

O'REILLY®

Werde ein Data Head

Data Science, Machine Learning und Statistik
verstehen und datenintensive Jobs meistern



Alex J. Gutman, Jordan Goldmeier
Übersetzung von Jørgen W. Lang

Papier
plus⁺
PDF.

Zu diesem Buch – sowie zu vielen weiteren O’Reilly-Büchern – können Sie auch das entsprechende E-Book im PDF-Format herunterladen. Werden Sie dazu einfach Mitglied bei oreilly.plus⁺:

www.oreilly.plus

Werde ein Data Head

*Data Science, Machine Learning und
Statistik
verstehen und datenintensive Jobs
meistern*

***Alex J. Gutman, Jordan
Goldmeier***

*Deutsche Übersetzung von
Jørgen W. Lang*

O'REILLY®

Alex J. Gutman, Jordan Goldmeier

Lektorat: Alexandra Follenius

Übersetzung: Jørgen W. Lang

Fachgutachten: Marcus Fraaß

Korrektur: Sibylle Feldmann, www.richtiger-text.de

Satz: III-satz, www.drei-satz.de

Herstellung: Stefanie Weidner

Umschlaggestaltung: Michael Oréal, www.oreal.de, unter Verwendung der
iStock-Illustration

ID 1173117448 von Vertigo3d/Getty Images

Bibliografische Information der Deutschen Nationalbibliothek

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der
Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im
Internet über <http://dnb.d-nb.de> abrufbar.

ISBN:

Print 978-3-96009-191-2

PDF 978-3-96010-667-8

ePub 978-3-96010-668-5

mobi 978-3-96010-669-2

1. Auflage 2022

Translation Copyright für die deutschsprachige Ausgabe © 2022 dpunkt.verlag
GmbH

Wieblinger Weg 17

69123 Heidelberg

Authorized German translation of the English edition of *Becoming a Data Head: How to Think, Speak, and Understand Data Science, Statistics, and Machine Learning* by Alex J. Gutman and Jordan Goldmeier, ISBN 9781119741749 © 2021 John Wiley & Sons, Inc., Indianapolis, Indiana. All rights reserved.

Dieses Buch erscheint in Kooperation mit O'Reilly Media, Inc. unter dem
Imprint »O'REILLY«.

O'REILLY ist ein Markenzeichen und eine eingetragene Marke von O'Reilly
Media, Inc. und wird mit Einwilligung des Eigentümers verwendet.



Hinweis:

Dieses Buch wurde auf PEFC-zertifiziertem Papier aus nachhaltiger Waldwirtschaft gedruckt. Der Umwelt zuliebe verzichten wir zusätzlich auf die Einschweißfolie.

Schreiben Sie uns:

Falls Sie Anregungen, Wünsche und Kommentare haben, lassen Sie es uns wissen: komentar@oreilly.de.

Die vorliegende Publikation ist urheberrechtlich geschützt. Alle Rechte vorbehalten. Die Verwendung der Texte und Abbildungen, auch auszugsweise, ist ohne die schriftliche Zustimmung des Verlags urheberrechtswidrig und daher strafbar. Dies gilt insbesondere für die Vervielfältigung, Übersetzung oder die Verwendung in elektronischen Systemen.

Es wird darauf hingewiesen, dass die im Buch verwendeten Soft- und Hardware-Bezeichnungen sowie Markennamen und Produktbezeichnungen der jeweiligen Firmen im Allgemeinen warenzeichen-, marken- oder patentrechtlichem Schutz unterliegen.

Alle Angaben und Programme in diesem Buch wurden mit größter Sorgfalt kontrolliert. Weder Autor noch Verlag noch Übersetzer können jedoch für Schäden haftbar gemacht werden, die in Zusammenhang mit der Verwendung dieses Buches stehen.

5 4 3 2 1 0

Inhalt

Vorwort

Einleitung

Die Data-Science-Industrie

Warum uns das Thema so wichtig ist

Die Krise auf dem US-amerikanischen Subprime-Hypothekenmarkt

Die US-Präsidentenwahl von 2016

Unsere Hypothese

Daten am Arbeitsplatz

Die berühmte Sitzungssaal-Szene

Sie können das große Ganze verstehen

Restaurants klassifizieren

Ja und?

Für wen dieses Buch geschrieben wurde

Warum wir dieses Buch geschrieben haben

Was Sie lernen werden

Wie dieses Buch strukturiert ist

Ein letzter Punkt, bevor es wirklich losgeht

Teil I Denken wie ein Data Head

1 Was ist das Problem?

Fragen, die ein Data Head stellen sollte

Warum ist das Problem wichtig?

Wen betrifft das Problem?

Was ist, wenn wir nicht die richtigen Daten haben?

Wann ist das Projekt zu Ende?

Was tun wir, wenn uns die Ergebnisse nicht gefallen?

Verstehen, warum Datenprojekte scheitern

Szenario: Kundenwahrnehmung

Diskussion

An den wichtigen Problemen arbeiten

Zusammenfassung

2 Was sind Daten?

Daten oder Informationen?

Ein Beispiel-Datensatz

Datentypen

Wie Daten gesammelt und strukturiert werden

Beobachtungsbasierte versus experimentelle Daten

Strukturierte versus unstrukturierte Daten

Die Basics der zusammenfassenden Statistik

Zusammenfassung

3 Vorbereitungen für das statistische Denken

Stellen Sie Fragen!

In allen Dingen ist Variation

Szenario: Kundenwahrnehmung (die Fortsetzung)

Fallstudie: Nierenkrebsraten

Wahrscheinlichkeitsrechnung und Statistik

Wahrscheinlichkeit oder Intuition

Entdeckungen mit Statistiken

Teil II Sprechen wie ein Data Head

4 Daten infrage stellen

Was würden Sie tun?

Katastrophe durch fehlende Daten

Erzählen Sie mir die Herkunftsgeschichte der Daten

Wer hat die Daten gesammelt?

Wie wurden die Daten gesammelt?

Sind die Daten repräsentativ?

Gibt es eine Stichprobenverzerrung?

Wie wurde mit Ausreißern umgegangen?

Welche Daten sehe ich nicht?

Wie gehen Sie mit fehlenden Werten um?

Können die Daten abbilden, was Sie mit ihnen messen wollen?

Stellen Sie Daten infrage, egal wie groß die Datenmenge ist

Zusammenfassung

5 Daten erkunden

Ihre Rolle in der explorativen Datenanalyse

Wie ein Forscher denken

Leitfragen

Der Versuchsaufbau

Können die Daten Ihre Frage beantworten?

Legen Sie Erwartungen fest und benutzen Sie Ihren gesunden Menschenverstand

Ergeben die Werte intuitiv einen Sinn?

Achtung: Ausreißer und fehlende Werte

Sind Ihnen irgendwelche Beziehungen aufgefallen?

Korrelation verstehen

Achtung: Korrelation falsch interpretieren
Achtung: Korrelation bedeutet nicht Kausalität
Haben Sie in den Daten neue Einsatzmöglichkeiten
oder unentdeckte Potenziale gefunden?
Zusammenfassung

6 Wahrscheinlichkeiten untersuchen

Raten Sie mal
Die Spielregeln
Schreibweise
Bedingte Wahrscheinlichkeit und unabhängige
Ereignisse
Die Wahrscheinlichkeit mehrfacher Ereignisse
Gedankenexperiment zur Wahrscheinlichkeit
Die nächsten Schritte
Seien Sie vorsichtig bei der Annahme von
Abhängigkeiten
Fallen Sie nicht auf den Spieler-Fehlschluss herein
Alle Wahrscheinlichkeiten unterliegen bestimmten
Bedingungen
Vertauschen Sie Abhängigkeiten nicht
Der Satz von Bayes
Stellen Sie sicher, dass die Wahrscheinlichkeiten einen
Sinn ergeben
Kalibrierung
Seltene Ereignisse können und werden eintreffen
Zusammenfassung

7 Hinterfragen Sie Statistiken

Kleine Einführung in die statistische Inferenz
Schaffen Sie sich etwas Spielraum
Mehr Daten, mehr Evidenz
Hinterfragen Sie den Status quo
Beweise für das Gegenteil (Evidenz)
Entscheidungsfehler ausgleichen

Die Vorgehensweise der statistischen Inferenz
Die Fragen, die Sie stellen sollten, um Statistiken zu hinterfragen

Was ist der Kontext für diese Statistik?

Wie groß ist der Stichprobenumfang?

Was testen Sie?

Wie lautet die Nullhypothese?

Wie hoch ist das Signifikanzniveau?

Wie viele Tests führen Sie durch?

Kann ich bitte die Konfidenzintervalle sehen?

Ist dies von praktischer Bedeutung?

Gehen Sie von einer Kausalität aus?

Zusammenfassung

Teil III Den Werkzeugkasten des Data Scientist verstehen

8 Nach versteckten Gruppen suchen

Unüberwachtes Lernen

Dimensionsreduktion

Zusammengefasste Features erstellen

Hauptkomponentenanalyse

Beispiel: HKA für die sportliche Leistungsfähigkeit

Zusammenfassung zur HKA

Mögliche Fallen

Clustering

Clustering mit dem k-Means-Algorithmus

Beispiel: Clustering von Verkaufsfilialen

Mögliche Fallen

Zusammenfassung

9 Das Regressionsmodell verstehen

Überwachtes Lernen

Was macht die lineare Regression?

Kleinste-Quadrate-Regression: mehr als nur ein hübscher Name

Vorteile der linearen Regression

Auf mehrere Features erweitern

Probleme und Fallstricke der linearen Regression

Unberücksichtigte Variablen

Multikollinearität

Data Leakage

Extrapolationsfehler

Viele Beziehungen sind nicht linear

Erklärst du noch, oder machst du schon

Vorhersagen?

Leistungsfähigkeit der Regression

Andere Regressionsmodelle

Zusammenfassung

10 Das Klassifikationsmodell verstehen

Einführung in die Klassifikation

Was Sie lernen werden

Klassifikationsproblem: Versuchsaufbau

Logistische Regression

Logistische Regression: Na und?

Entscheidungsbäume

Ensemblemethoden

Zufallswälder

Gradientenverstärkte Bäume

Interpretierbarkeit von Ensemblemethoden

Achten Sie auf Fallstricke

Falsche Anwendung des Problems

Data Leakage

Keine Aufteilung der Daten

Den richtigen Cut-off-Wert wählen

Falsch verstandene Genauigkeit

Konfusionsmatrizen

Zusammenfassung

11 Textanalyse verstehen

Erwartungen an die Textanalyse

Wie aus Text Zahlen werden

Ein großer Sack voll Wörter

N-Gramme

Worteinbettungen

Topic Modeling

Textklassifikation

Naive Bayes

Sentimentanalyse

Praktische Überlegungen bei der Arbeit mit Text

Die großen Technologiekonzerne haben die

Oberhand

Zusammenfassung

12 Konzepte des Deep Learning

Neuronale Netzwerke

Worin besteht die Ähnlichkeit zwischen neuronalen
Netzwerken und dem Gehirn?

Ein einfaches neuronales Netzwerk

Wie ein neuronales Netzwerk lernt

Ein etwas komplexeres neuronales Netzwerk

Anwendungen des Deep Learning

Die Vorteile des Deep Learning

Wie Computer Bilder »sehen«

Neuronale Konvolutionsnetze

Deep Learning für Sprache und Wortsequenzen

Deep Learning in der Praxis

Haben Sie Daten?

Sind Ihre Daten strukturiert?

Wie wird das Netzwerk aussehen?

Die künstliche Intelligenz und Sie

Die großen Technologiekonzerne haben die
Oberhand
Ethik im Deep Learning
Zusammenfassung

Teil IV Den Erfolg sichern

13 Achten Sie auf Fallstricke

Bias und seltsame Datenphänomene
Survivorship Bias
Regression zur Mitte
Das Simpson-Paradoxon
Confirmation Bias
Effort Bias
Algorithmischer Bias
Weitere Formen von Bias
Die große Liste möglicher Fallstricke
Fallstricke der Statistik und des Machine Learning
Projektbezogene Fallstricke
Zusammenfassung

14 Menschen und Persönlichkeiten kennen

Sieben Szenarien typischer Kommunikationspannen
Das Postmortem
Märchenstunde
Stille Post
Verzettelt
Der Realitätsabgleich
Die Übernahme
Der Angeber
Datenpersönlichkeiten
Datenenthusiasten
Datenzyniker

Data Heads
Zusammenfassung

15 Was kommt danach?

Danksagungen

Index

Für meine Kinder Allie, William und Ellen

Allie war gerade drei Jahre alt, als sie herausfand, dass ihr Vater ein »Doktor« ist. Etwas irritiert sah sie mich an und sagte: »Aber du hilfst den Menschen doch gar nicht ...« In diesem Sinne widme ich dieses Buch auch Ihnen, den Leserinnen und Lesern.

Ich hoffe, dass es Ihnen hilft.

- Alex

Für Stephen und Melissa

- Jordan

Vorwort

Werde ein Data Head kommt angesichts der aktuellen Situation der Daten und Analysen in vielen Organisationen genau zur richtigen Zeit. Werfen wir einen kurzen Blick auf die nähere Vergangenheit. Einige wenige führende Unternehmen setzen seit mehreren Jahrzehnten, genauer gesagt seit den 1970er-Jahren, Daten und Analysen effektiv als Orientierungshilfe für ihre Entscheidungen ein. Die meisten anderen Unternehmen haben diese wertvolle Ressource dagegen schlicht ignoriert oder versteckten sie in irgendeinem Hinterzimmer, ohne ihr viel Beachtung zu schenken.

Das begann sich Anfang der 2000er-Jahre zu ändern. Unternehmen begeisterten sich für die Möglichkeiten, ihr Geschäft auf Basis von Daten und Analysemethoden neu aufzustellen. In den frühen 2010er-Jahren verschob sich diese Begeisterung in Richtung *Big Data*, einem Begriff, der ursprünglich von Internetunternehmen stammt, sich aber schnell auf alle fortschrittlichen Wirtschaftszweige ausbreitete. Um mit der ständig steigenden Menge und Komplexität der Daten umgehen zu können, entstand in vielen Unternehmen die Rolle des *Data Scientist* - auch hier zuerst im Silicon Valley und dann überall.

Aber gerade als viele Unternehmen begannen, sich mit Big Data auseinanderzusetzen, verschob sich zwischen 2015 und 2018 das Hauptaugenmerk erneut, und zwar auf künstliche Intelligenz. Das Sammeln, Speichern und Analysieren großer Datenmengen musste Machine Learning, natürlicher Sprachverarbeitung (engl. *Natural Language Processing*, NLP) und Automatisierung weichen.

Eingebettet in diese sehr schnellen Verschiebungen der Aufmerksamkeit, entstand eine Reihe von Annahmen über Daten und Datenanalyse in Unternehmen. Ich bin froh, dass *Werde ein Data Head* viele dieser Annahmen über den Haufen wirft, denn das ist schon lange fällig. Viele, die mit diesen Trends arbeiten oder sie genau beobachten, geben langsam zu, dass sie durch diese Annahmen immer unproduktiver wurden. Im Rest dieses Vorworts werde ich daher fünf miteinander verbundene Annahmen beschreiben und zeigen, auf welche Weise die Ideen in diesem Buch ihnen zu Recht widersprechen.

Annahme 1: Datenanalyse, Big Data und KI sind vollkommen unterschiedliche Phänomene.

Außenstehende gehen oft davon aus, dass »traditionelle« Analysemethoden, Big Data und KI vollkommen eigenständige und unterschiedliche Themenbereiche sind. *Werde ein Data Head* zeigt dagegen, dass diese Themen sogar sehr eng miteinander verknüpft sind. Bei allen geht es um statistisches Denken. Traditionelle Analysemethoden wie die Regressionsanalyse kommen ebenfalls in allen drei Bereichen zum Einsatz und auch in Techniken zur Datenvisualisierung. Die prädiktive Analyse ist im Grunde nichts anderes als überwacht Machine Learning, und die meisten Techniken zur Datenanalyse

funktionieren auf Datensätzen beliebiger Größe. Kurz gesagt: Ein guter Data Head bewegt sich effektiv in allen drei Bereichen und weiß, dass es wenig produktiv ist, sich zu sehr mit den Unterschieden zu beschäftigen.

Annahme 2: Data Scientists sind die einzigen Personen, die in diesem Sandkasten spielen können.

Wir haben Data Scientists über den grünen Klee gelobt und sind oft davon ausgegangen, sie seien die einzigen Menschen, die effektiv mit Daten und Analysen arbeiten können. Tatsächlich findet aktuell eine überaus wichtige Bewegung in Richtung Demokratisierung dieser Ideen statt. Immer mehr Unternehmen setzen auf sogenannte »Laien-Data-Scientists«. Automatische Werkzeuge für das Machine Learning erleichtern die Erstellung hervorragender Vorhersagemodelle. Natürlich gibt es auch weiterhin Bedarf an professionellen Data Scientists, die neue Algorithmen entwickeln und die Arbeit der Laien überwachen, besonders wenn komplexe Analysen durchgeführt werden. Doch Unternehmen, die Analysen und Data Science demokratisieren und »Laien-Data-Heads« einsetzen, können die gesamte Nutzung dieser wichtigen Fähigkeiten deutlich steigern.

Annahme 3: Ein Data Scientist ist eine »Eier legende Wollmilchsau«, die alle Fähigkeiten besitzt, die für diese Aufgaben nötig sind.

Oft wird davon ausgegangen, dass Data Scientists, also Personen, deren Ausbildungs- und Arbeitsschwerpunkt auf der Entwicklung und Programmierung von Modellen liegt, auch alle anderen Aufgaben ausführen können, die für eine vollständige Implementierung dieser Modelle nötig sind. Anders gesagt: Wir stellen sie uns als eine

Art Eier legende Wollmilchsau vor, also als wahre Alleskönner. Aber diese Alleskönner gibt es nicht, oder sie sind nur sehr selten anzutreffen. Data Heads, die nicht nur die Feinheiten der Data Science, sondern auch den geschäftlichen Teil kennen, können effektiv Projekte leiten, haben zudem ausgezeichnete Fähigkeiten im Aufbau von Geschäftsbeziehungen und sind eine sehr wertvolle Ressource in Data-Science-Projekten. Sie können produktive Mitglieder von Data-Science-Teams sein und erhöhen die Wahrscheinlichkeit, dass Data-Science-Projekte den Geschäftswert steigern.

Annahme 4: Sie brauchen einen wirklich hohen Intelligenzquotienten und eine Menge Training, um mit Daten und Analysen erfolgreich zu sein.

Eine weitere verwandte Annahme besagt, dass man sehr gut in Data Science ausgebildet sein muss, um in diesem Bereich zu bestehen, und dass ein Data Head sehr gut mit Zahlen umgehen können muss. Training im Umgang mit Zahlen und Sachverstand sind sicher eine Hilfe. *Werde ein Data Head* vertritt allerdings die Meinung (der ich übrigens zustimme), dass man mit etwas Ehrgeiz durchaus in der Lage ist, sich die nötigen Fähigkeiten zu Daten und Analysen anzueignen, um in Data-Science-Projekten nützlich zu sein. Das liegt zum Teil daran, dass die Grundprinzipien statistischer Analyse durchaus keine Raketenwissenschaften sind. Man muss sich nicht einmal extrem gut mit Daten und Analysen auskennen, um in Data-Science-Projekten »nützlich zu sein«. Für die Arbeit mit ausgebildeten Data Scientists oder automatisierten KI-Programmen muss man nur die Fähigkeit und die Neugier besitzen, gute Fragen zu stellen und die Verbindungen zwischen geschäftlichen Themen und quantitativen Ergebnissen

herzustellen – ohne dabei auf zweifelhafte Annahmen hereinzufallen.

Annahme 5: Wenn Sie keine quantitativen Fächer (Algebra, Statistik etc.) studiert haben, ist es schon zu spät, sich das für die Arbeit mit Daten und Analysen nötige Wissen anzueignen.

Diese Annahme wird von Umfragedaten gestützt. In einer Umfrage der Internetplattform Splunk aus dem Jahr 2019 unter rund 1.300 Führungskräften weltweit gaben praktisch alle Befragten (98 %) an, dass Fähigkeiten im Umgang mit Daten für die Arbeitsplätze von morgen eine wichtige Rolle spielen.¹ 81 % der Führungskräfte waren außerdem der Ansicht, dass Fähigkeiten im Umgang mit Daten nötig sein werden, um höhere Führungspositionen einzunehmen. 85 % waren sich einig, dass diese Kenntnisse für ihre Unternehmen immer wertvoller werden. Trotzdem gaben 67 % an, sich nicht wohl dabei zu fühlen, selbst auf Daten zuzugreifen oder diese zu nutzen. 73 % glaubten, dass Kenntnisse im Umgang mit Daten schwerer erlernbar sind als andere Geschäftsfähigkeiten. 53 % waren der Meinung, sie wären zu alt, um Fähigkeiten im Umgang mit Daten noch zu erlernen. Dieser »Daten-Defätismus« ist schädlich für Einzelpersonen wie für Unternehmen. Weder die Autoren dieses Buchs noch ich selbst halte ihn für gerechtfertigt. Sehen Sie sich die Seiten nach diesem Vorwort an, und Sie werden feststellen, dass wirklich keine Raketenwissenschaft nötig ist.

Vergessen Sie also diese falschen Annahmen und werden Sie zum Data Head. Ihr Wert als Mitarbeiterin oder Mitarbeiter wird steigen, und Sie werden Ihr Unternehmen

erfolgreicher machen. Das ist der Lauf der Welt – es ist an der Zeit, sich damit vertraut zu machen und mehr über Daten und Analysen zu lernen. Ich bin überzeugt, Sie werden den Prozess und die Lektüre von *Werde ein Data Head* als lohnender und angenehmer empfinden, als Sie es sich vorstellen können.

Thomas H. Davenport
Distinguished Professor, Babson College Visiting Professor,
Oxford Saïd Business School Research Fellow, MIT
Initiative on the Digital Economy Autor von *Competing on
Analytics, Big Data @ Work* und *The AI Advantage*

Einleitung

Ob Sie wollen oder nicht: Daten sind wahrscheinlich der wichtigste Aspekt Ihrer Arbeit. Und sehr wahrscheinlich lesen Sie dieses Buch, um verstehen zu können, worum es überhaupt geht.

Zu Beginn lohnt es sich, noch einmal auszusprechen, was fast schon ein Klischee ist: Wir erzeugen und konsumieren mehr Informationen als jemals zuvor. Wir befinden uns ohne Zweifel im Zeitalter der Daten. Und dieses Zeitalter hat einen ganz eigenen Wirtschaftszweig mit Versprechen, Buzzwords und Produkten hervorgebracht, die Sie, Ihre Vorgesetzten, Ihre Kolleginnen und Kollegen sowie Ihre Mitarbeitenden benutzen oder benutzen werden. Aber trotz aller Behauptungen und weitverbreiteten Datenversprechen und -produkten schlagen Data-Science-Projekte mit alarmierender Häufigkeit fehl.¹

Damit wollen wir nicht sagen, dass alle Datenversprechen leer und alle Produkte furchtbar sind. Es geht eher darum, dass Sie eine grundsätzliche Wahrheit erkennen müssen, um das Thema wirklich verstehen zu können: Dieses Zeug ist wirklich komplex. Bei der Arbeit mit Daten geht es um Zahlen, feine Unterschiede und Unsicherheit. Sicher, Daten sind wichtig, aber selten einfach. Und trotzdem gibt es eine ganze Branche, die versucht, uns etwas anderes zu

erzählen. Eine Branche, die uns Sicherheit in einer unsicheren Welt verspricht und mit der Angst der Unternehmen spielt, etwas zu verpassen. Wir, die Autoren, nennen dies die Data-Science-Industrie.

Die Data-Science-Industrie

Dieses Problem betrifft alle Beteiligten. Unternehmen suchen ständig nach Produkten, die ihnen das Denken abnehmen. Manager stellen Analyseprofis ein, die in Wirklichkeit keine sind. Data Scientists werden von Unternehmen angeheuert, die eigentlich noch gar nicht dafür bereit sind. Führungskräfte werden gezwungen, sich technologisches Fachchinesisch anzuhören und so zu tun, als verstünden sie alles Gesagte. Projekte geraten in Stocken, Geld wird verschwendet.

Gleichzeitig spuckt die Data-Science-Industrie schneller neue Konzepte aus, als wir in der Lage sind, die neu geschaffenen Möglichkeiten (und Probleme) zu erfassen und auf den Punkt zu bringen. Ein Augenblick - und schon ist wieder eine Chance verpasst. Als die Autoren ihre Zusammenarbeit begannen, war *Big Data* das große Zauberwort. Im Laufe der Zeit wurde dann *Data Science* das neue Thema. Mittlerweile liegt das Hauptaugenmerk auf Dingen wie *Machine Learning*, *Deep Learning* und *künstlicher Intelligenz*.

Für die neugierigen und kritischen Denker unter uns scheint hier irgendetwas nicht zu stimmen. Sind diese Problemstellungen wirklich neu? Oder sind die neuen Begriffe nur alter Wein in neuen Schläuchen?

Die Antwort lautet für beide Fragen natürlich: Ja.

Die größere und wichtigere Frage, die Sie sich hoffentlich stellen, lautet allerdings: *Wie kann ich kritisch über Daten denken und sprechen?*

Genau das wollen wir Ihnen hier beibringen.

Mit diesem Buch geben wir Ihnen die Werkzeuge, Fachbegriffe und Denkweisen an die Hand, die nötig sind, um sich in der Data-Science-Branche zu orientieren und die gesteckten Ziele zu erreichen. Sie werden ein tieferes Verständnis für Daten und ihre Herausforderungen entwickeln. Sie werden lernen, kritisch über Daten und die gefundenen Ergebnisse zu denken, und Sie werden in der Lage sein, informiert und klug über alles zu sprechen, was mit Daten zu tun hat.

Kurz gesagt, Sie werden ein *Data Head*.

Warum uns das Thema so wichtig ist

Bevor wir uns mit den Details befassen, ist es sinnvoll, zu verstehen, warum Ihren Autoren Alex und Jordan dieses Thema so sehr am Herzen liegt. In diesem Abschnitt zeigen wir Ihnen zwei wichtige Beispiele dafür, wie Daten Einfluss auf große Teile der Gesellschaft und uns persönlich genommen haben.

Die Krise auf dem US-amerikanischen Subprime-Hypothekenmarkt

Wir kamen gerade frisch vom College, als die Subprime-Hypothekenkrise über uns hereinbrach. 2009, in einer Zeit, in der es schwer war, überhaupt einen Job zu bekommen, schafften wir es beide, Arbeit bei der Air Force zu

bekommen. Wir hatten beide Glück, weil wir eine sehr gefragte Fähigkeit besaßen: Wir konnten mit Daten umgehen. Tagein, tagaus arbeiteten wir mit Daten, um die Forschung von Air-Force-Analysten und -Wissenschaftlern in Produkte zu verwandeln, mit denen die Regierung etwas anfangen konnte. Unsere Anstellung sollte zu einem Vorboten der Aufmerksamkeit werden, die das ganze Land bald den von uns ausgefüllten Rollen widmen sollte. Als zwei Datenanalysten betrachteten wir die Hypothekenkrise mit Interesse und Neugier.

Zum Entstehen der Subprime-Hypothekenkrise trug eine Reihe verschiedener Faktoren bei.² In unserem Versuch, sie als Beispiel zu verwenden, wollen wir weitere Faktoren nicht ignorieren. Dennoch sehen wir, vereinfacht gesagt, die Krise als einen großen Datenfehler. Banken und Investoren erstellten Modelle, um den Wert von hypothekarisch abgesicherten Schuldverschreibungen (engl. *Mortgage-backed Collateralized Debt Obligations*, CDOs) zu verstehen. Vielleicht erinnern Sie sich, dass genau dieses Investitionsmodell für den Zusammenbruch der Märkte in den Vereinigten Staaten verantwortlich war.

CDOs wurden als sichere Investition angesehen, weil das Kreditausfallrisiko auf mehrere Investitionseinheiten verteilt wird. Der Gedanke war, dass in einem Portfolio von Hypotheken der Ausfall einiger weniger Hypotheken keine wesentlichen Auswirkungen auf den zugrunde liegenden Wert des gesamten Portfolios haben würde.

Und trotzdem wissen wir mittlerweile, dass einige grundlegende Annahmen falsch waren. Am schwersten wog die Fehleinschätzung, dass Kreditausfälle voneinander unabhängige Ereignisse waren. Wenn Person A ihren Kredit nicht zurückzahlen kann, hat das keinen Einfluss auf Person B – dachte man. Wenig später mussten wir lernen,

dass Kreditausfälle eher wie Dominosteine funktionieren, bei denen ein vorheriger Ausfall ein Anzeichen für weitere Ausfälle ist. Sobald eine Hypothek geplatzt war, sanken in der Folge die Immobilienpreise in der Umgebung, und das Risiko für weitere Ausfälle in dieser Wohngegend stieg. Durch den Kreditausfall wurden die benachbarten Häuser mit in den Abgrund gerissen.

Von Unabhängigkeit auszugehen, wenn die Ereignisse tatsächlich einen Zusammenhang haben, ist ein häufig anzutreffender Fehler in der Statistik.

Aber tauchen wir noch etwas tiefer in die Geschichte ein. Investmentbanken hatten ein Modell geschaffen, das Investitionen überbewertete. Ein Modell ist ein absichtlich stark vereinfachtes Abbild einer realen Situation. Es basiert auf Annahmen über die echte Welt, um bestimmte Phänomene besser zu verstehen und Vorhersagen darüber zu treffen. Auf Modelle werden wir weiter unten im Buch noch genauer eingehen.

Und wer waren die Leute, die dieses Modell erstellt und verstanden haben? Das waren genau diejenigen, die die Grundlagen für ein Berufsbild geschaffen haben, das wir heute als *Data Scientist* bezeichnen. Leute wie wir. Statistiker, Ökonomen, Physiker - Leute, die sich mit Machine Learning, künstlicher Intelligenz und Statistik befassen. Sie arbeiteten mit Daten. Sie waren schlau. Superschlau.

Und trotzdem ging etwas schief. Haben sie nicht die richtigen Fragen zu ihrer Arbeit gestellt? Gingen die Risikoeinschätzungen bei einer Runde »Stille Post« in den Telefonaten zwischen Analysten und Entscheidungsträgern verloren? Wurde die Unsicherheit in jeder Runde des Spiels immer weiter zur Seite geschoben, bis der Eindruck eines perfekt vorhersagbaren Wohnungsmarkts entstand? Oder

haben die Beteiligten über die tatsächlichen Ereignisse einfach gelogen?

Für uns persönlich ist die Frage viel wichtiger, wie wir ähnliche Fehler bei unserer eigenen Arbeit vermeiden können.

Wir hatten viele Fragen und konnten über die Antworten nur spekulieren. Eine Sache aber war klar: Hier geschah eine flächendeckende Datenkatastrophe. Und es würde nicht die letzte sein.

Die US-Präsidentschaftswahl von 2016

Am 8. November 2016 gewann der republikanische Kandidat Donald J. Trump die Präsidentschaftswahl in den USA gegen die vermeintliche Spitzenkandidatin und demokratische Herausforderin Hillary Clinton. Für die politischen Meinungsforscher war das ein Schock. Ihre Modelle hatten seinen Sieg nicht vorhergesagt. Und ausgerechnet das sollte das Jahr der Wahlvorhersagen sein.

Im Jahr 2008 gelang dem Blog *FiveThirtyEight* von Nate Silver - damals noch Teil der New York Times - eine erstaunlich genaue Vorhersage von Barack Obamas Wahlgewinn. Zu der Zeit waren die Experten noch skeptisch, dennoch sagte Silvers Algorithmus das Wahlergebnis korrekt voraus. 2012 stand Silver erneut im Rampenlicht, weil er einen weiteren Sieg für Barack Obama richtig vorhergesagt hatte.

Zu dieser Zeit begann die Geschäftswelt, Daten als wichtig anzusehen und Data Scientists einzustellen. Die erfolgreiche Vorhersage der Wiederwahl von Barack Obama durch Nate Silver verstärkte noch die Bedeutung der fast orakelhaften Fähigkeiten datenbasierter

Vorhersagen. Artikel in Businessmagazinen warnten Führungskräfte vor der Gefahr, von Mitbewerbern geschluckt zu werden, wenn diese ihr Geschäft datenbasiert betrieben, das eigene Unternehmen aber nicht. Die Data-Science-Industrie nahm richtig Fahrt auf.

Bis zum Jahr 2016 hatte jede größere Nachrichtenagentur in Vorhersagealgorithmen investiert, um das Ergebnis der nächsten Präsidentschaftswahlen vorauszuberechnen. Die allergrößte Mehrheit der Modelle sah einen überwältigenden Sieg der demokratischen Kandidatin Hillary Clinton voraus. Oh, wie falsch sie lagen!

Vergleichen wir das mit der Subprime-Hypothekenkrise. Man sollte davon ausgehen, dass man viel aus der Vergangenheit hätte lernen können. Das Interesse an Data Science hätte dazu führen müssen, dass Fehler vermieden werden. Und das stimmt auch: Seit 2008 und 2012 haben Nachrichtenagenturen Data Scientists eingestellt, in Umfrageforschung investiert, Datenteams geschaffen und mehr Geld für gute Daten ausgegeben.

Das führt uns nun zu der Frage: Was ist trotz dieses Einsatzes an Zeit, Geld, Aufwand und Ausbildung denn nun wirklich passiert?³

Unsere Hypothese

Warum gibt es Datenprobleme wie diese? Wir sehen drei Gründe: schwer zu lösende Probleme, Mangel an kritischem Denken und schlechte Kommunikation.

Erstens, wie bereits gesagt: *Dieses Zeug ist komplex*. Viele Datenprobleme sind äußerst schwer zu lösen – selbst mit einer Menge Daten und den richtigen Werkzeugen. Auch mit den besten Vorgehensweisen und den schlauesten

Analysten treten Fehler auf. Vorhersagen können und werden danebenliegen. Das ist einfach so.

Zweitens haben einige Analysten und Entscheider aufgehört, kritisch über Datenprobleme nachzudenken. Die Data-Science-Industrie zeichnete in ihrer Selbstüberschätzung ein Bild von Sicherheit und Einfachheit, und einige Menschen nahmen einfach alles für bare Münze. Vielleicht ist es auch nur menschlich, nicht zugeben zu wollen, dass man keine Ahnung davon hat, was gerade wirklich passiert. Dabei darf man sich nichts vormachen: Beim Nachdenken über Daten und deren Einsatz kann es auch zu falschen Entscheidungen kommen. Das bedeutet, Risiken und Unwägbarkeiten müssen klar kommuniziert werden. Aus irgendeinem Grund ist diese Nachricht wohl untergegangen. Obwohl wir eigentlich gehofft hatten, dass der enorme Fortschritt bei der Erforschung und Anwendung von Datenanalysen das kritische Denken aller schärft, hat es bei einigen eher zu einer kompletten Abschaltung geführt.

Der dritte Grund, warum Datenprobleme unserer Meinung nach auftreten, ist schlechte Kommunikation zwischen Data Scientists und Entscheidern. Trotz bester Absichten gehen Ergebnisse oft auf dem Weg der Übersetzung verloren. Nur selten sprechen Entscheider die Sprache der Data Scientists, weil sich niemand die Arbeit gemacht hat, sie ihnen beizubringen. Und ganz ehrlich: Datenanalysten sind nicht unbedingt gut darin, Dinge zu erklären. Hier gibt es eine klare Kommunikationslücke.

Daten am Arbeitsplatz

Ihre Datenprobleme werden vielleicht nicht die Weltwirtschaft zum Einsturz bringen oder den nächsten