them," she later told me. How could a good teacher get such dismal scores? What was the value-added model measuring?

Well, she learned, it was complicated. The district had hired a consultancy, Princeton-based Mathematica Policy Research, to come up with the evaluation system. Mathematica's challenge was to measure the educational progress of the students in the district and then to calculate how much of their advance or decline could be attributed to their teachers. This wasn't easy, of course. The researchers knew that many variables, from students' socioeconomic backgrounds to the effects of learning disabilities, could affect student outcomes. The algorithms had to make allowances for such differences, which was one reason they were so complex.

Indeed, attempting to reduce human behavior, performance, and potential to algorithms is no easy job. To understand what Mathematica was up against, picture a ten-year-old girl living in a poor neighborhood in southeastern Washington, D.C. At the end of one school year, she takes her fifth-grade standardized test. Then life goes on. She may have family issues or money problems. Maybe she's moving from one house to another or worried about an older brother who's in trouble with the law. Maybe she's unhappy about her weight or frightened by a bully at school. In any case, the following year she takes another standardized test, this

one designed for sixth graders.

If you compare the results of the tests, the scores should stay stable, or hopefully, jump up. But if her results sink, it's easy to calculate the gap between her performance and that of the successful students.

But how much of that gap is due to her teacher? It's hard to know, and Mathematica's models have only a few numbers to compare. At Big Data companies like Google, by contrast, researchers run constant tests and monitor thousands of variables. They can change the font on a single advertisement from blue to red, serve each version to ten million people, and keep track of which one gets more clicks. They use this feedback to hone their algorithms and fine-tune their operation. While I have plenty of issues with Google, which we'll get to, this type of testing is an effective use of statistics.

Attempting to calculate the impact that one person may have on another over the course of a school year is much more complex. "There are so many factors that go into learning and teaching that it would be very difficult to measure them all," Wysocki says. What's more, attempting to score a teacher's effectiveness by analyzing the test results of only twenty-five or thirty students is statistically unsound, even laughable. The numbers are far too small given all the things that could go wrong. Indeed, if we were to analyze teachers with the statistical rigor of a search engine, we'd have to test them on thousands or even millions of randomly selected students. Statisticians count on large numbers to balance out exceptions and anomalies. (And WMDs, as we'll see, often punish individuals who happen to *be* the exception.)

Equally important, statistical systems require feedback—something to tell them when they're off track. Statisticians use errors to train their models and make them smarter. If Amazon.com, through a faulty correlation, started recommending lawn care books to teenage girls, the clicks would plummet, and the algorithm would be tweaked until it got it right. Without feedback, however, a statistical engine can continue spinning out faulty and damaging analysis while never learning from its mistakes.

Many of the WMDs I'll be discussing in this book, including the Washington school district's value-added model, behave like that. They define their own reality and use it to justify their results. This type of

model is self-perpetuating, highly destructive—and very common.

When Mathematica's scoring system tags Sarah Wysocki and 205 other teachers as failures, the district fires them. But how does it ever learn if it was right? It doesn't. The system itself has determined that they were failures, and that is how they are viewed. Two hundred and six "bad" teachers are gone. That fact alone appears to demonstrate how effective the value-added model is. It is cleansing the district of underperforming teachers. Instead of searching for the truth, the score comes to embody it.

This is one example of a WMD feedback loop. We'll see many of them throughout this book. Employers, for example, are increasingly using credit scores to evaluate potential hires. Those who pay their bills promptly, the thinking goes, are more likely to show up to work on time and follow the rules. In fact, there are plenty of responsible people and good workers who suffer misfortune and see their credit scores fall. But the belief that bad credit correlates with bad job performance leaves those with low scores less likely to find work. Joblessness pushes them toward poverty, which further worsens their scores, making it even harder for them to land a job. It's a downward spiral. And employers never learn how many good employees they've missed out on by focusing on credit scores. In WMDs, many poisonous assumptions are camouflaged by math and go largely untested and unquestioned.

This underscores another common feature of WMDs. They tend to punish the poor. This is, in part, because they are engineered to evaluate large numbers of people. They specialize in bulk, and they're cheap. That's part of their appeal. The wealthy, by contrast, often benefit from personal input. A white-shoe law firm or an exclusive prep school will lean far more on recommendations and face-to-face interviews than will a fast-food chain or a cash-strapped urban school district. The privileged, we'll see time and again, are processed more by people, the masses by machines.

Wysocki's inability to find someone who could explain her appalling score, too, is telling. Verdicts from WMDs land like dictates from the algorithmic gods. The model itself is a black box, its contents a fiercely guarded corporate secret. This allows consultants like Mathematica to charge more, but it serves another purpose as well: if the people being evaluated are kept in the dark, the thinking goes, they'll be less likely to

attempt to game the system. Instead, they'll simply have to work hard, follow the rules, and pray that the model registers and appreciates their efforts. But if the details are hidden, it's also harder to question the score or to protest against it.

For years, Washington teachers complained about the arbitrary scores and clamored for details on what went into them. It's an algorithm, they were told. It's very complex. This discouraged many from pressing further. Many people, unfortunately, are intimidated by math. But a math teacher named Sarah Bax continued to push the district administrator, a former colleague named Jason Kamras, for details. After a back-and-forth that extended for months, Kamras told her to wait for an upcoming technical report. Bax responded: "How do you justify evaluating people by a measure for which you are unable to provide explanation?" But that's the nature of WMDs. The analysis is outsourced to coders and statisticians. And as a rule, they let the machines do the talking.

Even so, Sarah Wysocki was well aware that her students' standardized test scores counted heavily in the formula. And here she had some suspicions. Before starting what would be her final year at MacFarland Middle School, she had been pleased to see that her incoming fifth graders had scored surprisingly well on their year-end tests. At Barnard Elementary School, where many of Sarah's students came from, 29 percent of the students were ranked at an "advanced reading level." This was five times the average in the school district.

Yet when classes started she saw that many of her students struggled to read even simple sentences. Much later, investigations by the *Washington Post* and *USA Today* revealed a high level of erasures on the standardized tests at forty-one schools in the district, including Barnard. A high rate of corrected answers points to a greater likelihood of cheating. In some of the schools, as many as 70 percent of the classrooms were suspected.

What does this have to do with WMDs? A couple of things. First, teacher evaluation algorithms are a powerful tool for behavioral modification. That's their purpose, and in the Washington schools they featured both a stick and a carrot. Teachers knew that if their students stumbled on the test their own jobs were at risk. This gave teachers a strong motivation to ensure their students passed, especially as the Great

Recession battered the labor market. At the same time, if their students outperformed their peers, teachers and administrators could receive bonuses of up to \$8,000. If you add those powerful incentives to the evidence in the case—the high number of erasures and the abnormally high test scores—there were grounds for suspicion that fourth-grade teachers, bowing either to fear or to greed, had corrected their students' exams.

It is conceivable, then, that Sarah Wysocki's fifth-grade students started the school year with artificially inflated scores. If so, their results the following year would make it appear that they'd lost ground in fifth grade—and that their teacher was an underperformer. Wysocki was convinced that this was what had happened to her. That explanation would fit with the observations from parents, colleagues, and her principal that she was indeed a good teacher. It would clear up the confusion. Sarah Wysocki had a strong case to make.

But you cannot appeal to a WMD. That's part of their fearsome power. They do not listen. Nor do they bend. They're deaf not only to charm, threats, and cajoling but also to logic—even when there is good reason to question the data that feeds their conclusions. Yes, if it becomes clear that automated systems are screwing up on an embarrassing and systematic basis, programmers will go back in and tweak the algorithms. But for the most part, the programs deliver unflinching verdicts, and the human beings employing them can only shrug, as if to say, “Hey, what can you do?”

And that is precisely the response Sarah Wysocki finally got from the school district. Jason Kamras later told the *Washington Post* that the erasures were “suggestive” and that the numbers might have been wrong in her fifth-grade class. But the evidence was not conclusive. He said she had been treated fairly.

Do you see the paradox? An algorithm processes a slew of statistics and comes up with a probability that a certain person *might* be a bad hire, a risky borrower, a terrorist, or a miserable teacher. That probability is distilled into a score, which can turn someone's life upside down. And yet when the person fights back, “suggestive” countervailing evidence simply won't cut it. The case must be ironclad. The human victims of WMDs, we'll see time and again, are held to a far higher standard of evidence

than the algorithms themselves.

After the shock of her firing, Sarah Wysocki was out of a job for only a few days. She had plenty of people, including her principal, to vouch for her as a teacher, and she promptly landed a position at a school in an affluent district in northern Virginia. So thanks to a highly questionable model, a poor school lost a good teacher, and a rich school, which didn't fire people on the basis of their students' scores, gained one.

■ ■ ■

Following the housing crash, I woke up to the proliferation of WMDs in banking and to the danger they posed to our economy. In early 2011 I quit my job at the hedge fund. Later, after rebranding myself as a data scientist, I joined an e-commerce start-up. From that vantage point, I could see that legions of other WMDs were churning away in every conceivable industry, many of them exacerbating inequality and punishing the poor. They were at the heart of the raging data economy.

To spread the word about WMDs, I launched a blog, MathBabe. My goal was to mobilize fellow mathematicians against the use of sloppy statistics and biased models that created their own toxic feedback loops. Data specialists, in particular, were drawn to the blog, and they alerted me to the spread of WMDs in new domains. But in mid-2011, when Occupy Wall Street sprang to life in Lower Manhattan, I saw that we had work to do among the broader public. Thousands had gathered to demand economic justice and accountability. And yet when I heard interviews with the Occupiers, they often seemed ignorant of basic issues related to finance. They clearly hadn't been reading my blog. (I should add, though, that you don't need to understand all the details of a system to know that it has failed.)

I could either criticize them or join them, I realized, so I joined them. Soon I was facilitating weekly meetings of the Alternative Banking Group at Columbia University, where we discussed financial reform. Through this process, I came to see that my two ventures outside academia, one in finance, the other in data science, had provided me with fabulous access to the technology and culture powering WMDs.

Ill-conceived mathematical models now micromanage the economy, from advertising to prisons. These WMDs have many of the same

characteristics as the value-added model that derailed Sarah Wysocki's career in Washington's public schools. They're opaque, unquestioned, and unaccountable, and they operate at a scale to sort, target, or "optimize" millions of people. By confusing their findings with on-the-ground reality, most of them create pernicious WMD feedback loops.

But there's one important distinction between a school district's value-added model and, say, a WMD that scouts out prospects for extortionate payday loans. They have different payoffs. For the school district, the payoff is a kind of political currency, a sense that problems are being fixed. But for businesses it's just the standard currency: money. For many of the businesses running these rogue algorithms, the money pouring in seems to prove that their models are working. Look at it through their eyes and it makes sense. When they're building statistical systems to find customers or manipulate desperate borrowers, growing revenue appears to show that they're on the right track. The software is doing its job. The trouble is that profits end up serving as a stand-in, or proxy, for truth. We'll see this dangerous confusion crop up again and again.

This happens because data scientists all too often lose sight of the folks on the receiving end of the transaction. They certainly understand that a data-crunching program is bound to misinterpret people a certain percentage of the time, putting them in the wrong groups and denying them a job or a chance at their dream house. But as a rule, the people running the WMDs don't dwell on those errors. Their feedback is money, which is also their incentive. Their systems are engineered to gobble up more data and fine-tune their analytics so that more money will pour in. Investors, of course, feast on these returns and shower WMD companies with more money.

And the victims? Well, an internal data scientist might say, no statistical system can be *perfect*. Those folks are collateral damage. And often, like Sarah Wysocki, they are deemed unworthy and expendable. Forget about them for a minute, they might say, and focus on all the people who get helpful suggestions from recommendation engines or who find music they love on Pandora, the ideal job on LinkedIn, or perhaps the love of their life on Match.com. Think of the astounding scale, and ignore the imperfections.

Big Data has plenty of evangelists, but I'm not one of them. This book

will focus sharply in the other direction, on the damage inflicted by WMDs and the injustice they perpetuate. We will explore harmful examples that affect people at critical life moments: going to college, borrowing money, getting sentenced to prison, or finding and holding a job. All of these life domains are increasingly controlled by secret models wielding arbitrary punishments.

Welcome to the dark side of Big Data.

NOTES

INTRODUCTION

one out of every two: Robert Stillwell, *Public School Graduates and Dropouts from the Common Core of Data: School Year 2006–07*, NCES 2010-313 (Washington, DC: National Center for Education Statistics, Institute of Education Sciences, US Department of Education, 2009), 5, <http://nces.ed.gov/pubsearch/pubsinfo.asp?pubid=2010313>.

8 percent of eighth graders: Jihyun Lee, Wendy S. Grigg, and Gloria S. Dion, *The Nation's Report Card Mathematics 2007*, NCES 2007-494 (Washington, DC: National Center for Education Statistics, Institute of Education Sciences, US Department of Education, 2007), 32, <https://nces.ed.gov/nationsreportcard/pdf/main2007/2007494.pdf>.

Rhee developed a teacher assessment tool: Bill Turque, “Rhee Dismisses 241 D.C. Teachers; Union Vows to Contest Firings,” *Washington Post*, July 24, 2010, www.washingtonpost.com/wp-dyn/content/article/2010/07/23/AR2010072303093.html.

the district fired all the teachers: Steven Sawchuck, “Rhee to Dismiss Hundreds of Teachers for Poor Performance,” *Education Week Blog*, July 23, 2010, http://blogs.edweek.org/edweek/teacherbeat/2010/07/_states_and_districts_across.html.

another 5 percent, or 205 teachers: Bill Turque, “206 Low-Performing D.C. Teachers Fired,” *Washington Post*, July 15, 2011, www.washingtonpost.com/local/education/206-low-performing-dc-teachers-fired/2011/07/15/gIQANEj5GI_story.html.

Sarah Wysocki, a fifth-grade teacher: Bill Turque, “‘Creative...Motivating’ and Fired,” *Washington Post*, March 6, 2012, www.washingtonpost.com/local/education/creative-motivating-and-fired/2012/02/04/gIQAwzZpvR_story.html.

One evaluation praised her: Ibid.

Wysocki received a miserable score: Ibid.

represented half of her overall evaluation: Ibid.

The district had hired a consultancy: Ibid.

“There are so many factors”: Sarah Wysocki, e-mail interview by author, August 6, 2015.

a math teacher named Sarah Bax: Guy Brandenburg, “DCPS Administrators Won’t or Can’t Give a DCPS Teacher the IMPACT Value-Added Algorithm,” *GFBrandenburg’s Blog*, February 27, 2011, <https://gfbrandenburg.wordpress.com/2011/02/27/dcps-administrators-wont-or-cant-give-a-dcps-teacher-the-impact-value-added-algorithm/>.

29 percent of the students: Turque, “‘Creative...Motivating’ and Fired.”

USA Today revealed a high level: Jack Gillum and Marisol Bello, “When Standardized Test Scores Soared in D.C., Were the Gains Real?,” *USA Today*, March 30, 2011, http://usatoday30.usatoday.com/news/education/2011-03-28-1Aschooltesting28_CV_N.htm.

bonuses of up to \$8,000: Ibid.

the erasures were “suggestive”: Turque, “‘Creative...Motivating’ and Fired.”

Sarah Wysocki was out of a job: Ibid.

CHAPTER 1

Boudreau, perhaps out of desperation: David Waldstein, “Who’s on Third? In Baseball’s Shifting Defenses, Maybe Nobody,” *New York Times*, May 12, 2014, www.nytimes.com/2014/05/13/sports/baseball/whos-on-third-in-baseballs-shifting-defenses-maybe-nobody.html?_r=0.

Moneyball: Michael Lewis, *Moneyball: The Art of Winning an Unfair Game* (New York: W. W. Norton, 2003).

In 1997, a convicted murderer: Manny Fernandez, “Texas Execution Stayed Based on Race Testimony,” *New York Times*, September 16, 2011, www.nytimes.com/2011/09/17/us/experts-testimony-on-race-led-to-stay-of-execution-in-texas.html?pagewanted=all.

made a reference to Buck’s race: Ibid.

“It is inappropriate to allow race”: Alan Berlow, “See No Racism, Hear No Racism: Despite Evidence, Perry About to Execute Another Texas Man,” *National Memo*, September 15, 2011, www.nationalmemo.com/perry-might-let-another-man-die/.

Buck never got a new hearing: NAACP Legal Defense Fund, “Texas Fifth Circuit Rejects Appeal in Case of Duane Buck,” NAACP LDF website, August 21, 2015, www.naacpldf.org/update/texas-fifth-circuit-rejects-appeal-case-duane-buck.

prosecutors were three times more likely: OpenFile, “TX: Study Finds Harris County Prosecutors Sought Death Penalty 3-4 Times More Often Against Defendants of Color,” *Open File, Prosecutorial Misconduct and Accountability*, March 15, 2013, www.prosecutorialaccountability.com/2013/03/15/tx-study-finds-harris-county-prosecutors-sought-death-penalty-3-4-times-more-often-against-defendants-of-color/.

sentences imposed on black men: American Civil Liberties Union, *Racial Disparities in Sentencing, Hearing on Reports of Racism in the Justice System of the United States*, submitted to the Inter-American Commission on Human Rights, 153rd Session, October 27, 2014, www.aclu.org/sites/default/files/assets/141027_iachr_racial_disparities_aclu_submission_o.pdf.

blacks fill up 40 percent of America’s prison cells: Federal Bureau of Prisons, Statistics web page, accessed January 8, 2016, www.bop.gov/about/statistics/statistics_inmate_race.jsp.

courts in twenty-four states: Sonja Starr, “Sentencing, by the Numbers,” *New York Times*, August 10, 2014, www.nytimes.com/2014/08/11/opinion/sentencing-by-the-numbers.html.

average of \$31,000 a year: Christian Henrichson and Ruth Delaney, *The Price of Prisons: What Incarceration Costs Taxpayers* (New York: VERA Institute of Justice, 2012), [www](http://www.vera-institute.org/).