

FOM-Edition

Sebastian Sauer

Moderne Datenanalyse mit R

Daten einlesen, aufbereiten,
visualisieren, modellieren
und kommunizieren



Springer Gabler

FOM-Edition

FOM Hochschule für Oekonomie & Management

Reihenherausgeber

FOM Hochschule für Oekonomie & Management, Essen, Deutschland

Dieses Werk erscheint in der FOM-Edition, herausgegeben von der FOM Hochschule für Oekonomie & Management.

Weitere Bände in der Reihe

<http://www.springer.com/series/12753>

Sebastian Sauer

Moderne Datenanalyse mit R

Daten einlesen, aufbereiten,
visualisieren, modellieren
und kommunizieren

 Springer Gabler



Sebastian Sauer
FOM Hochschule für Oekonomie &
Management
Nürnberg, Deutschland

FOM-Edition

ISBN 978-3-658-21586-6

<https://doi.org/10.1007/978-3-658-21587-3>

ISBN 978-3-658-21587-3 (eBook)

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

Springer Gabler

© Springer Fachmedien Wiesbaden GmbH, ein Teil von Springer Nature 2019

Das Werk einschließlich aller seiner Teile ist urheberrechtlich geschützt. Jede Verwertung, die nicht ausdrücklich vom Urheberrechtsgesetz zugelassen ist, bedarf der vorherigen Zustimmung des Verlags. Das gilt insbesondere für Vervielfältigungen, Bearbeitungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Verarbeitung in elektronischen Systemen.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften. Der Verlag, die Autoren und die Herausgeber gehen davon aus, dass die Angaben und Informationen in diesem Werk zum Zeitpunkt der Veröffentlichung vollständig und korrekt sind. Weder der Verlag noch die Autoren oder die Herausgeber übernehmen, ausdrücklich oder implizit, Gewähr für den Inhalt des Werkes, etwaige Fehler oder Äußerungen. Der Verlag bleibt im Hinblick auf geografische Zuordnungen und Gebietsbezeichnungen in veröffentlichten Karten und Institutionsadressen neutral.

Springer Gabler ist ein Imprint der eingetragenen Gesellschaft Springer Fachmedien Wiesbaden GmbH und ist ein Teil von Springer Nature.

Die Anschrift der Gesellschaft ist: Abraham-Lincoln-Str. 46, 65189 Wiesbaden, Germany

Vorwort

Wir fühlen, dass selbst, wenn alle möglichen wissenschaftlichen Fragen beantwortet sind, unsere Lebensprobleme noch gar nicht berührt sind.

– Wittgenstein, *Tractatus*, 6.52

Dystopie eines Datenzeitalters

Wir leben im frühen Zeitalter der Daten und Algorithmen, dem *Algorithmozän*. Gebetsmühlenartig haben uns Unternehmensberatungen, Politiker und Google-Ingenieure dieses Mantra vorgetragen, so dass es längt zum Fundus säuberlich abgehefteter Binsenweisheiten gehört (vgl. Chui et al. 2018). Untermalt wird dieses sonore Flüstern durch ein Stakkato von Eilmeldungen wie kürzlich von *AlphaGo Zero* (D. Silver et al. 2017b). Das ist ein Programm, das sehr gut im Brettspiel *Go* ist (und wohl auch in einigen anderen Spielen laut D. Silver et al. (2017a)). Wie sein Vorgänger aus dem letzten Jahr, *AlphaGo*, basiert das Programm auf sog. neuronalen Netzen, einer wohlbekannten Methode des statistischen Modellierens (Scherer 2013). Das neue Programm spielt deutlich besser als das alte: Bei einem gemütlichen Treffen schlug der Neue den Alten vernichtend: Mit 100 zu 0 fegte der Frischling den alten Hasen vom Platz. Dabei hatte der Alte einiges vorzuweisen. Im Vorjahr hatte er einen Meister des *Go*-Spiels, einen Menschen, klar besiegt (D. Silver et al. 2017b). Interessant ist, dass *AlphaGo Zero* ohne Lernmaterial von außen auskam, im Gegensatz zu früheren Programmen wie *AlphaGo*. Das legt nahe, dass Maschinen grundsätzlich in der Lage sind, ohne Hilfe von Menschen zu lernen – und „übermenschliche“ Leistung im *Go*-Spiel und vielleicht auch anderswo zu erzielen.

Ähnlich spektakulär: Eine Reihe von Fachartikeln zeigte, dass Algorithmen – ausreichend mit Daten gefüttert – die Persönlichkeit einer Person besser einschätzen können als deren Freunde (Kosinski et al. 2013; Quercia, Kosinski, Stillwell, & Crowcroft 2011; Youyou et al. 2015). Auch hier wurden moderne Modelle des statistischen Modellierens verwendet. Auf analoge Art schlägt Ihnen ein Algorithmus auf einer Webseite vor, für welche Produkte Sie sich noch interessieren könnten. Nicht immer ist der Algorithmus auch nur ansatzweise clever: Wer hat noch nicht erlebt, im Internet einen Gegenstand erworben zu haben, ein Fahrrad zum Beispiel, und danach noch wochenlang von Werbung für Fahrräder drangsaliiert zu werden (mir reicht *ein* Rad, Google).

Bei aller Dystopie, die mit der Digitalisierung zusammenhängt und mit Chinas „Kreditwürdigkeitspunkten“ bisher am konsequentesten weitergedacht wurde (Botsman 2017), es gibt auch Nutzen. Moderne Daten-Technologien stecken im Smartphone, in medizinischen Anwendungen und in deutschen Autos. Wer möchte auf diesen Fortschritt verzichten? Die wenigsten offenbar.

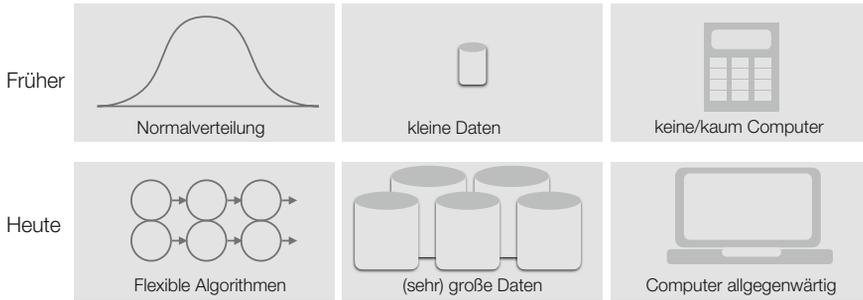
Diese Beispiele ließen sich noch länger mit Unterhaltungswert fortsetzen. Wie man den Fortschritt der Algorithmen auch einschätzt – wünschenswert, durchwachsen oder bedrohlich – man muss zum gleichen Schluss kommen: An der intensiven Beschäftigung mit dieser Technik kommen wir (jeder Einzelne) nicht vorbei. Ein Baustein, um die unheimliche Bedrohung durch panoptische, orwelleske Überwachung abzuwenden, ist das Verständnis der modernen Datentechnik. Gleichermäßen gilt: Um die offenbar gewaltigen ökonomischen Potenziale für die Unternehmen urbar zu machen, müssen wir die Technik verstehen (Brynjolfsson und McAfee 2016). Die Digitalisierung ist der bestimmende Trend des Wirtschaftslebens – wahrscheinlich (Vorsicht mit Vorhersagen); daher ist es beruflich, gesellschaftlich und politisch geboten, sich dem Algorithmozän zu stellen. Das heißt nicht, dass jeder Programmierer und Statistiker werden muss. Aber ein gewisses Grundverständnis sollte zum Bildungsstandard gehören.

Das Statistikcurriculum ist veraltet

Die Lehrpläne der Hochschulen geben sich von dramatischen Meldungen und neuen Technologien noch weitgehend unbescholten. Zumeist gilt in Lehrplänen für Statistik: Über den t -Test geht nichts. Der Wirklichkeit außerhalb der Alma Mater wird das kaum gerecht. Die Gründe für diese Gemächlichkeit können darin liegen, dass sich einige Hochschullehrer¹ mit neuen Technologien schwertun und mit Daten operieren (wollen), für die die alten Methoden wie der t -Test geeignet sind. Frei nach Max Planck kann man behaupten, dass alte Lehrmeinungen dann erst das Zeitliche segnen, wenn das auch die Professoren tun, die die Lehrmeinungen vertreten. Vielleicht liegt es auch schlicht daran, dass unser Alltag in den meisten Belangen wenig von der Digitalisierung und von Algorithmen berührt scheint: Beim Bäcker grüßt man wie seit Altvaterzeiten; deutsche Autos rollen vom Band wie seit dem Wirtschaftswunder; Schüler und Studenten lesen in ihren Büchern, wie es ihnen in Preußen, als das deutsche Schulwesen seine Anfänge fand, eingebläut wurde (Foucault 1994). Die Revolution der Daten ist kaum spürbar; sie fühlt sich weit weg an.

Aber ein Wechsel „unter den Talaren“ zeichnet sich ab: Statistiker mit Renommee rufen dazu auf, das Datenzeitalter im Unterricht einzuläuten (Cobb 2007; Hardin et al. 2015): Datenanalyse heute ist anders als gestern (s. folgende Abbildung). Immer mehr Lehrbücher zu moderner Statistik und Datenanalysen erscheinen, auch richtig gute (Baumer et al. 2017; z. B. James et al. 2013; McElreath 2015; Wickham und Golemund 2016). Bislang zumeist im englischen Sprachraum, aber es gibt auch zunehmend mehr deutschsprachige Bücher (z. B. Wickham und Golemund (2017)).

¹ Aus Gründen der Lesbarkeit wird in diesem Buch das generische Maskulinum („der Leser“) verwendet; immer sind alle Geschlechter gleichermaßen gemeint.



Das vorliegende Buch versucht, einen Teil der Lücke im deutschsprachigen Raum zu schließen. Sie werden in diesem Buch die grundlegenden Ideen der modernen Datenanalyse lernen. Ziel ist es, Sie – in Grundzügen – mit moderner Statistik vertraut zu machen.²

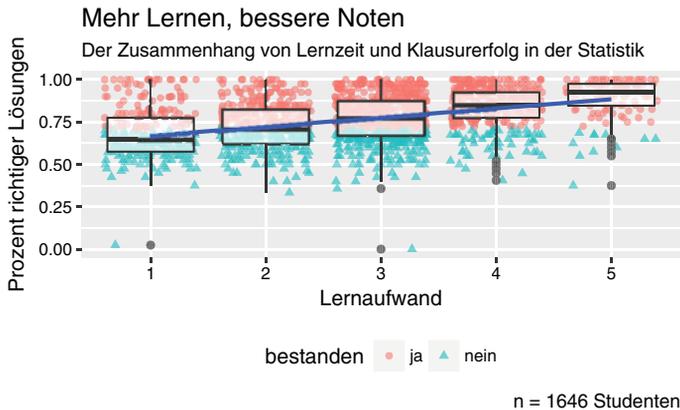
Lernziele

Zielgruppe dieses Buches sind Einsteiger; Formeln und mathematische Hintergründe sucht man meist vergebens. Das liegt zum einen daran, dass keine oder kaum Vorkenntnisse in Datenanalyse vorausgesetzt werden. Zum anderen beruht das Buch auf einem didaktischen Ansatz, der das *Tun* vor das *Wissen* setzt. Das bedeutet nicht, dass Wissen geringer geschätzt würde als Handeln. Vielmehr steht dahinter die Idee, dass es dem Verstehen und dem *statistical thinking* (C. J. Wild und Pfannkuch 1999) hilft, sich frühzeitig mit dem Ausprobieren auseinanderzusetzen. Der Werkzeug- oder Problemlöse-Charakter des Denkens steht im Mittelpunkt des Lernansatzes dieses Buches (vgl. Trilling und Fadel (2012)). Kein Inhalt dieses Buches bleibt ohne Umsetzung, ohne Anwendung; es ist ein Buch für Praktiker. Wer eine tiefere, mathematisch anspruchsvollere Einführung sucht, sei an das exzellente Buch von Hastie et al. (2013) verwiesen.

Reines Lesen dieses Buches wird dem Anfänger in etwa so viel bringen wie die Lektüre einer Schwimmfibel. Umgekehrt ist Üben die Grundlage für Fortschritt in der Kunst der Datenanalyse (s. folgende Abbildung).³ Nutzen Sie die Übungsangebote: die reichhaltige R-Syntax, die Daten, die Aufgaben und die Verweise zu weiterführender Literatur. Wesentlich ist das Durcharbeiten der Syntax-Beispiele. Die Kapitel sind zum Teil in sich abgeschlossen; die Grundlagen (bis einschließlich Kap. 7) werden durchgängig benötigt. Eine Lektüre von vorne nach hinten ist ratsam, aber nicht zwangsläufig nötig, gerade für fortgeschrittene Leser. Vergleichsweise schwierig sind die Kap. 28 und 29.

² Statistiker, die diesem Buch als Vorbild Pate standen, sind: Roger D. Peng: <http://www.biostat.jhsph.edu/~rpeng/>, Hadley Wickham: <http://hadley.nz>, Jennifer Bryan: <https://github.com/jennybc>.

³ Die Abbildung zeigt den Zusammenhang von Klausurerfolg und Vorbereitungsaufwand. Man sieht, dass der Klausurerfolg (Y-Achse) tendenziell steigt, wenn der Vorbereitungsaufwand (X-Achse) steigt.



Nach der Lektüre dieses Buches können Sie:

- den Ablauf eines Projekts aus der Datenanalyse in wesentlichen Schritten nachvollziehen
- Daten aufbereiten und ansprechend visualisieren
- Inferenzstatistik anwenden und kritisch hinterfragen
- klassische Vorhersagemethoden (Regression) anwenden
- moderne Methoden der angewandten Datenanalyse anwenden (z. B. Textmining)
- (wirtschaftliche) Fragestellungen mittels datengetriebener Vorhersagemodelle beantworten.

Zur Didaktik

Im Gegensatz zu anderen vergleichbaren Kursen steht hier die Umsetzung mit R (R Core Team 2018) im Vordergrund. Dies hat pragmatische Gründe: Möchte man Daten einer statistischen Analyse unterziehen, so muss man sie zumeist erst aufbereiten, und zwar oft mühselig. Selten kann man den Luxus genießen, einfach „nur“, nach Herzenslust sozusagen, ein Feuerwerk an multivariater Statistik abzubrennen. Zuvor gilt es, die Daten umzuformen, zu prüfen und zusammenzufassen. Für beide Anforderungen ist R bestens geeignet. Dem Teil des Aufbereitens der Daten ist hier ausführlich Rechnung getragen. Außerdem spielt in diesem Kurs die Visualisierung von Daten eine große Rolle. Ein Grund ist, dass Menschen bekanntlich Augentiere sind. Zum anderen bieten Diagramme bei umfangreichen Daten Einsichten, die sonst leicht wortwörtlich übersehen würden.

Lovett und Greenhouse (2000) leiten aus der kognitiven Theorie fünf Prinzipien zur Didaktik des Statistikerunterrichts ab; nach diesen Prinzipien ist dieses Buch ausgerichtet. (1) *Menschen lernen am meisten durch das und von dem, was sie selber ausprobieren*: Das Selber-Tun steht im Zentrum dieses Buches. (2) *Wissen ist situiert, kontextspezifisch*: Im Unterricht bzw. in einem Buch sollte daher lebensnah und alltagsrelevantes Wissen

vermittelt werden; die Beispiele und Methoden dieses Buches sind aus typischen oder verbreiteten Fragestellungen des Wirtschaftslebens entnommen. (3) *Direktes Feedback verbessert das Lernergebnis*: Das sofortige Ausprobieren anhand der R-Syntax gibt unmittelbares Feedback, ob ein Plan aufgegangen ist. Wie beim Jonglieren: Wenn ein Ball zu Boden fällt, weiß man, es ist ein Fehler passiert. Ähnlich verhält es sich, wenn R nichts oder Kauderwelsch ausspuckt. (4) *Lernen geschieht beim Verbinden von Bekanntem mit Neuem*: Der Sprachduktus ist informell, da viele Themen gerade in Bereichen nahe der Mathematik nicht wegen des Inhalts, sondern wegen der Formalisierung kompliziert werden. Freilich setzt man mit informeller Sprache Genauigkeit und Detailtiefe aufs Spiel. Da es sich aber um ein für die avisierte Leserschaft neues Thema handelt, neigt sich die Waage hier zugunsten der informellen, intuitiven Herangehensweise. Für fortgeschrittene Leser ist dieses Buch daher weniger geeignet. (5) *Die mentale Belastung sollte ausgewogen sein*: Das Buch beginnt mit grundlegenden Themen; die vergleichsweise schwierigen warten gen Ende. Jedes Kapitel behandelt nur ein Thema, um geistige Ressourcen effektiv zu nutzen. Die „R-Philosophie“ dieses Buches orientiert sich am „Tidyverse“ (vgl. Wickham und Grolemund (2017)); alle Kapitel und alle R-Syntax sind diesem Paradigma verhaftet. Sie werden schnell den ähnlichen Aufbau der Syntax in allen Kapiteln erkennen. Erfahrenen R-Programmierern wird der ausgiebige Gebrauch der „Pfeife“ aus `magrittr` auffallen; genauso wie der ausgiebige Gebrauch von `dplyr` und anderen Figuren aus dem „Tidyverse“.

Icons

R spricht zu Ihnen; sie versucht es jedenfalls in diesem Buch, und zwar mit folgenden Icons (Fonticons 2018).



R-Pseudo-Syntax: An vielen Stellen dieses Buches findet sich R-Syntax. Neue oder kompliziertere Syntax ist Zeile für Zeile ins Deutsche übersetzt.



Achtung, aufgepasst: Schwierige, merkwürdige oder fehlerträchtige Stellen sind mit diesem Symbol markiert.



Übungsaufgaben: In jedem Kapitel finden sich Übungsaufgaben. Auf diese wird mit diesem Icon verwiesen oder die Übungen sind in einem Abschnitt mit einschichtigem Titel zu finden.

Hinweise

Kunstwerke (Bilder) sind genau wie Standard-Literatur im Text zitiert, die verwendeten R-Pakete nur im Anhang. Alle Werke (auch Daten und Software) finden sich im Literaturverzeichnis. Dieses Buch wurde mit dem R-Paket `bookdown` basierend auf R (R Core Team 2018) in RStudio (RStudio 2018) geschrieben. `bookdown` basiert wiederum u. a. auf den R-Paketen `knitr` und `rmarkdown`. Norman Markgrafs Typografie-Paket hat den

Textsatz geschliffen (Markgraf 2018). Diese Pakete stellen großartige Funktionalität zur Verfügung, kostenlos und quelloffen.

Ein Hinweis für Dozenten: Der Text enthält an vielen Stellen Abbildungen oder kurze Zusammenfassungen, die sich für Folien eignen. Getreu der Philosophie, nach der Folien keine Konkurrenz zum Buch sind, sondern nur Wesentliches illustrieren, sollten vor allem diese Bilder in einen möglichen Foliensatz für eine Lehrveranstaltung eingehen. Das hat den charmanten Vorteil, dass die Erstellung der Folien einfach von der Hand geht. Im Anhang findet sich ein Vorschlag für das Curriculum einer Einführungsveranstaltung (s. Abschn. A.2.4).

Zu diesem Buch gibt es eine Webseite: <https://sebastiansauer.github.io/modar/>. Dort finden sich die Abbildungen, die Syntax und mehr. Feedback, Fragen und Hinweise können Sie dort auch einstellen.

Danke

Dieses Buch ist mit der Hilfe vieler und dank glücklicher Umstände entstanden. Von meinen Kollegen mit ihren Hinweisen, ihrem Ansporn, ihrem Lob und ihrer Kritik habe ich vielfältig profitiert; das hat dieses Buch geschliffen. Herausgreifen möchte ich Karsten Lübke, von dem ich viel gelernt habe. Norman Markgraf hat sein Wissen großzügig mit mir geteilt; für viele, gerade technische Hilfestellungen bin ich ihm dankbar. Weiter danke ich meiner Frau Sabrina, Christoph Kurz, Felix Bauer, Moritz Körber und Oliver Gansser für ihr Feedback zum Manuskript. Der Austausch mit den Kollegen vom ifes-Institut war fruchtbar und schwungvoll. Dieses Buch baut vielfach auf dieser Zusammenarbeit auf. Ich danke der Hochschulleitung der FOM für ihre Unterstützung und den Vorschuss an Vertrauen, was mein Arbeiten freudvoll machte; Gleiches gilt für den Dekan der Wirtschaftspsychologie an der FOM, Christoph Berg, der den Weg ebnete für neue Ideen, wie sie auch in diesem Buch Eingang fanden. Mein Dank gilt außerdem Kai Stumpp für die Begleitung bei der Erstellung des Buches. Nicht zuletzt danke ich meinen Studenten; ich weiß nicht, wer von wem mehr gelernt hat: sie von mir oder ich von ihnen, dank vieler Fragen und Hinweise. Vermeintlich „dumme“ Fragen zielen oft in die Mitte des Wesentlichen und die meisten komplexen Dinge kann man in einfachen Worten erklären, wenn man sie verstanden hat. Ohne umfangreiche Open-Source-Software wäre dieses Buch nicht entstanden; viele Menschen haben unentgeltlich mitgewirkt. Solcher Reichtum verblüfft mich immer wieder.

Ich hoffe, dass Sie mit diesem Buch einiges Handwerkszeug der modernen Datenanalyse lernen; dass Sie Gefallen an der „Kunst und Wissenschaft“ der Datenanalyse finden. Für Ihre Anregungen, Hinweise zu Fehlern und Ideen bin ich dankbar; am besten stellen Sie sie hier ein: <https://github.com/sebastiansauer/modar/issues>.

Inhaltsverzeichnis

Teil I Rahmen

1	Statistik heute	3
1.1	Datenanalyse, Statistik, Data Science und Co.	4
1.2	Wissensgebiete der Datenanalyse	6
1.3	Einige Grundbegriffe	8
1.4	Signal und Rauschen	9
2	Hallo, R	13
2.1	Eine kurze Geschichte von R	13
2.2	Warum R? Warum, R?	15
2.2.1	Warum R?	15
2.2.2	Warum, R?	17
3	R starten	21
3.1	R und RStudio installieren	21
3.2	Pakete	23
3.2.1	Pakete von CRAN installieren	23
3.2.2	Pakete installieren vs. Pakete starten (laden)	24
3.2.3	Pakete wie <code>pradadata</code> von Github installieren	25
3.3	Hilfe! R startet nicht!	25
3.4	Zuordnung von Paketen zu Befehlen	27
3.5	R-Skript-Dateien	29
3.6	Daten	29
3.6.1	Datensätze aus verschiedenen R-Paketen	29
3.6.2	Datensätze aus dem R-Paket <code>pradadata</code>	30
3.7	Grundlagen der Arbeit mit RStudio	30
3.7.1	Das Arbeitsverzeichnis	31
3.7.2	RStudio-Projekte	32
3.8	Hier werden Sie geholfen	33
3.8.1	Wo finde ich Hilfe?	33
3.8.2	Einfache, reproduzierbare Beispiele (ERBies)	33

4	Erstkontakt	37
4.1	R ist pingelig	37
4.2	Variablen zuweisen und auslesen	38
4.3	Funktionen aufrufen	39
4.4	Logische Prüfungen	40
4.5	Vektorielle Funktionen	42
4.6	Literaturempfehlungen	43

Teil II Daten einlesen

5	Datenstrukturen	47
5.1	Überblick über die wichtigsten Objekttypen	47
5.2	Objekttypen in R	49
5.2.1	Vektoren	49
5.2.2	Faktoren	51
5.2.3	Listen	53
5.2.4	Matrizen und Arrays	53
5.2.5	Dataframes	53
5.3	Daten auslesen und indizieren	55
5.3.1	Reine Vektoren	55
5.3.2	Matrizen und Arrays	57
5.3.3	Listen	57
5.3.4	Dataframes	59
5.4	Namen geben	60
6	Datenimport und -export	63
6.1	Daten in R importieren	63
6.1.1	Excel-Dateien importieren	64
6.1.2	Daten aus R-Paketen importieren	64
6.1.3	Daten im R-Format laden	65
6.1.4	CSV-Dateien importieren	65
6.2	Textkodierung	68
6.3	Daten exportieren	69

Teil III Daten aufbereiten

7	Datenjudo	75
7.1	Daten aufbereiten mit <code>dplyr</code>	77
7.2	Zentrale Bausteine von <code>dplyr</code>	78
7.2.1	Zeilen filtern mit <code>filter()</code>	78
7.2.2	Fortgeschrittene Beispiele für <code>filter()</code>	80

7.2.3	Spalten wählen mit <code>select()</code>	81
7.2.4	Zeilen sortieren mit <code>arrange()</code>	82
7.2.5	Einen Datensatz gruppieren mit <code>group_by()</code>	84
7.2.6	Eine Spalte zusammenfassen mit <code>summarise()</code>	86
7.2.7	Zeilen zählen mit <code>n()</code> und <code>count()</code>	89
7.3	Die Pfeife	91
7.4	Spalten berechnen mit <code>mutate()</code>	93
7.5	Bedingte Analysen mit den Suffixen von <code>dplyr</code>	96
7.5.1	Suffix <code>_if</code>	96
7.5.2	Suffix <code>_all</code>	97
7.5.3	Suffix <code>_at</code>	97
7.6	Tabellen zusammenführen (<code>join</code>)	99
8	Deskriptive Statistik	103
8.1	Univariate Statistik	104
8.1.1	Deskriptive Statistik mit <code>mosaic</code>	107
8.1.2	Deskriptive Statistik mit <code>dplyr</code>	108
8.1.3	Relative Häufigkeiten	109
8.2	Korrelationen berechnen	112
9	Praxisprobleme der Datenaufbereitung	117
9.1	Fehlende Werte	118
9.1.1	Ursachen von fehlenden Werten	118
9.1.2	Auf fehlende Werte prüfen	119
9.1.3	Umgang mit fehlenden Werten	119
9.1.4	Fälle mit fehlenden Werten löschen	119
9.1.5	Fehlende Werte einer Spalte zählen	121
9.1.6	Fehlende Werte ersetzen	122
9.1.7	-99 in NA umwandeln	124
9.2	Datenanomalien	125
9.2.1	Doppelte Fälle löschen	125
9.2.2	Nach Anomalien suchen	126
9.2.3	Ausreißer identifizieren	127
9.2.4	Hochkorrelierte Variablen finden	128
9.2.5	Quasi-Konstante finden	129
9.2.6	Auf Normalverteilung prüfen	130
9.3	Daten umformen	130
9.3.1	Aufgeräumte Dataframes	130
9.3.2	Langes vs. breites Format	131
9.3.3	z-Standardisieren	133
9.3.4	Spaltennamen ändern	134
9.3.5	Variablentypen ändern	135

9.4	Werte umkodieren und partitionieren	136
9.4.1	Umkodieren und partitionieren mit <code>car::recode()</code>	137
9.4.2	Einfaches Umkodieren mit einer Logik-Prüfung	138
9.4.3	Binnen mit <code>cut()</code>	139
9.5	Vektoren zu Skalaren zusammenfassen	141
9.5.1	Mittelwerte pro Zeile berechnen	141
9.5.2	Beliebige Statistiken pro Zeile berechnen mit <code>rowwise()</code>	142
10	Fallstudie: Datenjudo	145
10.1	Deskriptive Statistiken zu den New Yorker Flügen	146
10.2	Visualisierungen zu den deskriptiven Statistiken	149
10.2.1	Maximale Verspätung	149
10.2.2	Durchschnittliche Verspätung	151
10.2.3	Korrelate der Verspätung	152
 Teil IV Daten visualisieren		
11	Datenvisualisierung mit ggplot2	157
11.1	Einstieg in ggplot2	158
11.1.1	Ein Bild sagt mehr als 1000 Worte	158
11.1.2	Diagramme mit ggplot2 zeichnen	158
11.1.3	Die Anatomie eines Diagramms	159
11.1.4	Schnell Diagramme erstellen mit <code>qplot()</code>	162
11.1.5	ggplot-Diagramme mit <code>mosaic</code>	164
11.2	Häufige Arten von Diagrammen (Geomen)	166
11.2.1	Eine kontinuierliche Variable – Histogramme und Co.	166
11.2.2	Zwei kontinuierliche Variablen	167
11.2.3	Eine oder zwei nominale Variablen	169
11.2.4	Zusammenfassungen zeigen	174
11.3	Die Gefühlswelt von ggplot2	178
11.4	<code>ggplot()</code> , der große Bruder von <code>qplot()</code>	179
12	Fortgeschrittene Themen der Visualisierung	187
12.1	Farbwahl	187
12.1.1	Die Farben von Cynthia Brewer	188
12.1.2	Die Farben von Wes Anderson	190
12.1.3	Viridis	193
12.2	ggplot2-Themen	194
12.2.1	Schwarz-Weiß-Druck	195
12.3	Interaktive Diagramme	197
12.3.1	Plotly	197
12.3.2	Weitere interaktive Diagramme	198

13	Fallstudie: Visualisierung	201
13.1	Umfragedaten visualisieren mit „likert“	202
13.2	Umfragedaten visualisieren mit ggplot	203
13.2.1	Daten aufbereiten	203
13.2.2	Daten umstellen	204
13.2.3	Diagramme für Anteile	204
13.2.4	Rotierte Balkendiagramme	207
13.2.5	Text-Labels	208
13.2.6	Diagramm beschriften	211
13.2.7	Balken mit Häufigkeitswerten	211
13.2.8	Sortieren der Balken	212
14	Geovisualisierung	215
14.1	Kartendaten	216
14.1.1	Geo-Daten der deutschen Verwaltungsgebiete	216
14.1.2	Daten der Wahlkreise	218
14.2	Unterschiede in Kartensegmenten visualisieren	219
14.2.1	Karte der Wahlkreise gefärbt nach Arbeitslosigkeit	219
14.2.2	Wahlergebnisse nach Wahlkreisen	220
14.2.3	Zusammenhang von Arbeitslosigkeit und AfD-Wahlergebnis	222
14.2.4	Ein komplexeres Modell	223
14.3	Weltkarten	224
14.3.1	rworldmap	224
14.3.2	rworldmap mit geom_sf	227
14.4	Anwendungsbeispiel: Konkordanz von Kulturwerten und Wohlbefinden	229
14.5	Interaktive Karten	234
14.5.1	Karten mit „leaflet“	234
14.5.2	Karten mit googleVis	235

Teil V Modellieren

15	Grundlagen des Modellierens	245
15.1	Was ist ein Modell? Was ist Modellieren?	246
15.2	Abduktion als Erkenntnisfigur im Modellieren	248
15.3	Ein Beispiel zum Modellieren in der Datenanalyse	250
15.4	Taxonomie der Ziele des Modellierens	251
15.5	Die vier Schritte des statistischen Modellierens	254
15.6	Einfache vs. komplexe Modelle: Unter- vs. Überanpassung	255

15.7	Bias-Varianz-Abwägung	256
15.8	Trainings- vs. Test-Stichprobe	257
15.9	Resampling und Kreuzvalidierung	259
15.10	Wann welches Modell?	260
15.11	Modellgüte	260
15.11.1	Modellgüte in numerischen Vorhersagemodellen	261
15.11.2	Modellgüte bei Klassifikationsmodellen	261
15.12	Der Fluch der Dimension	262
16	Inferenzstatistik	267
16.1	Wozu Inferenzstatistik?	268
16.2	Der p -Wert	269
16.2.1	Was sagt der p -Wert?	269
16.2.2	Der zwielichtige Statistiker – ein einführendes Beispiel zur Inferenzstatistik	271
16.2.3	Von Männern und Päpsten – Was der p -Wert nicht sagt	274
16.2.4	Der p -Wert ist eine Funktion der Stichprobengröße	275
16.2.5	Mythen zum p -Wert	276
16.3	Wann welcher Inferenztest?	277
16.4	Beispiele für häufige Inferenztests	278
16.4.1	χ^2 -Test	278
16.4.2	t-Test	280
16.4.3	Einfache Varianzanalyse	281
16.4.4	Korrelationen (nach Pearson) auf Signifikanz prüfen	282
16.4.5	Regression	283
16.4.6	Wilcoxon-Test	284
16.4.7	Kruskal-Wallis-Test	284
16.4.8	Shapiro-Test	285
16.4.9	Logistische Regression	285
16.4.10	Spearman's Korrelation	286
16.5	Alternativen zum p -Wert	286
16.5.1	Konfidenzintervalle	286
16.5.2	Effektstärke	289
16.5.3	Power-Analyse	292
16.5.4	Bayes-Statistik	293
17	Simulationsbasierte Inferenz	301
17.1	Stichproben, Statistiken und Population	301
17.2	Die Stichprobenverteilung	304
17.3	Der Bootstrap	308
17.4	Nullhypothesen auf Signifikanz testen	311

Teil VI Geleitetes Modellieren

18	Lineare Modelle	321
18.1	Die Idee der klassischen Regression	321
18.2	Modellgüte	324
18.2.1	Mittlere Quadratfehler	325
18.2.2	R-Quadrat (R^2)	325
18.3	Die Regression an einem Beispiel erläutert	327
18.4	Überprüfung der Annahmen der linearen Regression	329
18.5	Regression mit kategorialen Prädiktoren	331
18.6	Multiple Regression	333
18.7	Interaktionen	335
18.8	Prädiktorenrelevanz	337
18.9	Anwendungsbeispiel zur linearen Regression	339
18.9.1	Overfitting	339
18.9.2	Konfidenzintervalle der Parameter	341
19	Klassifizierende Regression	345
19.1	Normale Regression für ein binäres Kriterium	346
19.2	Die logistische Funktion	347
19.3	Interpretation des Logits	350
19.4	Kategoriale Prädiktoren	351
19.5	Multiple logistische Regression	352
19.6	Modellgüte	353
19.6.1	Vier Arten von Ergebnissen einer Klassifikation	353
19.6.2	Kennzahlen der Klassifikationsgüte	355
19.7	Vorhersagen	356
19.8	ROC-Kurven und Fläche unter der Kurve (AUC)	357
19.8.1	Cohens Kappa	360
20	Fallstudie: Titanic	365
20.1	Explorative Analyse	366
20.1.1	Univariate Häufigkeiten	366
20.1.2	Bivariate Häufigkeiten	367
20.2	Inferenzstatistik	368
20.2.1	χ^2 -Test	368
20.2.2	Effektstärke	369
20.2.3	Logistische Regression	373
21	Baumbasierte Verfahren	377
21.1	Entscheidungsbäume	378
21.1.1	Einführendes Beispiel	378
21.1.2	Tuningparameter	383

21.2	Entscheidungsbäume mit <code>caret</code>	384
21.2.1	Vorhersagegüte	388
21.3	Der Algorithmus der Entscheidungsbäume	391
21.4	Regressionsbäume	391
21.5	Stärken und Schwächen von Bäumen	391
21.6	Bagging	393
21.7	Grundlagen von Random Forests	394
21.7.1	Grundlagen	394
21.8	Variablenrelevanz bei Baummodellen	398
22	Fallstudie: Kreditwürdigkeit mit <code>caret</code>	401
22.1	Zwei Arten der prädiktiven Modellierung	402
22.2	Daten aufbereiten	403
22.2.1	Fehlende Werte	403
22.2.2	Trainings- und Test-Sample aufteilen	403
22.2.3	Variablen ohne Varianz	404
22.2.4	Hochkorrelierte Variablen entfernen	405
22.2.5	Parallele Verarbeitung	406
22.3	Modelle anpassen	407
22.3.1	Kreuzvalidierung	407
22.3.2	Modell im Trainings-Sample anpassen mit <code>train()</code>	407
22.3.3	Ein einfaches Modell	408
22.3.4	Random Forest	411
22.3.5	Support Vector Machines	414
22.3.6	Penalisierte lineare Modelle	416
22.4	Modellgüte bestimmen	418
22.4.1	Modellgüte in der Test-Stichprobe	418
22.4.2	Modellgüte in der Kreuzvalidierung	421
22.5	Wichtigkeit der Prädiktoren bestimmen	426
22.5.1	Modellunabhängige Variablenwichtigkeit für Klassifikation	427
22.5.2	Modellunabhängige Prädiktorenrelevanz bei numerischen Vorhersagen	429
22.5.3	Modellabhängige Variablenwichtigkeit	430
 Teil VII Ungeleitetes Modellieren		
23	Clusteranalyse	437
23.1	Grundlagen der Clusteranalyse	437
23.1.1	Intuitive Darstellung der Clusteranalyse	438
23.1.2	Euklidische Distanz	440
23.1.3	k-Means-Clusteranalyse	442

23.2	Beispiel für eine einfache Clusteranalyse	443
23.2.1	Distanzmaße berechnen	443
23.2.2	Fallstudie: kmeans für den Extraversionsdatensatz	444
24	Textmining	449
24.1	Grundlegende Analyse	450
24.1.1	Tidytext-Datframes	450
24.1.2	Regulärausdrücke	453
24.1.3	Textdaten einlesen	455
24.1.4	Worthäufigkeiten auszählen	456
24.1.5	Visualisierung	457
24.2	Sentimentanalyse	459
25	Fallstudie: Twitter-Mining	463
25.1	Zum Einstieg: Moderne Methoden der Sentimentanalyse	464
25.2	Grundlagen des Twitter-Minings	465
25.2.1	Authentifizierung bei der Twitter-API	466
25.2.2	Hashtags und Nutzer suchen	467
25.2.3	Tweets einer Nutzermenge auslesen	469
25.2.4	Aufbau einer Tweets-Datenbank	471
 Teil VIII Kommunizieren		
26	RMarkdown	475
26.1	Forderungen an Werkzeuge zur Berichterstellung	476
26.2	Start mit RMarkdown	478
26.3	RMarkdown in Action	480
26.4	Aufbau einer Markdown-Datei	482
26.5	Syntax-Grundlagen von Markdown	483
26.6	Tabellen	484
26.7	Zitieren	487
26.8	Format-Vorlagen für RMarkdown	489
 Teil IX Rahmen 2		
27	Projektmanagement am Beispiel einer Fallstudie	495
27.1	Was ist Populismus?	496
27.2	Forschungsfrage und Operationalisierung	497
27.3	Emotionslexikon	498
27.4	Daten, Stichprobe und Analysekontext	499

27.5	Prozess der Datenanalyse	499
27.6	Zentrale Ergebnisse	501
27.7	Projektmanagement	504
27.7.1	Gliederung eines Projektverzeichnis	504
27.7.2	Faustregeln zur Struktur eines Projekts	505
27.7.3	Versionierung mit Git	508
28	Programmieren mit R	511
28.1	Funktionen schreiben	511
28.2	Wiederholungen	514
28.2.1	Wiederholungen für Elemente eines Vektors	515
28.2.2	Wiederholungen für Spalten eines Dataframes	516
28.2.3	Dateien wiederholt einlesen	518
28.2.4	Anwendungsbeispiele für map	520
28.2.5	Einige Rechtschreibregeln für map ()	523
28.3	Defensives Programmieren	523
29	Programmieren mit dplyr	527
29.1	Wie man mit dplyr nicht sprechen darf	527
29.2	Standard-Evaluation vs. Non-Standard-Evaluation	528
29.3	NSE als Backen	530
29.4	Wie man Funktionen mit dplyr-Verben schreibt	534
29.5	Beispiele für NSE-Funktionen	537
29.5.1	Funktionen für ggplot	539
Anhang A		541
Literatur		547
Sachverzeichnis		559

Der Autor

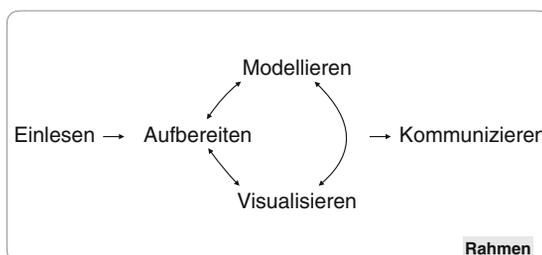
Professor Dr. habil. Sebastian Sauer arbeitet als Hochschullehrer für Wirtschaftspsychologie an der FOM und versteht sich als Data Scientist. Datenanalyse mit R ist sein zentrales Thema im Moment. Neben dem „Wie“ der Datenanalyse beschäftigen ihn die Grenzen sowie die Gefahren, die die moderne Datenwissenschaft mit sich bringt. Außerdem engagiert er sich für das Thema Open Science und interessiert sich für die Frage, wie die Psychologie zur Klärung von Problemen mit gesellschaftlicher Relevanz beitragen kann. Sein Blog <https://data-se.netlify.com/> dient ihm als Notizbuch sich entwickelnder Gedanken. Data Science für die Wirtschaft bietet er auf <https://www.data-divers.com/> an.

Teil I

Rahmen

Datenanalyse, praktisch betrachtet, kann man in fünf Schritte einteilen (Wickham und Grolemund 2017), s. Abb. 1.1. Analog zu diesem Modell der Datenanalyse ist dieses Buch aufgebaut. Zuerst muss man die Daten *einlesen*, die Daten also in R (oder einer anderen Software) verfügbar machen (laden). Fügen wir hinzu: In schöner Form verfügbar machen; das man nennt auch Tidy Data (hört sich cooler an). Sobald die Daten in geeigneter Form in R geladen sind, folgt das *Aufbereiten*. Das beinhaltet das Zusammenfassen, Umformen oder Anreichern der Daten, je nach Bedarf. Ein nächster wesentlicher Schritt ist das *Visualisieren* der Daten. Ein Bild sagt bekanntlich mehr als tausend Worte. Schließlich folgt das *Modellieren* oder das Prüfen von Hypothesen: Man überlegt sich, wie sich die Daten erklären lassen könnten. Zu beachten ist, dass diese drei Schritte – Aufbereiten, Visualisieren, Modellieren – keine starre Abfolge sind, sondern eher ein munteres Hin- und-Her-Springen, ein aufeinander aufbauendes Abwechseln. Der letzte Schritt ist das *Kommunizieren* der Ergebnisse der Analyse – nicht der Daten. Niemand ist an Zahlenwüsten interessiert; es gilt, spannende Einblicke zu vermitteln. Die Datenanalyse als solche ist in einen *Rahmen* eingebettet; das beinhaltet philosophische und technische Grundlagen. Entsprechend diesen fünf Schritten sowie dem einbettenden Rahmen ist dieses Buch in Teile gegliedert. Zu Beginn jedes Teiles ist ein Diagramm analog zu 1.1 dargestellt, um einen Überblick über den jeweiligen Schritt der Datenanalyse zu geben.

Abb. 1.1 Der Rahmen als Bestandteil der Datenanalyse





Lernziele

- Wissen, was Statistik ist, bzw. einige Aspekte einer Definition kennen
- Statistik zu Data Science und anderen verwandten Begriffen abgrenzen können
- Grundkonzepte wie Daten, Variable und Beobachtung definieren können
- Die Begriffe *Signal* und *Rauschen* in Verbindung bringen können
- Die Wissensgebiete der Datenanalyse aufzählen und erläutern können

1.1 Datenanalyse, Statistik, Data Science und Co.

Was ist Statistik? Eine Antwort dazu ist, dass Statistik die Wissenschaft von Sammlung, Analyse, Interpretation und Kommunikation von Daten mithilfe mathematischer Verfahren ist und zur Entscheidungshilfe beitragen sollte (Romeijn 2016; *The Oxford Dictionary of Statistical Terms* 2006). E. N. Brown und Kass (2009) rücken die Wahrscheinlichkeitsrechnung ins Zentrum einer Definition von Statistik oder „Statistik-Denken“. Cobb (2015, S. 3) spricht kurz und bündig von „thinking with and about data“ als dem Wesensmerkmal von Statistik.

Wie lässt sich Statistik von Datenanalyse abgrenzen? Tukey (1962) definiert Datenanalyse als den Prozess des Erhebens von Daten, ihrer Auswertung und Interpretation. Unschwer zu sehen, dass sich diese beiden Definitionen nur um die Betonung der stochastischen Modellierung, der Anwendung der Wahrscheinlichkeitsrechnung, unterscheiden. Betrachtet man Lehrbücher der Statistik (Bortz 2013; Freedman et al. 2007), so fällt der stärkere Fokus auf mathematische Ableitung und Eigenschaften von Objekten der Statistik ins Auge; Datenanalyse scheint einen stärkeren Anwendungsfokus zu haben (im Gegensatz zu einem mathematischen Fokus). Statistik wird häufig in die zwei Gebiete *deskriptive* und *inferierende* Statistik eingeteilt (vgl. Abb. 1.2). Erstere fasst viele Zahlen zusammen (s. Kap. 8), so dass wir den Wald statt vieler Bäume sehen. Eine *Statistik* bezeichnet dabei die zusammenfassende Kenngröße; eine prototypische Statistik ist der Mittelwert. Die Inferenzstatistik verallgemeinert von den vorliegenden Daten auf eine zugrunde liegende Grundgesamtheit (Population; s. Kap. 16). So zieht man etwa eine Stichprobe von einigen College-Studenten und schließt auf dieser Basis auf alle Menschen dieser Welt. Ein abenteuerlicher Schluss, aber leider kein seltener (Henrich et al. 2010). Da *Analyse von Daten* ein allgemeiner Begriff ist, der wenig mit bestimmten Methoden aufgeladen ist, wird Datenanalyse im Folgenden als gemeinsamer Kern aller einschlägigen Disziplinen oder Begrifflichkeiten verwendet; der Fokus ist dabei als angewandt gedacht. In diesem Buch werden Statistik und Datenanalyse im Folgenden lose als Synonyme betrachtet.

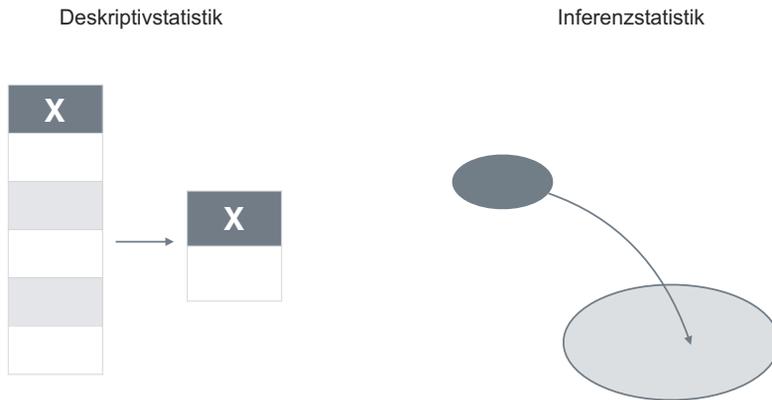


Abb. 1.2 Sinnbild für die Deskriptiv- und die Inferenzstatistik

Aufgabe der deskriptiven Statistik ist es primär, Daten prägnant zusammenzufassen. Aufgabe der Inferenzstatistik ist es, zu prüfen, ob Daten einer Stichprobe auf eine Grundgesamtheit verallgemeinert werden können.



Liegt der Schwerpunkt der Datenanalyse auf computerintensiven Methoden, so wird auch von *Data Science* gesprochen, wobei der Begriff nicht einheitlich verwendet wird (Hardin et al. 2015; Wickham und Golemund 2017). Einige Statistiker sehen *Data Science* als „Statistik am Computer“, und plädieren für ein „Neudenken“ der Statistik bzw. des Statistikerunterrichts, so dass Computermethoden eine zentrale Rolle spielen (Cobb 2015). Andere Statistiker wiederum grenzen *Data Science* von der Statistik ab, mit dem Argument, dass bei Ersterer Fragen der angewandten Informatik zentral sind, bei der Statistik nicht (Baumer et al. 2017). Die Popularität von *Data Science* ist dem Fortschritt in der Rechen- und Speicherkapazität der Computer zu verdanken. Heutzutage sind ganz andere Daten und Datenmengen verarbeitbar und von Interesse. Sicherlich bedingt die technische Machbarkeit auch, welche Forschungsfragen hoch im Kurs rangieren. Eine Echtzeit-Analyse von Twitter-Daten wäre bis vor einiger Zeit kaum möglich gewesen, da die Hardware nicht leistungsfähig genug war. (Wir übersehen hier geflissentlich, dass die Hardware heute immer noch an Grenzen kommt und dass es früher kein Twitter gab.) Die Datenmengen erfordern Arbeitsschritte wie z. B. Authentifizierung in die Twitter-Schnittstelle, wiederholtes Abfragen der Schnittstelle unter Beachtung der erlaubten Download-Obergrenzen, Umwandeln von einem Datenformat in ein anderes, Einlesen in eine Datenbank, Prüfung auf Programmfehler, Automatisierung des Prozesses, Verwendung mehrerer Rechenkerne, Bereinigung des Textmaterials von Artefakten, Aufbereiten der Daten und so weiter. Keine dieser Aufgaben war bei Statistikern vor 100 Jahren verbreitet. Da konnte man sich noch ganz auf die mathematischen Eigenschaften

des t-Tests konzentrieren. Die Verschiebung innerhalb der Datenanalyse spiegelt einfach die technische Entwicklung allgemein in der Gesellschaft wider. Das heißt nicht, dass heute jeder Programmierer sein muss; auch für die Datenanalyse nicht. Glücklicherweise gibt es eine Reihe von Werkzeugen, die die Handhabung der Daten einfacher macht. Diese sind hinzugekommen zum Handwerkskoffer des Datenschreiners. Nichtsdestotrotz ist ein grundlegendes Verständnis von computergestützter Datenverarbeitung wichtig und wird zunehmend wichtiger für die Datenanalyse.

Schließlich kursieren u. a. noch die Begriffe *Data Mining*, *maschinelles Lernen* und *statistisches Lernen*. Data Mining und maschinelles Lernen sind Begriffe, die eher in informatiknahen Gebieten verwendet werden; entsprechend sind die Themen mehr in Richtung Informatik verschoben. Die technische Repräsentation und technische Aspekte der Datenmanipulation (Data Warehouse, Datenbanken) werden von IT-affinen Autoren stärker betont als von Autoren, die nicht aus IT-nahen Fakultäten stammen. Beim statistischem Lernen stehen Konzepte und Algorithmen des Modellierens (s. Kap. 15) im Vordergrund. Auf der anderen Seite: Vergleicht man die Inhaltsverzeichnisse von Büchern aus allen diesen Bereichen, so stellt man eine große Überschneidung des Inhaltsverzeichnisses fest. Typische Themen in allen diesen Büchern sind (Bishop 2006; Han et al. 2011; James et al. 2013; Tan 2013) baumbasierte Verfahren (s. Kap. 21), Support Vector Machines, Dimensionsreduktion, Clusteranalyse, Regression (s. Kap. 18) oder Datenexploration (s. Kap. 11). Kurz: Die Gebiete überlappen einander beträchtlich; mal stehen IT-nahe Themen im Vordergrund, mal wird die Mathematik betont, mal die Anbindung an empirisches Forschen. Immer geht es darum, aus Daten Wissen zu generieren bzw. Entscheidungen datenbasiert und rational zu begründen.

1.2 Wissensgebiete der Datenanalyse

Egal, ob man von Datenanalyse, Statistik, Data Mining, maschinellem Lernen, Data Science oder statistischem Lernen spricht: In der Schnittmenge des Analysierens von Daten gleichen sich die Anforderungen. Wenn sich die Nuancen zwischen den Fachgebieten auch verschieben, so sind doch stets die folgenden Wissensgebiete gefragt:

1. Philosophische Grundlagen

Dazu gehören die Annahmen, die im Alltag meist unhinterfragt für bare Münzen genommen werden. Annahmen, die das Fundament des Gebäudes der Datenanalyse stützen. Ist dieses Fundament auf Sand gebaut, so ist nicht zu erwarten, dass das Gebäude seinen Zweck erfüllt. Zu den Grundlagenfragen gehört die Frage, was Wahrscheinlichkeit, Erkenntnis, Kausalität und Unendlichkeit sind. So wird zum Beispiel kontrovers diskutiert, ob Wahrscheinlichkeit besser als Grenzwert der relativen Häufigkeit, als subjektive Angelegenheit oder als Erweiterung der Aussagenlogik zu betrachten ist (Briggs 2016; Jaynes 2003; Keynes 2013; Rucker 2004). Wichtig ist weiter die Frage, was eine Messung genau ist, woran man erkennt, ob ein Variable quantitative Aus-

sagen erlaubt und woran man die Güte eines Messinstruments festmacht (Saint-Mont 2011).

2. Mathematisch-statistische Anwendungen

Zentrale Theorie für die Statistik oder für Wissenschaft allgemein ist die Wahrscheinlichkeitsrechnung bzw. die Theorie der Wahrscheinlichkeit (Jaynes 2003). In den meisten Lehrbüchern der Statistik, auch in den Einsteigerbüchern, findet sich eine mal schmalere, mal ausführlichere Einführung in Aspekte verschiedener Verteilungen, die gewisse Zufallsprozesse, berechenbare Zufallsprozesse, voraussetzen.¹ Häufig wird angenommen, dass sich eine bestimmte Variable nach einem bekannten stochastischen Modell verhält, so dass aus dem Modell Aussagen ableitbar sind. Besondere Berühmtheit hat die Normalverteilung erlangt, die wohl an keinem Studenten eines empirisch orientierten Faches vorbeigegangen ist.² Neben der Stochastik spielen aber noch weitere Felder der angewandten Mathematik eine Rolle; maßgeblich sind das die lineare Algebra und die Infinitesimalrechnung (J. D. Brown 2015). Ein eigener Zweig, der stark mit der Wahrscheinlichkeitslehre verbunden ist, ist die Bayes-Statistik (Wagenmakers et al. 2016).

3. Computerwissenschaftliche Anwendungen

Die zunehmende Digitalisierung der Gesellschaft macht vor der Datenanalyse keinen Halt. Im Gegenteil; es liegt in der Natur der Datenanalyse, computeraffin zu sein – sind doch Daten Gegenstand sowohl der Statistik als auch der Computerwissenschaft. War es vor einigen Jahren noch ausreichend oder beeindruckend, den Knopf für den t-Test zu kennen, sind die Anforderungen bei modernen Daten meist höher. Einlesen, Aufbereiten und Speichern können leicht den größten Zeitanteil einer Datenanalyse ausmachen. In einigen Anwendungen kommt noch der Anspruch dazu, dass die Analyse schnell gehen muss: Wir haben gerade viele Kunden auf der Webseite und müssen wissen, wem wir welches Produkt und welchen Preis vorschlagen müssen. Nein, wir wollen die Antwort nicht morgen, wir brauchen sie in ein paar Sekunden; Korrektur: jetzt. Ach ja, der Datensatz ist ein paar Terabyte³ groß.

4. Fach- und Branchenkenntnis

Möchte man die Zufriedenheit eines Kunden vorhersagen, so ist es hilfreich, etwas über die Ursachen von Kundenzufriedenheit zu kennen – also Wissen über den Gegenstand Kundenzufriedenheit zu haben. Wenn man schon die Ursachen nicht kennt, so ist zumindest Wissen über zusammenhängende Variablen (Korrelate) sinnvoll. Wenn ein Arzt aus Erfahrung weiß, was die Risikofaktoren einer Erkrankung sind, dann sollten diese Informationen in das statistische Modell einfließen. Sach- inkl. Branchenkennt-

¹ Zufall wird hier verstanden als ein nicht näher bekannter Prozess, dessen Ergebnisse zwar nicht sicher sind, aber doch ein gewisses Muster erwarten lassen.

² Micceri (1989) zeigt auf, dass Normalverteilungen seltener sind als gemeinhin angenommen. McElreath (2015) bietet einen gut verständlichen Einblick in die Informationsentropie; diese Darstellung zeigt, dass eine Normalverteilung eine konservative Annahme für eine Verteilung darstellt.

³ 1 Terabyte sind 10^{12} Byte.

nis ist zentral für gute (genaue) statistische Modelle (Shearer 2000). Wissen über den Sachgegenstand ist schon deshalb unerlässlich, weil Entscheidungen keine Frage der Statistik sind: Ob ein Kunde zufrieden ist mit einer Vorhersage durch ein statistisches Modell oder ein Patient gesund genug ist nach Aussage eines statistisches Modells, muss der Anwender entscheiden. Ob ein Betrugsversuch mit 90 %, 99 % oder 99.9 % Sicherheit erkannt werden soll, kann einen Unterschied machen – für einen bestimmten Anwender, in einer bestimmten Situation. Ein wichtiger, vielleicht der wichtigste Punkt der Datenanalyse ist es, Entscheidungen für Handlungen zu begründen. Daher muss die Datenanalyse immer wieder auf die Entscheidung und damit auf die Präferenz des Nutzers zurückgeführt werden.

1.3 Einige Grundbegriffe

Daten (die Einzahl *Datum* ist ungewöhnlich) kann man definieren als Informationen, die in einem Kontext stehen (Moore 1990), wobei eine numerische Konnotation mitschwingt. Häufig sind Daten in *Tabellen* bzw. tabellenähnlichen Strukturen gespeichert; die Excel-Tabelle ist der Prototyp davon. Tabellen, so wie sie hier verstanden werden, zeichnen sich dadurch aus, dass sie *rechteckig* sind und aus Zeilen und Spalten bestehen. Rechteckig impliziert, dass alle Zeilen gleich lang sind und alle Spalten gleich lang sind (die Tabelle muss aber nicht quadratisch sein). Knotenpunkte von Zeilen und Spalten heißen *Zellen* oder *Elemente*; die Zellen dürfen auch leer sein oder mit einem Symbol für „kein Wert vorhanden“ gefüllt sein.

Daten sind ein Produkt von *Variablen* (Merkmalen) und *Beobachtungseinheiten* (Fällen, Beobachtungen). Eine Tabelle mit ihren zwei rechtwinkligen Achsen der Zeilen und Spalten verdeutlicht das (s. Abschn. 9.3.1 und Abb. 9.6). Die Beobachtungseinheit ist das Objekt, das die untersuchten Merkmale aufweist. Oft sind es Personen, es können aber auch Firmen, Filme, Filialen oder Flüge sein. Ein Merkmal einer Beobachtungseinheit kann ihre Schuhgröße, der Umsatz, die Verspätung oder das Budget sein. Betrachten wir das Merkmal *Schuhgröße* der Beobachtungseinheit *Person S* und finden wir *46*, so ist *46* der *Wert* oder die *Ausprägung* dieses Merkmals dieser Beobachtungseinheit. Um nicht so viel schreiben zu müssen, wird der Wert der Beobachtungseinheit *i* in der Variablen *k* häufig als x_{ik} bezeichnet. Die Gesamtheit der verfügbaren und zusammengehörigen Daten eines Sachverhalts bezeichnen wir als *Datensatz*; meist ist es eine Stichprobe aus einer Population. Nehmen wir an, es gäbe nur zwei *verschiedene* Schuhgrößen: 36 und 46. Dann hat die Variable Schuhgröße zwei *Ausprägungen*.

Natürlich gibt es auch Daten, die sich nicht (so einfach) in das enge Korsett einer Tabelle pressen lassen; Textdokumente, Sprachdaten oder Bilder zum Beispiel. Nicht tabellarisierte Daten werden auch als *unstrukturiert*, Daten in Tabellenform als *strukturiert* bezeichnet. Rein mengenmäßig überwiegen unstrukturierte Daten in der Welt. Allerdings sind strukturierte Daten einfacher zu verarbeiten; wir werden uns auf diese konzentrieren.

1.4 Signal und Rauschen

Die Aufgabe der Wissenschaft – oder sogar jeglichen Erkennens – kann man als zweistufigen Prozess betrachten: erstens Signale (Phänomene) erkennen und zweitens diese dann erklären (Bogen und Woodward 1988; Silver 2012). *Signale erkennen* ist der Versuch, aus Daten ein Muster, Regularitäten, herauszulesen. Das impliziert, dass Daten dieses Muster nicht direkt offenbaren und nicht identisch mit dem Muster sind. Mit dem alten Bild, im Rauschen ein (leises) Geräusch, das Signal, herauszuhören, ist der Sachverhalt gut beschrieben (Haig 2014). Daten sind vergänglich, veränderlich, vorübergehend – und für den erkennenden Verstand ohne Interesse. Das Muster ist es, welches von alleinigem Interesse ist. Schält sich ein Muster heraus, ist es über die Zeit, Erhebungsmethode und Situation hinweg stabil, so hat man ein *Phänomen* identifiziert. Im Gegensatz zu Daten ist ein Phänomen unbeobachtbar; es ist die Abstrahierung der Gemeinsamkeit aus der Konkretheit der Daten. Ein Phänomen könnte sein, dass „Männer, die Windeln kaufen, auch Bier kaufen“. Ist ein Phänomen identifiziert, so ist die wichtigste Frage zumeist, was die Ursache des Phänomens ist. Wissenschaftliche und Alltagstheorien untersuchen Phänomene, nicht Daten. Warum ist es wichtig, Ursachen von Phänomen zu kennen? Ein Grund ist, dass man das Auftreten eines Phänomens beeinflussen kann, wenn man seine Ursache kennt. Überspitzt formuliert: Die Ursache der Entzündung sind Bakterien? Entferne die Bakterien, und die Entzündung klingt ab.

Abb. 1.3 stellt den Unterschied (und den Zusammenhang) von Rauschen und Signal dar. Der linke Teil des Diagramms zeigt die Körpergröße einer Reihe von Frauen und Männern (und damit von zwei Variablen). Die Beobachtungseinheit, als schwarzer Punkt dargestellt, ist eine Person. Wie man sieht, unterscheiden sich die Körpergrößen der Personen; einige sind größer, andere kleiner. Auch zwischen den *Geschlechtern* gibt es Unterschiede. Unser Auge ist schnell mit der Erkenntnis des Musters, dass „Männer größer sind als Frauen“. Nicht alle Männer sind größer als alle Frauen; einige Frauen sind größer als einige Männer. Aber in den meisten Fällen gilt: Der Mann ist größer als die Frau. Anders betrachtet: Der „mittlere“ Mann ist größer als die „mittlere“ Frau (als Quadrat bzw. Linie dargestellt im rechten Teil der Abb. 1.4). Als Beobachtung in den Daten existiert aber weder die mittlere Frau noch der mittlere Mann; wir erkennen dies als Phänomen in den Daten bzw. aus den Daten heraus.

Das Rauschen in den Daten kann vielerlei Ursachen haben: Messfehler, Besonderheiten der ausgewählten Merkmalsträger oder der Situation. Zufall ist keine Ursache im strengen Sinne des Wortes; in der Regel wird dieser Begriff verwendet, wenn man die wahre Ursache nicht kennt.⁴ Experimentieren ist nichts anderes als die Kunst, Rauschen *vor der Messung* zu verringern. Genauer gesagt solches Rauschen zu verringern, welches das Signal sonst überlagert. Analog kann man sagen, dass Datenanalyse das Ziel hat, Rau-

⁴ Diesem Gedanken hinterliegt ein deterministisches Weltbild; strittige quantentheoretische Phänomene (Jaynes 2003) und der (meiner Meinung nach) freie menschliche Wille sind davon ausgeklammert.