

REVISED AND UPDATED

PREDICTIVE ANALYTICS

"Mesmerizing & fascinating..."

—*The Seattle Post-Intelligencer*

AN INTRODUCTION
FOR EVERYONE



THE POWER TO PREDICT WHO WILL
CLICK, BUY, LIE, OR DIE

ERIC SIEGEL

WILEY

PREDICTIVE ANALYTICS



**THE POWER TO PREDICT WHO WILL
CLICK, BUY, LIE, OR DIE**

ERIC SIEGEL

WILEY

Cover image: Winona Nelson
Cover design: Wiley
Interior image design: Matt Kornhaas

Copyright © 2016 by Eric Siegel. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.
Published simultaneously in Canada.

Jeopardy!® is a registered trademark of Jeopardy Productions, Inc.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600, or on the Web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at www.wiley.com/go/permissions.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with the respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor the author shall be liable for damages arising herefrom.

For general information about our other products and services, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley publishes in a variety of print and electronic formats and by print-on-demand. Some material included with standard print versions of this book may not be included in e-books or in print-on-demand. If this book refers to media such as a CD or DVD that is not included in the version you purchased, you may download this material at <http://booksupport.wiley.com>. For more information about Wiley products, visit www.wiley.com.

Library of Congress Cataloging-in-Publication Data:

Names: Siegel, Eric, 1968-

Title: Predictive analytics : the power to predict who will click, buy, lie, or die / Eric Siegel.

Description: Revised and Updated Edition. | Hoboken : Wiley, 2016. | Revised edition of the author's *Predictive analytics*, 2013. | Includes index.

Identifiers: LCCN 2015031895 (print) | LCCN 2015039877 (ebook) | ISBN 9781119145677 (paperback) | ISBN 9781119145684 (pdf) | ISBN 9781119153658 (epub)

Subjects: LCSH: Social sciences—Forecasting. | Economic forecasting | Prediction (Psychology) | Social prediction. | Human behavior. | BISAC: BUSINESS & ECONOMICS / Consumer Behavior. | BUSINESS & ECONOMICS / Econometrics. | BUSINESS & ECONOMICS / Marketing / General.

Classification: LCC H61.4 .S54 2016 (print) | LCC H61.4 (ebook) | DDC 303.49—dc23

LC record available at <http://lcn.loc.gov/2015031895>

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

Foreword

This book deals with quantitative efforts to predict human behavior. One of the earliest efforts to do that was in World War II. Norbert Wiener, the father of “cybernetics,” began trying to predict the behavior of German airplane pilots in 1940—with the goal of shooting them from the sky. His method was to take as input the trajectory of the plane from its observed motion, consider the pilot’s most likely evasive maneuvers, and predict where the plane would be in the near future so that a fired shell could hit it. Unfortunately, Wiener could predict only one second ahead of a plane’s motion, but 20 seconds of future trajectory were necessary to shoot down a plane.

In Eric Siegel’s book, however, you will learn about a large number of prediction efforts that are much more successful. Computers have gotten a lot faster since Wiener’s day, and we have a lot more data. As a result, banks, retailers, political campaigns, doctors and hospitals, and many more organizations have been quite successful of late at predicting the behavior of particular humans. Their efforts have been helpful at winning customers, elections, and battles with disease.

My view—and Siegel’s, I would guess—is that this predictive activity has generally been good for humankind. In the context of healthcare, crime, and terrorism, it can save lives. In the context of advertising, using predictions is more efficient and could conceivably save both trees (for direct mail and catalogs) and the time and attention of the recipient. In politics, it seems to reward those candidates who respect the scientific method (some might disagree, but I see that as a positive).

However, as Siegel points out—early in the book, which is admirable—these approaches can also be used in somewhat harmful ways. “With great power comes great responsibility,” he notes in quoting *Spider-Man*. The implication is that we must be careful as a society about how we use predictive models, or we may be restricted from using and benefiting from them. Like other powerful technologies or disruptive human innovations, predictive analytics is essentially amoral and can be used for good or evil. To avoid the evil applications, however, it is certainly important to understand what is possible with predictive analytics, and you will certainly learn that if you keep reading.

This book is focused on predictive analytics, which is not the only type of analytics, but the most interesting and important type. I don’t think we need more books anyway on purely descriptive analytics, which only describe the past and don’t provide any insight as to why it happened. I also often refer in my own writing to a third type of analytics—“prescriptive”—that tells its users what to do through controlled experiments or optimization. Those quantitative methods are much less popular, however, than predictive analytics.

This book and the ideas behind it are a good counterpoint to the work of Nassim Nicholas Taleb. His books, including *The Black Swan*, suggest that many efforts at prediction are doomed to fail because of randomness and the inherent unpredictability of complex events. Taleb is no doubt correct that some events are black swans that are beyond prediction, but the fact is that most human behavior is quite regular and predictable. The many examples that Siegel provides of successful prediction remind us that most swans are white.

Siegel also resists the blandishments of the “big data” movement. Certainly some of the examples he mentions fall into this category—data that is too large or unstructured to be easily managed by conventional relational databases. But the point of predictive analytics is not the relative size or unruliness of your data, but what you do with it. I have found that “big data often equals small math,” and many big data practitioners are content just to use their data to create some appealing visual analytics. That’s not nearly as valuable as creating a predictive model.

Siegel has fashioned a book that is both sophisticated and fully accessible to the non-quantitative reader. It's got great stories, great illustrations, and an entertaining tone. Such non-quants should definitely read this book, because there is little doubt that their behavior will be analyzed and predicted throughout their lives. It's also quite likely that most non-quants will increasingly have to consider, evaluate, and act on predictive models at work.

In short, we live in a predictive society. The best way to prosper in it is to understand the objectives, techniques, and limits of predictive models. And the best way to do that is simply to keep reading this book.

—**Thomas H. Davenport**

Thomas H. Davenport is the President's
Distinguished Professor at Babson College,
a fellow of the MIT Center for Digital Business,
Senior Advisor to Deloitte Analytics,
and cofounder of the International Institute for Analytics.

He is the coauthor of *Competing on Analytics*,
Big Data @ Work, and several other books on analytics.

Preface to the Original Edition

Yesterday is history, tomorrow is a mystery, but today is a gift. That's why we call it the present.

—Attributed to A. A. Milne, Bil Keane, and Oogway,
the wise turtle in *Kung Fu Panda*

People look at me funny when I tell them what I do. It's an occupational hazard.

The Information Age suffers from a glaring omission. This claim may surprise many, considering we are actively recording Everything That Happens in the World. Moving beyond history books that document important events, we've progressed to systems that log every click, payment, call, crash, crime, and illness. With this in place, you would expect lovers of data to be satisfied, if not spoiled rotten.

But this apparent infinity of information excludes the very events that would be most valuable to know of: *things that haven't happened yet*.

Everyone craves the power to see the future; we are collectively obsessed with prediction. We bow to prognostic deities. We empty our pockets for palm readers. We hearken to horoscopes, adore astrology, and feast upon fortune cookies.

But many people who salivate for psychics also spurn science. Their innate response says “yuck”—it's either too hard to understand or too boring. Or perhaps many believe prediction by its nature is just impossible without supernatural support.

There's a lighthearted TV show I like premised on this very theme, *Psych*, in which a sharp-eyed detective—a modern-day, data-driven Sherlock Holmesian hipster—has perfected the art of observation so masterfully, the cops believe his spot-on deductions must be an admission of guilt. The hero gets out of this pickle by conforming to the norm: He simply informs the police he is psychic, thereby managing to stay out of prison and continuing to fight crime. Comedy ensues.

I've experienced the same impulse, for example, when receiving the occasional friendly inquiry as to my astrological sign. But, instead of posing as a believer, I turn to humor: "I'm a Scorpio, and Scorpios don't believe in astrology."

The more common cocktail party interview asks what I do for a living. I brace myself for eyes glazing over as I carefully enunciate: *predictive analytics*. Most people have the luxury of describing their job in a single word: doctor, lawyer, waiter, accountant, or actor. But, for me, describing this largely unknown field hijacks the conversation every time. Any attempt to be succinct falls flat:

I'm a business consultant in technology. They aren't satisfied and ask, "What kind of technology?"

I make computers predict what people will do. Bewilderment results, accompanied by complete disbelief and a little fear.

I make computers learn from data to predict individual human behavior. Bewilderment, plus nobody wants to talk about data at a party.

I analyze data to find patterns. Eyes glaze over even more; awkward pauses sink amid a sea of abstraction.

I help marketers target which customers will buy or cancel. They sort of get it, but this wildly undersells and pigeonholes the field.

I predict customer behavior, like when Target famously predicted whether you are pregnant. Moonwalking ensues.

So I wrote this book to demonstrate for you why predictive analytics is intuitive, powerful, and awe-inspiring.

I have good news: *A little prediction goes a long way.* I call this The Prediction Effect, a theme that runs throughout the book. The potency of prediction is

pronounced—as long as the predictions are better than guessing. This effect renders predictive analytics believable. We don't have to do the impossible and attain true clairvoyance. The story is exciting yet credible: Putting odds on the future to lift the fog just a bit off our hazy view of tomorrow means pay dirt. In this way, predictive analytics combats risk, boosts sales, cuts costs, fortifies healthcare, streamlines manufacturing, conquers spam, toughens crime fighting, optimizes social networks, and wins elections.

Do you have the heart of a scientist or a businessperson? Do you feel more excited by the very idea of prediction, or by the value it holds for the world?

I was struck by the notion of *knowing the unknowable*. Prediction seems to defy a law of nature: You cannot see the future because it isn't here yet. We find a workaround by building machines that learn from experience. It's the regimented discipline of using what we *do* know—in the form of data—to place increasingly accurate odds on what's coming next. We blend the best of math and technology, systematically tweaking until our scientific hearts are content to derive a system that peers right through the previously impenetrable barrier between today and tomorrow.

Talk about boldly going where no one has gone before!

Some people are in sales; others are in politics. I'm in prediction, and it's awesome.

Introduction

The Prediction Effect

I'm just like you. I succeed at times, and at others I fail. Some days good things happen to me, some days bad. We always wonder how things could have gone differently. I begin with seven brief tales of woe:

1. In 2009 I just about destroyed my right knee downhill skiing in Utah. The jump was no problem; it was landing that presented an issue. For knee surgery, I had to pick a graft source from which to reconstruct my busted ACL (the knee's central ligament). The choice is a tough one and can make the difference between living with a good knee or a bad knee. I went with my hamstring. *Could the hospital have selected a medically better option for my case?*
2. Despite all my suffering, it was really my health insurance company that paid dearly—knee surgery is expensive. *Could the company have better anticipated the risk of accepting a ski jumping fool as a customer and priced my insurance premium accordingly?*
3. Back in 1995 another incident caused me suffering, although it hurt less. I fell victim to identity theft, costing me dozens of hours of bureaucratic baloney and tedious paperwork to clear up my damaged credit rating. *Could the creditors have prevented the fiasco by detecting*

that the accounts were bogus when they were filed under my name in the first place?

4. With my name cleared, I recently took out a mortgage to buy an apartment. Was it a good move, or *should my financial adviser have warned me the property could soon be outvalued by my mortgage?*
5. While embarking on vacation, I asked the neighboring airplane passenger what price she'd paid for her ticket, and it was much less than I'd paid. *Before I booked the flight, could I have determined the airfare was going to drop?*
6. My professional life is susceptible, too. My business is faring well, but a company always faces the risk of changing economic conditions and growing competition. *Could we protect the bottom line by foreseeing which marketing activities and other investments will pay off, and which will amount to burnt capital?*
7. Small ups and downs determine your fate and mine, every day. A precise spam filter has a meaningful impact on almost every working hour. We depend heavily on effective Internet search for work, health (e.g., exploring knee surgery options), home improvement, and most everything else. We put our faith in personalized music and movie recommendations from Spotify and Netflix. After all these years, my mailbox wonders why companies don't know me well enough to send less junk mail (and sacrifice fewer trees needlessly).

These predicaments matter. They can make or break your day, year, or life. But what do they all have in common?

These challenges—and many others like them—are best addressed with *prediction*. Will the patient's outcome from surgery be positive? Will the credit applicant turn out to be a fraudster? Will the homeowner face a bad mortgage? Will the airfare go down? Will the customer respond if mailed a brochure? By predicting these things, it is possible to fortify healthcare, combat risk, conquer spam, toughen crime fighting, boost sales, and cut costs.

PREDICTION IN BIG BUSINESS—THE DESTINY OF ASSETS

There's another angle. Beyond benefiting you and me as consumers, prediction serves the organization, empowering it with an entirely new form of competitive armament. Corporations positively pounce on prediction.

In the mid-1990s, an entrepreneurial scientist named Dan Steinberg delivered predictive capabilities unto the nation's largest bank, Chase, to assist with their management of millions of mortgages. This mammoth enterprise put its faith in Dan's predictive technology, deploying it to drive transactional decisions across a tremendous mortgage portfolio. What did this guy have on his résumé?

Prediction is power. Big business secures a killer competitive stronghold by predicting the future destiny and value of individual assets. In this case, by driving mortgage decisions with predictions about the future payment behavior of homeowners, Chase curtailed risk, boosted profit, and witnessed a windfall.

INTRODUCING . . . THE CLAIRVOYANT COMPUTER

Compelled to grow and propelled to the mainstream, predictive technology is commonplace and affects everyone, every day. It impacts your experiences in undetectable ways as you drive, shop, study, vote, see the doctor, communicate, watch TV, earn, borrow, or even steal.

This book is about the most influential and valuable achievements of computerized prediction, and the two things that make it possible: the people behind it, and the fascinating science that powers it.

Making such predictions poses a tough challenge. Each prediction depends on multiple factors: The various characteristics known about each patient, each homeowner, each consumer, and each e-mail that may be spam. How shall we attack the intricate problem of putting all these pieces together for each prediction?

The idea is simple, although that doesn't make it easy. The challenge is tackled by a systematic, scientific means to develop and continually improve prediction—to literally *learn* to predict.

The solution is *machine learning*—computers automatically developing new knowledge and capabilities by furiously feeding on modern society's greatest and most potent *unnatural* resource: data.

“FEED ME!”—FOOD FOR THOUGHT FOR THE MACHINE

Data is the new oil.

—European Consumer Commissioner Meglena Kuneva

The only source of knowledge is experience.

—Albert Einstein

In God we trust. All others must bring data.

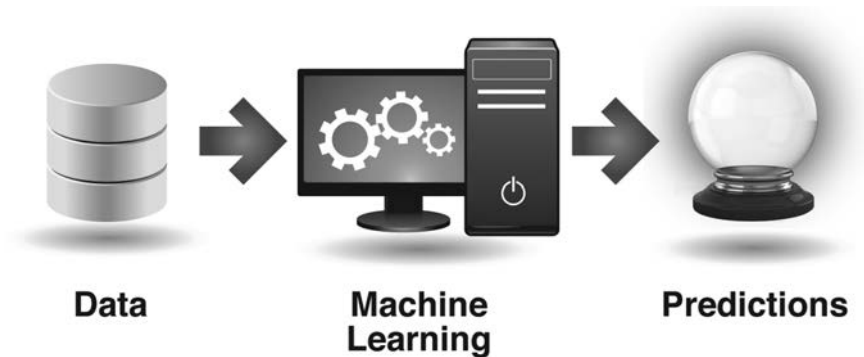
—William Edwards Deming (a business professor famous for work in manufacturing)

Most people couldn't be less interested in data. It can seem like such dry, boring stuff. It's a vast, endless regimen of recorded facts and figures, each alone as mundane as the most banal tweet, “I just bought some new sneakers!” It's the unsalted, flavorless residue deposited en masse as businesses churn away.

Don't be fooled! The truth is that data embodies a priceless collection of experience from which to learn. Every medical procedure, credit application, Facebook post, movie recommendation, fraudulent act, spammy e-mail, and purchase of any kind—each positive or negative outcome, each successful or failed sales call, each incident, event, and transaction—is encoded as data and warehoused. This glut grows by an estimated 2.5 quintillion bytes per day (that's a 1 with 18 zeros after it). And so a veritable Big Bang has set off, delivering an epic sea of raw materials, a plethora of examples so great in number, only a computer could manage to learn from them. Used correctly, computers avidly soak up this ocean like a sponge.

As data piles up, we have ourselves a genuine gold rush. But data isn't the gold. I repeat, data in its raw form is boring crud. The gold is what's discovered therein.

The process of machines learning from data unleashes the power of this exploding resource. It uncovers what drives people and the actions they take—what makes us tick and how the world works. With the new knowledge gained, prediction is possible.



This learning process discovers insightful gems such as:¹

- Early retirement decreases your life expectancy.
- Online daters more consistently rated as attractive receive *less* interest.
- Rihanna fans are mostly political Democrats.
- Vegetarians miss fewer flights.
- Local crime increases after public sporting events.

Machine learning builds upon insights such as these in order to develop predictive capabilities, following a number-crunching, trial-and-error process that has its roots in statistics and computer science.

¹ See Chapter 3 for more details on these examples.

I KNEW YOU WERE GOING TO DO THAT

With this power at hand, what do we want to predict? Every important thing a person does is valuable to predict, namely: *consume, think, work, quit, vote, love, procreate, divorce, mess up, lie, cheat, steal, kill, and die*. Let's explore some examples.²

PEOPLE CONSUME

- Hollywood studios predict the success of a screenplay if produced.
- Netflix awarded \$1 million to a team of scientists who best improved their recommendation system's ability to predict which movies you will like.
- The Hopper app helps you get the best deal on a flight by recommending whether you should buy or wait, based on its prediction as to whether the airfare will change.
- Australian energy company Energex predicts electricity demand in order to decide where to build out its power grid, and Con Edison predicts system failure in the face of high levels of consumption.
- Wall Street firms trade algorithmically, buying and selling based on the prediction of stock prices.
- Companies predict which customer will buy their products in order to target their marketing, from U.S. Bank down to small companies like Harbor Sweets (candy) and Vermont Country Store ("top quality and hard-to-find classic products"). These predictions dictate the allocations of precious marketing budgets. Some companies literally predict how to best influence you to buy more (the topic of Chapter 7).
- Prediction drives the coupons you get at the grocery cash register. U.K. grocery giant Tesco, the world's third-largest retailer, predicts which discounts will be redeemed in order to target more than

² For more examples and further detail, see this book's Central Tables.

100 million personalized coupons annually at cash registers across 13 countries. Similarly, Kmart, Kroger, Ralph's, Safeway, Stop & Shop, Target, and Winn-Dixie follow in kind.

- Predicting mouse clicks pays off massively. Since websites are often paid per click for the advertisements they display, they predict which ad you're mostly likely to click in order to instantly choose which one to show you. This, in effect, selects more relevant ads and drives millions in newly found revenue.
- Facebook predicts which of the thousands of posts by your friends will interest you most every time you view the news feed (unless you change the default setting). The social network also predicts the suggested "people you may know," not to mention which ads you're likely to click.

PEOPLE LOVE, WORK, PROCREATE, AND DIVORCE

- The leading career-focused social network, LinkedIn, predicts your job skills.
- Online dating leaders Match.com, OkCupid, and eHarmony predict which hottie on your screen would be the best bet at your side.
- Target predicts customer pregnancy in order to market relevant products accordingly. Nothing foretells consumer need like predicting the birth of a new consumer.
- Clinical researchers predict infidelity and divorce. There's even a self-help website tool to put odds on your marriage's long-term success (www.divorceprobability.com).

PEOPLE THINK AND DECIDE

- Obama was reelected in 2012 with the help of voter prediction. The Obama for America campaign predicted which voters would be positively persuaded by campaign contact (a call, door knock, flier, or TV ad), and which would actually be inadvertently influenced to

(continued)

(continued)

vote adversely by contact. Employed to drive campaign decisions for millions of swing state voters, this method was shown to successfully convince more voters to choose Obama than traditional campaign targeting. Hillary for America 2016 is positioning to apply the same technique.

- “What did you mean by that?” Systems have learned to ascertain the intent behind the written word. Citibank and PayPal detect the customer sentiment about their products, and one researcher’s machine can tell which Amazon.com book reviews are sarcastic.
- Student essay grade prediction has been developed for possible use to automatically grade. The system grades as accurately as human graders.
- There’s a machine that can participate in the same capacity as humans in the United States’ most popular broadcast celebration of human knowledge and cultural literacy. On the TV quiz show *Jeopardy!*, IBM’s Watson computer triumphed. This machine learned to work proficiently enough with English to predict the answers to free-form inquiries across an open range of topics and defeat the two all-time human champs.
- Computers can literally read your mind. Researchers trained systems to decode a scan of your brain and determine which type of object you’re thinking about—such as certain tools, buildings, and food—with over 80 percent accuracy for some human subjects.

PEOPLE QUIT

- Hewlett-Packard (HP) earmarks each and every one of its more than 300,000 worldwide employees according to “Flight Risk,” the expected chance he or she will quit their job, so that managers may intervene in advance where possible and plan accordingly otherwise.
- Ever experience frustration with your cell phone service? Your service provider endeavors to know. All major wireless carriers

predict how likely it is you will cancel and switch to a competitor—possibly before you have even conceived a plan to do so—based on factors such as dropped calls, your phone usage, billing information, and whether your contacts have already defected.

- FedEx stays ahead of the game by predicting—with 65 to 90 percent accuracy—which customers are at risk of defecting to a competitor.
- The American Public University System predicted student dropouts and used these predictions to intervene successfully; the University of Alabama, Arizona State University, Iowa State University, Oklahoma State University, and the Netherlands' Eindhoven University of Technology predict dropouts as well.
- Wikipedia predicts which of its editors, who work for free as a labor of love to keep this priceless online asset alive, are going to discontinue their valuable service.
- Researchers at Harvard Medical School predict that if your friends stop smoking, you're more likely to do so yourself as well. Quitting smoking is contagious.

PEOPLE MESS UP

- Insurance companies predict who is going to crash a car or hurt themselves another way (such as a ski accident). Allstate predicts bodily injury liability from car crashes based on the characteristics of the insured vehicle, demonstrating improvements to prediction that could be worth an estimated \$40 million annually. Another top insurance provider reported savings of almost \$50 million per year by expanding its actuarial practices with advanced predictive techniques.
- Ford is learning from data so its cars can detect when the driver is not alert due to distraction, fatigue, or intoxication and take action such as sounding an alarm.
- Researchers have identified aviation incidents that are five times more likely than average to be fatal, using data from the National Transportation Safety Board.

(continued)

(continued)

- All large banks and credit card companies predict which debtors are most likely to turn delinquent, failing to pay back their loans or credit card balances. Collection agencies prioritize their efforts with predictions of which tactic has the best chance to recoup the most from each defaulting debtor.

PEOPLE GET SICK AND DIE

I'm not afraid of death; I just don't want to be there when it happens.

—Woody Allen

- In 2013, the Heritage Provider Network handed over \$500,000 to a team of scientists who won an analytics competition to best predict individual hospital admissions. By following these predictions, proactive preventive measures can take a healthier bite out of the tens of billions of dollars spent annually on unnecessary hospitalizations. Similarly, the University of Pittsburgh Medical Center predicts short-term hospital readmissions, so doctors can be prompted to think twice before a hasty discharge.
- At Stanford University, a machine learned to diagnose breast cancer better than human doctors by discovering an innovative method that considers a greater number of factors in a tissue sample.
- Researchers at Brigham Young University and the University of Utah correctly predict about 80 percent of premature births (and about 80 percent of full-term births), based on peptide biomarkers, as found in a blood exam as early as week 24 of pregnancy.
- University researchers derived a method to detect patient schizophrenia from transcripts of their spoken words alone.
- A growing number of life insurance companies go beyond conventional actuarial tables and employ predictive technology to establish mortality risk. It's not called *death insurance*, but they calculate when you are going to die.

- Beyond life insurance, one top-five *health* insurance company predicts the probability that elderly insurance policyholders will pass away within 18 months, based on clinical markers in the insured's recent medical claims. Fear not—it's actually done for benevolent purposes.
- Researchers predict your risk of death in surgery based on aspects of you and your condition to help inform medical decisions.
- By following one common practice, doctors regularly—yet unintentionally—sacrifice some patients in order to save others, and this is done completely without controversy. But this would be lessened by predicting something besides diagnosis or outcome: healthcare *impact* (impact prediction is the topic of Chapter 7).

PEOPLE LIE, CHEAT, STEAL, AND KILL

- Most medium-size and large banks employ predictive technology to counter the ever-blooming assault of fraudulent checks, credit card charges, and other transactions. Citizens Bank developed the capacity to decrease losses resulting from check fraud by 20 percent. Hewlett-Packard saved \$66 million by detecting fraudulent warranty claims.
- Predictive computers help decide who belongs in prison. To assist with parole and sentencing decisions, officials in states such as Oregon and Pennsylvania consult prognostic machines that assess the risk a convict will offend again.
- Murder is widely considered impossible to predict with meaningful accuracy in general, but within at-risk populations predictive methods can be effective. Maryland analytically generates predictions as to which inmates will kill or be killed. University and law enforcement researchers have developed predictive systems that foretell murder among those previously convicted for homicide.
- One fraud expert at a large bank in the United Kingdom extended his work to discover a small pool of terror suspects based on their

(continued)

(continued)

banking activities. While few details have been disclosed publicly, it's clear that the National Security Agency also considers this type of analysis a strategic priority in order to automatically discover previously unknown potential suspects.

- Police patrol the areas predicted to spring up as crime hot spots in cities such as Chicago, Memphis, and Richmond, Va.
- Inspired by the TV crime drama *Lie to Me* about a microexpression reader, researchers at the University at Buffalo trained a system to detect lies with 82 percent accuracy by observing eye movements alone.
- As a professor at Columbia University in the late 1990s, I had a team of teaching assistants who employed cheating-detection software to patrol hundreds of computer programming homework submissions for plagiarism.
- The IRS predicts if you are cheating on your taxes.

THE LIMITS AND POTENTIAL OF PREDICTION

An economist is an expert who will know tomorrow why the things he predicted yesterday didn't happen.

—Earl Wilson

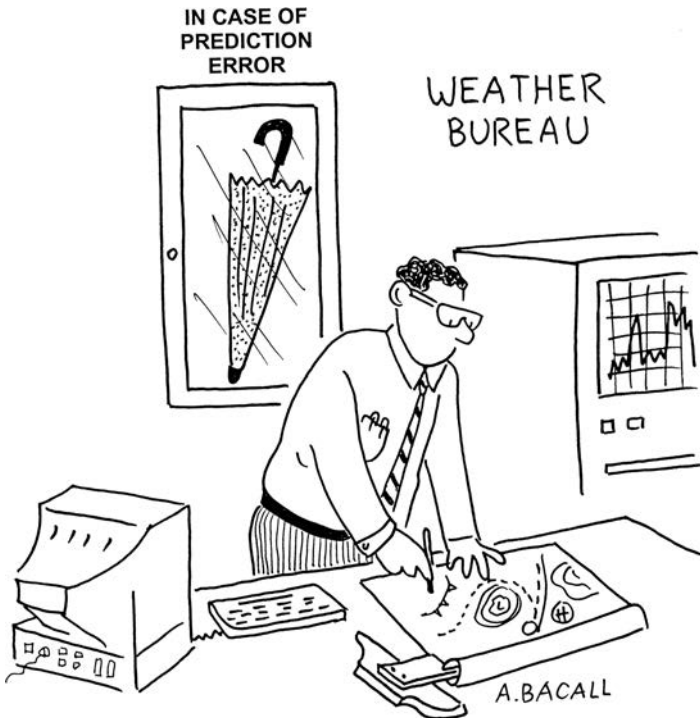
How come you never see a headline like "Psychic Wins Lottery"?

—Jay Leno

Each of the preceding accomplishments is powered by prediction, which is in turn a product of machine learning. A striking difference exists between these varied capabilities and science fiction: They aren't fiction. At this point, I predict that you won't be surprised to hear that those examples represent

only a small sample. You can safely predict that the power of prediction is here to stay.

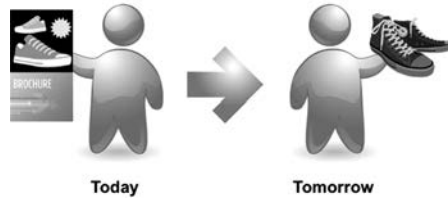
But are these claims too bold? As the Danish physicist Niels Bohr put it, “Prediction is very difficult, especially if it’s about the future.” After all, isn’t prediction basically impossible? The future is unknown, and uncertainty is the only thing about which we’re certain.



Let me be perfectly clear. It’s fuzzy. Accurate prediction is generally not possible. The weather is predicted with only about 50 percent accuracy, and it doesn’t get easier predicting the behavior of humans, be they patients, customers, or criminals.

Good news! Predictions need not be accurate to score big value. For instance, one of the most straightforward commercial applications of

predictive technology is deciding whom to target when a company sends direct mail. If the learning process identifies a carefully defined group of customers who are predicted to be, say, three times more likely than average to respond positively to the mail, the company profits big-time by preemptively removing likely *nonresponders* from the mailing list. And those non-responders in turn benefit, contending with less junk mail.



Prediction—A person who sees a sales brochure today buys a product tomorrow.

In this way the business, already playing a sort of numbers game by conducting mass marketing in the first place, tips the balance delicately yet significantly in its favor—and does so without highly accurate predictions. In fact, its utility withstands quite poor accuracy. If the overall marketing response is at 1 percent, the so-called hot pocket with three times as many would-be responders is at 3 percent. So, in this case, we can't confidently predict the response of any one particular customer. Rather, the value is derived from identifying a group of people who—in aggregate—will tend to behave in a certain way.

This demonstrates in a nutshell what I call *The Prediction Effect*. Predicting better than pure guesswork, even if not accurately, delivers real value. A hazy view of what's to come outperforms complete darkness by a landslide.

The Prediction Effect: *A little prediction goes a long way.*

This is the first of five Effects introduced in this book. You may have heard of the butterfly, Doppler, and placebo effects. Stay tuned here for the *Data*, *Induction*, *Ensemble*, and *Persuasion Effects*. Each of these Effects encompasses the fun part of science and technology: an intuitive hook that reveals how it works and why it succeeds.

THE FIELD OF DREAMS

People . . . operate with beliefs and biases. To the extent you can eliminate both and replace them with data, you gain a clear advantage.

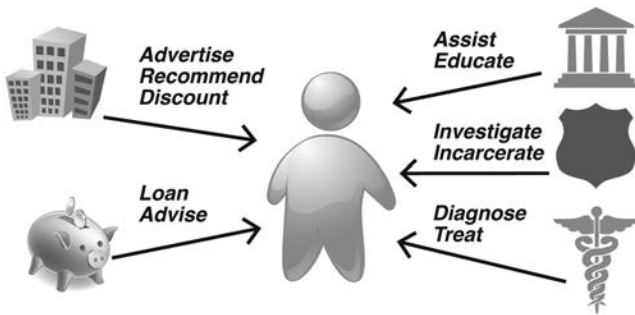
—Michael Lewis, *Moneyball: The Art of Winning an Unfair Game*

What field of study or branch of science are we talking about here? Learning how to predict from data is sometimes called *machine learning*—but it turns out this is mostly an academic term you find used within research labs, conference papers, and university courses (full disclosure: I taught the Machine Learning graduate course at Columbia University a couple of times in the late 1990s). These arenas are a priceless wellspring, but they aren't where the rubber hits the road. In commercial, industrial, and government applications—in the real-world usage of machine learning to predict—it's called something else, something that in fact is the very topic of this book:

Predictive analytics (PA)—*Technology that learns from experience (data) to predict the future behavior of individuals in order to drive better decisions.*³

³ In this definition, *individuals* is a broad term that can refer to people as well as other organizational elements. Most examples in this book involve predicting people, such as customers, debtors, applicants, employees, students, patients, donors, voters, taxpayers, potential suspects, and convicts. However, PA also applies to individual companies (e.g., for business-to-business), products, locations, restaurants, vehicles, ships, flights, deliveries, buildings, manholes, transactions, Facebook posts, movies, satellites, stocks, *Jeopardy!* questions, and much more. Whatever the domain, PA renders predictions over scalable numbers of individuals.

Built upon computer science and statistics and bolstered by devoted conferences and university degree programs, PA has emerged as its own discipline. But beyond a field of science, PA is a movement that exerts a forceful impact. Millions of decisions a day determine whom to call, mail, approve, test, diagnose, warn, investigate, incarcerate, set up on a date, and medicate. PA is the means to drive *per-person* decisions empirically, as guided by data. By answering this mountain of smaller questions, PA may in fact answer the biggest question of all: *How can we improve the effectiveness of all these massive functions across government, healthcare, business, nonprofit, and law enforcement work?*



Predictions drive how organizations treat and serve an individual, across the frontline operations that define a functional society.

In this way, PA is a completely different animal from *forecasting*. Forecasting makes aggregate predictions on a macroscopic level. How will the economy fare? Which presidential candidate will win more votes in Ohio? Whereas forecasting estimates the total number of ice cream cones to be purchased next month in Nebraska, PA tells you which *individual* Nebraskans are most likely to be seen with cone in hand.

PA leads within the growing trend to make decisions more “data driven,” relying less on one’s “gut” and more on hard, empirical evidence. Enter this fact-based domain and you’ll be attacked by buzzwords, including *analytics*, *big data*, *data science*, and *business intelligence*. While PA fits

underneath each of these umbrellas, these evocative terms refer more to the culture and general skill sets of technologists who do an assortment of creative, innovative things with data, rather than alluding to any specific technology or method. These areas are broad; in some cases, they refer simply to standard Excel reports—that is, to things that are important and require a great deal of craft, but may not rely on science or sophisticated math. And so they are more subjectively defined. As Mike Loukides, a vice president at the innovation publisher O’Reilly, once put it, “Data science is like porn—you know it when you see it.” Another term, *data mining*, is often used as a synonym for PA, but as an evocative metaphor depicting “digging around” through data in one fashion or another, it is often used more broadly as well.

ORGANIZATIONAL LEARNING

The powerhouse organizations of the Internet era, which include Google and Amazon . . . have business models that hinge on predictive models based on machine learning.

—Professor Vasant Dhar, Stern School of Business,
New York University

A breakthrough in machine learning would be worth 10 Microsofts.

—Bill Gates

An organization is sort of a “megaperson,” so shouldn’t it “megalearn”? A group comes together for the collective benefit of its members and those it serves, be it a company, government, hospital, university, or charity. Once formed, it gains from division of labor, mutually complementary skills, and the efficiency of mass production. The result is more powerful than the sum of its parts. Collective learning is the organization’s next logical step to further leverage this power. Just as a salesperson learns over time from her positive and negative interactions with sales leads, her successes, and failures, PA is the process by which an organization learns from the experience it has

collectively gained across its team members and computer systems. In fact, an organization that doesn't leverage its data in this way is like a person with a photographic memory who never bothers to think.

With only a few striking exceptions, we find that organizations, rather than individuals, benefit by employing PA. Organizations make the many, many operational decisions for which there's ample room for improvement; organizations are intrinsically inefficient and wasteful on a grand scale. Marketing casts a wide net—junk mail is marketing money wasted and trees felled to print unread brochures. An estimated 80 percent of all e-mail is spam. Risky debtors are given too much credit. Applications for government benefits are backlogged and delayed. And it's organizations that have the data to power the predictions that drive improvements in these operations.

In the commercial sector, profit is a driving force. You can well imagine the booming incentives intrinsic to rendering everyday routines more efficient, marketing more precisely, catching more fraud, avoiding bad debtors, and luring more online customers. Upgrading how business is done, PA rocks the enterprise's economies of scale, optimizing operations right where it makes the biggest difference.

THE NEW SUPER GEEK: DATA SCIENTISTS

The alternative [to thinking ahead] would be to think backwards . . . and that's just remembering.

—Sheldon, the theoretical physicist on *The Big Bang Theory*

Opportunities abound, but the profit incentive is not the only driving force. The source, the energy that makes it work, is Geek Power! I speak of the enthusiasm of technical practitioners. Truth be told, my passion for PA didn't originate from its value to organizations. I am in it for the fun. The idea of a machine that can actually learn seems so cool to me that I care more about what happens inside the magic box than its outer usefulness.

Indeed, perhaps that's the defining motivator that qualifies one as a geek. We love the technology; we're in awe of it. Case in point: The leading free, open-source software tool for PA, called R (a one-letter, geeky name), has a rapidly expanding base of users as well as enthusiastic volunteer developers who add to and support its functionalities. Great numbers of professionals and amateurs alike flock to public PA competitions with a tremendous spirit of "coopetition." We operate within organizations, or consult across them. We're in demand, so we fly a lot. But we fly coach, at best Economy Plus.

THE ART OF LEARNING

Whatcha gonna do with your CPU to reach its potentiality?

Use your noggin when you log in to crank it exponentially.

The endeavor that will render my obtuse computer clever:

Self-improve impeccably by way of trial and error.

Once upon a time, humanity created The Ultimate General Purpose Machine and, in an inexplicable fit of understatement, decided to call it "a computer" (a word that until this time had simply meant a person who did computations by hand). This automaton could crank through any demanding, detailed set of endless instructions without fail or error and with nary a complaint; within just a few decades, its speed became so blazingly brisk that humanity could only exclaim, "Gosh, we really cranked that!" An obviously much better name for this device would have been the appropriately grand *La Machine*, but a few decades later this name was hyperbolically bestowed upon a food processor (I am not joking). *Quel dommage*. "What should we do with the computer? What's its true potential, and how do we achieve it?" humanity asked of itself in wonderment.

A computer and your brain have something in common that renders them both mysterious, yet at the same time easy to take for granted. If while

pondering what this might be you heard a pin drop, you have your answer. They are both silent. Their mechanics make no sound. Sure, a computer may have a disk drive or cooling fan that stirs—just as one’s noggin may emit wheezes, sneezes, and snores—but the mammoth grunt work that takes place therein involves no “moving parts,” so these noiseless efforts go along completely unwitnessed. The smooth delivery of content on your screen—and ideas in your mind—can seem miraculous.⁴

They’re both powerful as heck, your brain and your computer. So could computers be successfully programmed to think, feel, or become truly intelligent? Who knows? At best these are stimulating philosophical questions that are difficult to answer, and at worst they are subjective benchmarks for which success could never be conclusively established. But thankfully we do have some clarity: There is one truly impressive, profound human endeavor computers *can* undertake. They can learn.

But how? It turns out that learning—generalizing from a list of examples, be it a long list or a short one—is more than just challenging. It’s a philosophically deep dilemma. Machine learning’s task is to find patterns that appear not only in the data at hand, but in general, so that what is learned will hold true in new situations never yet encountered. At the core, this ability to generalize is the magic bullet of PA. There is a true art in the design of these computer methods. We’ll explore more later, but for now I’ll give you a hint. The machine actually learns more about your next likely action by studying *others* than by studying *you*.

While I’m dispensing teasers that leave you hanging, here’s one more. This book’s final chapter answers the riddle: *What often happens to you that*

⁴ Silence is characteristic to solid state electronics, but computers didn’t have to be built that way. The idea of a general-purpose, instruction-following machine is abstract, not affixed to the notion of electricity. You could construct a computer of cogs and wheels and levers, powered by steam or gasoline. I mean, I wouldn’t recommend it, but you could. It would be slow, big, and loud, and nobody would buy it.

cannot be witnessed, and that you can't even be sure has happened afterward—but that can be predicted in advance?

Learning from data to predict is only the first step. To take the next step and *act on predictions* is to fearlessly gamble. Let's kick off Chapter 1 with a suspenseful story that shows why launching PA feels like blasting off in a rocket.

About the Author



Eric Siegel, PhD, founder of the Predictive Analytics World conference series and executive editor of *The Predictive Analytics Times*, makes the how and why of predictive analytics understandable and captivating. Eric is a former Columbia University professor—who used to sing educational songs to his students—and a renowned speaker, educator, and leader in the field.

Eric has appeared on Al Jazeera America, Bloomberg TV and Radio, Business News Network (Canada), Fox News, Israel National Radio, NPR Marketplace, Radio National (Australia), and TheStreet. He and this book have been featured in *Businessweek*, *CBS MoneyWatch*, *The Financial Times*, *Forbes*, *Forrester*, *Fortune*, *The Huffington Post*, *The New York Review of Books*, *Newsweek*, *The Seattle Post-Intelligencer*, *The Wall Street Journal*, *The Washington Post*, and *WSJ MarketWatch*.

Eric Siegel is available for select lectures. To inquire: www.ThePredictionBook.com

Interested in employing predictive analytics at your organization?

- Access the author's online, on-demand training workshop, Predictive Analytics Applied: www.businessprediction.com
- Get started with the Predictive Analytics Guide: www.pawcon.com/guide
- Follow Eric Siegel on Twitter: [@predictanalytic](https://twitter.com/predictanalytic)

REVISED AND UPDATED

PREDICTIVE ANALYTICS

"Mesmerizing & fascinating..."
—The Seattle Post-Intelligencer

AN INTRODUCTION
FOR EVERYONE



THE POWER TO PREDICT WHO WILL
CLICK, BUY, LIE, OR DIE

ERIC SIEGEL

WILEY

amazon[®]

BARNES & NOBLE
BOOKSELLERS

BAM!
BOOKS-A-MILLION

TRANSLATED INTO 9 LANGUAGES USED IN COURSES AT MORE THAN 30 UNIVERSITIES

In this rich, fascinating—and surprisingly accessible—introduction, leading expert Eric Siegel reveals how predictive analytics works, and how it affects everyone every day.

Trendsetters like Chase, Facebook, Google, Hillary for America, HP, IBM, Match.com, Netflix, the NSA, Pfizer, Target, and Uber are seizing upon the power of big data to predict human behavior—including yours.

Why? Predictive analytics reinvents industries and runs the world. Read on to discover how it combats risk, boosts sales, fortifies healthcare, optimizes social networks, toughens crime fighting, and wins elections.



Photo Credit: Dana Patrick

ERIC SIEGEL, PhD, is the founder of Predictive Analytics World and executive editor of *The Predictive Analytics Times*. A former Columbia University professor, he is a renowned speaker, educator, and leader in the field.

Learn more: www.ThePredictionBook.com

“What Nate Silver did for poker and politics, this does for everything else.”

—David Leinweber, author of *Nerds on Wall Street*

“The *Freakonomics* of big data.”

—Stein Kretzinger, founding executive, Advertising.com

“A deeply informative dive into a topic that is critical to virtually every sector of business today.”

—Geoffrey Moore, author of *Crossing the Chasm*

“*Moneyball* for business, government, and healthcare.”

—Jim Sterne, founder, eMetrics Summit




BUSINESS & ECONOMICS/
Popular Culture

Cover Design: Wiley
Cover Image: Winona Nelson

Subscribe to our free Business eNewsletter at wiley.com/enewsletters

Visit wiley.com/business

WILEY

 Also available as an e-book

ISBN 978-1-119-14567-7



9 781119 145677