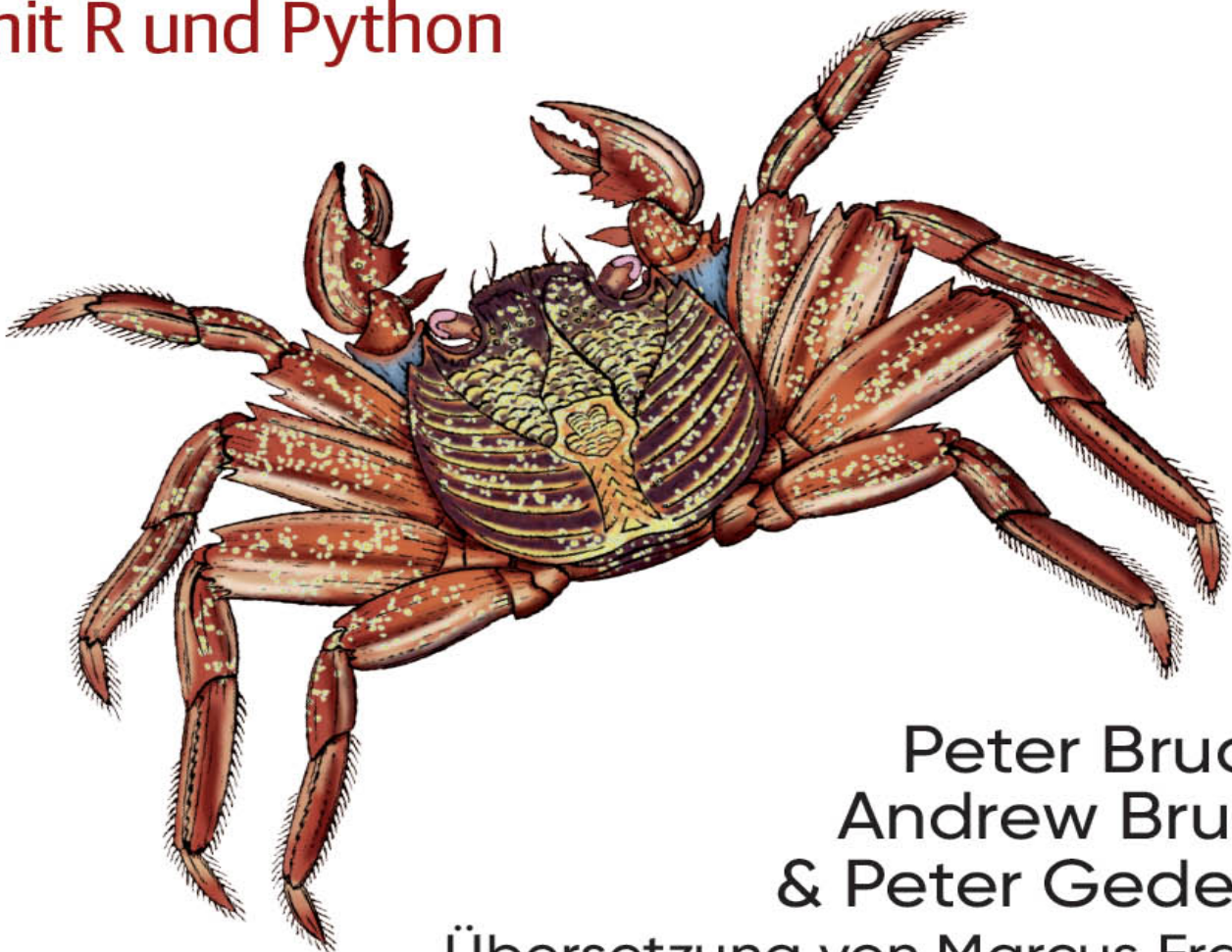


O'REILLY®

Übersetzung der
2. Auflage

Praktische Statistik für Data Scientists

50+ essenzielle Konzepte
mit R und Python



Peter Bruce,
Andrew Bruce
& Peter Gedeck

Übersetzung von Marcus Fraaß

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

Praktische Statistik für Data Scientists

50+ essenzielle Konzepte mit R und Python

Peter Bruce, Andrew Bruce & Peter Gedeck

Deutsche Übersetzung von Marcus Fraaß

O'REILLY®

Peter Bruce, Andrew Bruce, Peter Gedeck

Lektorat: Alexandra Follenius

Übersetzung: Marcus Fraaß

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

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

Herstellung: Stefanie Weidner

Umschlaggestaltung: Karen Montgomery, Michael Oréal, www.oreal.de

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 978-3-96009-153-0

PDF 978-3-96010-467-4

ePub 978-3-96010-468-1

mobi 978-3-96010-469-8

1. Auflage

Translation Copyright für die deutschsprachige Ausgabe © 2021 dpunkt.verlag GmbH

Wieblinger Weg 17
69123 Heidelberg

Authorized German translation of the English edition of *Practical Statistics for Data Scientists, 2nd Edition*, ISBN 9781492072942 © 2020 Peter Bruce, Andrew Bruce, and Peter Gedeck. This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

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: kommentar@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

Vorwort

1 Explorative Datenanalyse

Strukturierte Datentypen

Weiterführende Literatur

Tabellarische Daten

Data Frames und Tabellen

Nicht tabellarische Datenstrukturen

Weiterführende Literatur

Lagemaße

Mittelwert

Median und andere robuste Lagemaße

Beispiel: Lagemaße für Einwohnerzahlen und

Mordraten

Weiterführende Literatur

Streuungsmaße

Standardabweichung und ähnliche Maße

Streuungsmaße auf Basis von Perzentilen

Beispiel: Streuungsmaße für die Einwohnerzahlen

der Bundesstaaten in den USA

Weiterführende Literatur

Exploration der Datenverteilung

Perzentile und Box-Plots

Häufigkeitstabellen und Histogramme

Dichtediagramme und -schätzer

Weiterführende Literatur

Binäre und kategoriale Daten untersuchen

Modus

Erwartungswert

Wahrscheinlichkeiten

Weiterführende Literatur

Korrelation

Streudiagramme

Weiterführende Literatur

Zwei oder mehr Variablen untersuchen

Hexagonal-Binning- und Konturdiagramme
(Diagramme für mehrere numerische Variablen)

Zwei kategoriale Variablen

Kategoriale und numerische Variablen

Mehrere Variablen visualisieren

Weiterführende Literatur

Zusammenfassung

2 Daten- und Stichprobenverteilungen

Zufallsstichprobenziehung und Stichprobenverzerrung

Verzerrung

Zufallsauswahl

Größe versus Qualität: Wann spielt die
Stichprobengröße eine Rolle?

Unterschied zwischen dem Stichproben- und dem
Populationsmittelwert

Weiterführende Literatur

Auswahlverzerrung

Regression zur Mitte

Weiterführende Literatur

Stichprobenverteilung einer statistischen Größe

Zentraler Grenzwertsatz

Standardfehler

Weiterführende Literatur

Bootstrap-Verfahren

Unterschiede zwischen Resampling und dem
Bootstrap-Verfahren

Weiterführende Literatur
Konfidenzintervalle
 Weiterführende Literatur
Normalverteilung
 Standardnormalverteilung und Q-Q-Diagramme
Verteilungen mit langen Verteilungsenden
 Weiterführende Literatur
Studentsche t-Verteilung
 Weiterführende Literatur
Binomialverteilung
 Weiterführende Literatur
Chi-Quadrat-Verteilung
 Weiterführende Literatur
F-Verteilung
 Weiterführende Literatur
Poisson- und verwandte Verteilungen
 Poisson-Verteilung
 Exponentialverteilung
 Die Hazardrate schätzen
 Weibull-Verteilung
 Weiterführende Literatur
Zusammenfassung

3 Statistische Versuche und Signifikanztests

A/B-Test

 Warum eine Kontrollgruppe nutzen?

 Warum lediglich A/B? Warum nicht auch C, D usw.?

 Weiterführende Literatur

Hypothesentests

 Die Nullhypothese

 Die Alternativhypothese

 Einseitige und zweiseitige Hypothesentests

 Weiterführende Literatur

Resampling

 Permutationstest

Beispiel: Die Affinität von Nutzern zu einem Webinhalt messen (Web-Stickiness)
Exakte und Bootstrap-Permutationstests
Permutationstests: ein geeigneter Ausgangspunkt in der Data Science
Weiterführende Literatur
Statistische Signifikanz und p-Werte
p-Wert
Signifikanzniveau
Fehler 1. und 2. Art
Data Science und p-Werte
Weiterführende Literatur
t-Tests
Weiterführende Literatur
Testen mehrerer Hypothesen
Weiterführende Literatur
Die Anzahl der Freiheitsgrade
Weiterführende Literatur
Varianzanalyse (ANOVA)
F-Statistik
Zweifaktorielle Varianzanalyse
Weiterführende Literatur
Chi-Quadrat-Test
Chi-Quadrat-Test: ein Resampling-Ansatz
Chi-Quadrat-Test: die statistische Theorie
Exakter Test nach Fisher
Relevanz in der Data Science
Weiterführende Literatur
Mehrarmige Banditen
Weiterführende Literatur
Trennschärfe und Stichprobengröße
Stichprobengröße
Weiterführende Literatur
Zusammenfassung

4 Regression und Vorhersage

Lineare Einfachregression

- Die Regressionsgleichung

- Angepasste Werte und Residuen

- Die Methode der kleinsten Quadrate

- Unterschied zwischen Vorhersage- und erklärenden Modellen

- Weiterführende Literatur

Multiple lineare Regression

- Beispiel: Die King-County-Immobilien­daten

- Das Modell bewerten

- Kreuzvalidierung

- Modellauswahl und schrittweise Regression

- Gewichtete Regression

- Weiterführende Literatur

Vorhersage mittels Regression

- Risiken bei der Extrapolation

- Konfidenz- und Prognoseintervalle

Regression mit Faktorvariablen

- Darstellung durch Dummy-Variablen

- Faktorvariablen mit vielen Stufen

- Geordnete Faktorvariablen

Interpretieren der Regressionsgleichung

- Korrelierte Prädiktorvariablen

- Multikollinearität

- Konfundierende Variablen

- Interaktions- und Haupteffekte

Regressionsdiagnostik

- Ausreißer

- Einflussreiche Beobachtungen

- Heteroskedastische, nicht normalverteilte und korrelierte Fehler

- Partielle Residuendiagramme und Nichtlinearität

Polynomiale und Spline-Regression

- Polynome

- Splines

- Verallgemeinerte additive Modelle

Weiterführende Literatur
Zusammenfassung

5 Klassifikation

Naiver Bayes-Klassifikator

Warum eine exakte bayessche Klassifikation nicht praktikabel ist

Die naive Lösung

Numerische Prädiktorvariablen

Weiterführende Literatur

Diskriminanzanalyse

Kovarianzmatrix

Lineare Diskriminanzanalyse nach Fisher

Ein einfaches Beispiel

Weiterführende Literatur

Logistische Regression

Logistische Antwortfunktion und Logit-Funktion

Logistische Regression und verallgemeinerte lineare Modelle

Verallgemeinerte lineare Modelle

Vorhergesagte Werte aus der logistischen Regression

Interpretation der Koeffizienten und Odds-Ratios

Lineare und logistische Regression:

Gemeinsamkeiten und Unterschiede

Das Modell prüfen und bewerten

Weiterführende Literatur

Klassifikationsmodelle bewerten

Konfusionsmatrix

Die Problematik seltener Kategorien

Relevanz, Sensitivität und Spezifität

ROC-Kurve

Fläche unter der ROC-Kurve (AUC)

Lift

Weiterführende Literatur

Strategien bei unausgewogenen Daten

- Undersampling
- Oversampling und Up/Down Weighting
- Generierung von Daten
- Kostenbasierte Klassifikation
- Die Vorhersagen untersuchen
- Weiterführende Literatur
- Zusammenfassung

6 Statistisches maschinelles Lernen

- K-Nächste-Nachbarn

- Ein kleines Beispiel: Vorhersage von Kreditausfällen
 - Distanzmaße
 - 1-aus-n-Codierung
 - Standardisierung (Normierung, z-Werte)
 - K festlegen
 - KNN zur Merkmalskonstruktion

- Baummodelle

- Ein einfaches Beispiel
 - Der Recursive-Partitioning-Algorithmus
 - Homogenität und Unreinheit messen
 - Den Baum daran hindern, weiterzuwachsen
 - Vorhersage eines kontinuierlichen Werts
 - Wie Bäume verwendet werden
 - Weiterführende Literatur

- Bagging und Random Forests

- Bagging
 - Random Forest
 - Variablenwichtigkeit
 - Hyperparameter

- Boosting

- Der Boosting-Algorithmus
 - XGBoost
 - Regularisierung: Überanpassung vermeiden
 - Hyperparameter und Kreuzvalidierung

- Zusammenfassung

7 Unüberwachtes Lernen

Hauptkomponentenanalyse

Ein einfaches Beispiel

Die Hauptkomponenten berechnen

Die Hauptkomponenten interpretieren

Korrespondenzanalyse

Weiterführende Literatur

K-Means-Clustering

Ein einfaches Beispiel

Der K-Means-Algorithmus

Die Cluster interpretieren

Die Anzahl von Clustern bestimmen

Hierarchische Clusteranalyse

Ein einfaches Beispiel

Das Dendrogramm

Der agglomerative Algorithmus

Ähnlichkeitsmaße

Modellbasierte Clusteranalyse

Multivariate Normalverteilung

Zusammengesetzte Normalverteilungen (gaußsche Mischverteilungen)

Die Anzahl der Cluster bestimmen

Weiterführende Literatur

Skalierung und kategoriale Variablen

Variablen skalieren

Dominierende Variablen

Kategoriale Daten und die Gower-Distanz

Probleme bei der Clusteranalyse mit verschiedenen Datentypen

Zusammenfassung

Quellenangaben

Index

Vorwort

Dieses Buch richtet sich an Data Scientists, die mit den Programmiersprachen *R* und/oder *Python* vertraut sind und sich bereits früher (wenn auch nur punktuell oder zeitweise) mit Statistik beschäftigt haben. Zwei der Autoren entstammen der Welt der Statistik, ehe sie sich in den weiten Raum der Data Science begeben haben, und schätzen den Beitrag, den die Statistik zur Datenwissenschaft zu leisten vermag, sehr. Gleichzeitig sind wir uns der Grenzen des traditionellen Statistikerunterrichts durchaus bewusst: Statistik als Disziplin ist anderthalb Jahrhunderte alt, und die meisten Statistiklehrbücher und -kurse sind nicht gerade von Dynamik geprägt, sondern erinnern eher an die Trägheit eines Ozeanriesen. Alle Methoden in diesem Buch haben einen gewissen historischen oder methodologischen Bezug zur Disziplin der Statistik. Methoden, die sich hauptsächlich aus der Informatik entwickelt haben, wie z.B. neuronale Netze, werden nicht behandelt.

Diesem Buch liegen zwei Ziele zugrunde:

- Schlüsselbegriffe aus der Statistik, die für die Data Science relevant sind, in zugänglicher, übersichtlich gegliederter und leicht referenzierbarer Form darzulegen.

- Eine Erläuterung dazu zu geben, welche Konzepte aus datenwissenschaftlicher Sicht wichtig und nützlich sind, welche weniger wichtig sind und warum.

In diesem Buch verwendete Konventionen

Die folgenden typografischen Konventionen werden in diesem Buch verwendet:

Kursiv

Kennzeichnet neue Begriffe, URLs, E-Mail-Adressen, Dateinamen und Dateiendungen.

Konstante Zeichenbreite

Wird für Programmlistings und für Programmelemente in Textabschnitten wie Namen von Variablen und Funktionen, Datenbanken, Datentypen, Umgebungsvariablen, Anweisungen und Schlüsselwörter verwendet.

Konstante Zeichenbreite, fett

Kennzeichnet Befehle oder anderen Text, den der Nutzer wörtlich eingeben sollte.

Schlüsselbegriffe

Die Data Science baut auf mehreren Disziplinen auf, darunter Statistik, Informatik, Informationstechnologie und domänenspezifische Bereiche. Infolgedessen können mehrere unterschiedliche Begriffe verwendet werden, um auf ein bestimmtes Konzept zu verweisen. Schlüsselbegriffe und ihre Synonyme werden im gesamten Buch in einem Kasten wie diesem hervorgehoben.

Dieses Symbol steht für einen Tipp oder eine Empfehlung.



Dieses Symbol steht für einen allgemeinen Hinweis.



Dieses Symbol warnt oder mahnt zur Vorsicht.



Verwenden von Codebeispielen

Zu sämtlichen Beispielen zeigen wir in diesem Buch die entsprechenden Codebeispiele - zuerst immer in *R* und dann in *Python*. Um unnötige Wiederholungen zu vermeiden, zeigen wir im Allgemeinen nur Ausgaben und Diagramme, die durch den *R*-Code erzeugt wurden. Wir klammern auch den Code aus, der zum Laden der erforderlichen Pakete und Datensätze erforderlich ist. Den vollständigen Code sowie die Datensätze zum Herunterladen finden Sie unter <https://github.com/gedeck/practical-statistics-for-data-scientists>.

Dieses Buch dient dazu, Ihnen beim Erledigen Ihrer Arbeit zu helfen. Im Allgemeinen dürfen Sie die Codebeispiele aus diesem Buch in Ihren eigenen Programmen und der dazugehörigen Dokumentation verwenden. Sie müssen uns dazu nicht um Erlaubnis bitten, solange Sie nicht einen beträchtlichen Teil des Codes reproduzieren. Beispielsweise benötigen Sie keine Erlaubnis, um ein Programm zu schreiben, in dem mehrere Codefragmente aus diesem Buch vorkommen. Wollen Sie dagegen eine CD-ROM mit Beispielen aus Büchern von O'Reilly verkaufen oder verteilen, brauchen Sie eine Erlaubnis. Eine Frage zu

beantworten, indem Sie aus diesem Buch zitieren und ein Codebeispiel wiedergeben, benötigt keine Erlaubnis. Eine beträchtliche Menge Beispielcode aus diesem Buch in die Dokumentation Ihres Produkts aufzunehmen, bedarf hingegen unserer ausdrücklichen Zustimmung.

Wir freuen uns über Zitate, verlangen diese aber nicht. Ein Zitat enthält Titel, Autor, Verlag und ISBN, zum Beispiel: »*Praktische Statistik für Data Scientists* von Peter Bruce, Andrew Bruce und Peter Gedeck (O'Reilly). Copyright 2020 Peter Bruce, Andrew Bruce und Peter Gedeck, ISBN 978-3-96009-153-0.«

Wenn Sie glauben, dass Ihre Verwendung von Codebeispielen über die übliche Nutzung hinausgeht oder außerhalb der oben vorgestellten Nutzungsbedingungen liegt, kontaktieren Sie uns bitte unter kontakt@oreilly.de.

Danksagungen

Die Autoren danken den zahlreichen Menschen, die dazu beigetragen haben, dieses Buch Wirklichkeit werden zu lassen.

Gerhard Pilcher, CEO des Data-Mining-Unternehmens Elder Research, sah frühe Entwürfe des Buchs und half uns mit detaillierten und hilfreichen Korrekturen sowie Kommentaren. Ebenso gaben Anya McGuirk und Wei Xiao, Statistiker bei SAS, und Jay Hilfiger, ebenfalls Autor von O'Reilly, hilfreiches Feedback zu den ersten Entwürfen des Buchs. Toshiaki Kurokawa, der die erste Auflage ins Japanische übersetzte, leistete dabei umfassende Überarbeitungs- und Korrekturarbeit. Aaron Schumacher und Walter Paczkowski haben die zweite Auflage des Buchs gründlich überarbeitet und zahlreiche hilfreiche und wertvolle Anregungen gegeben, für die wir sehr dankbar

sind. Es versteht sich von selbst, dass alle noch verbleibenden Fehler allein auf uns zurückzuführen sind.

Bei O'Reilly begleitete uns Shannon Cutt mit guter Laune und der richtigen Portion Nachdruck durch den Publikationsprozess, während Kristen Brown unser Buch reibungslos durch den Produktionsprozess geführt hat. Rachel Monaghan und Eliahu Sussman korrigierten und verbesserten unser Buch mit Sorgfalt und Geduld, während Ellen Troutman-Zaig den Index erarbeitete. Nicole Tache übernahm das Lektorat der zweiten Auflage und hat den Prozess effektiv geleitet sowie viele gute redaktionelle Vorschläge gemacht, um die Lesbarkeit des Buchs für ein breites Publikum zu verbessern. Wir danken auch Marie Beaugureau, die unser Projekt bei O'Reilly initiiert hat, sowie Ben Bengfort, Autor von O'Reilly und Ausbilder bei [Statistics.com](https://www.statistics.com), der uns O'Reilly vorgestellt hat.

Wir und dieses Buch haben auch von den vielen Gesprächen profitiert, die Peter im Laufe der Jahre mit Galit Shmueli, Mitautorin bei anderen Buchprojekten, geführt hat.

Schließlich möchten wir besonders Elizabeth Bruce und Deborah Donnell danken, deren Geduld und Unterstützung dieses Vorhaben möglich gemacht haben.

Explorative Datenanalyse

Dieses Kapitel erläutert Ihnen den ersten Schritt in jedem datenwissenschaftlichen Projekt: die Datenexploration.

Die klassische Statistik konzentrierte sich fast ausschließlich auf die *Inferenz*, einen manchmal komplexen Satz von Verfahren, um aus kleinen Stichproben Rückschlüsse auf eine größere Grundgesamtheit zu ziehen. Im Jahr 1962 forderte John W. Tukey (<https://oreil.ly/LQw6q>) (siehe [Abbildung 1-1](#)) in seinem bahnbrechenden Aufsatz »The Future of Data Analysis« [Tukey-1962] eine Reform der Statistik. Er schlug eine neue wissenschaftliche Disziplin namens *Datenanalyse* vor, die die statistische Inferenz lediglich als eine Komponente enthielt. Tukey knüpfte Kontakte zu den Ingenieurs- und Informatikgemeinschaften (er prägte die Begriffe *Bit*, kurz für Binärziffer, und *Software*). Seine damaligen Ansätze haben bis heute überraschend Bestand und bilden einen Teil der Grundlagen der Data Science. Der Fachbereich der explorativen Datenanalyse wurde mit Tukeys im Jahr 1977 erschienenem und inzwischen als Klassiker geltendem Buch *Exploratory Data Analysis* [Tukey-1977] begründet. Tukey stellte darin einfache Diagramme (z.B. Box-Plots und Streudiagramme) vor, die in Kombination mit zusammenfassenden Statistiken (Mittelwert, Median,

Quantile usw.) dabei helfen, ein Bild eines Datensatzes zu zeichnen.

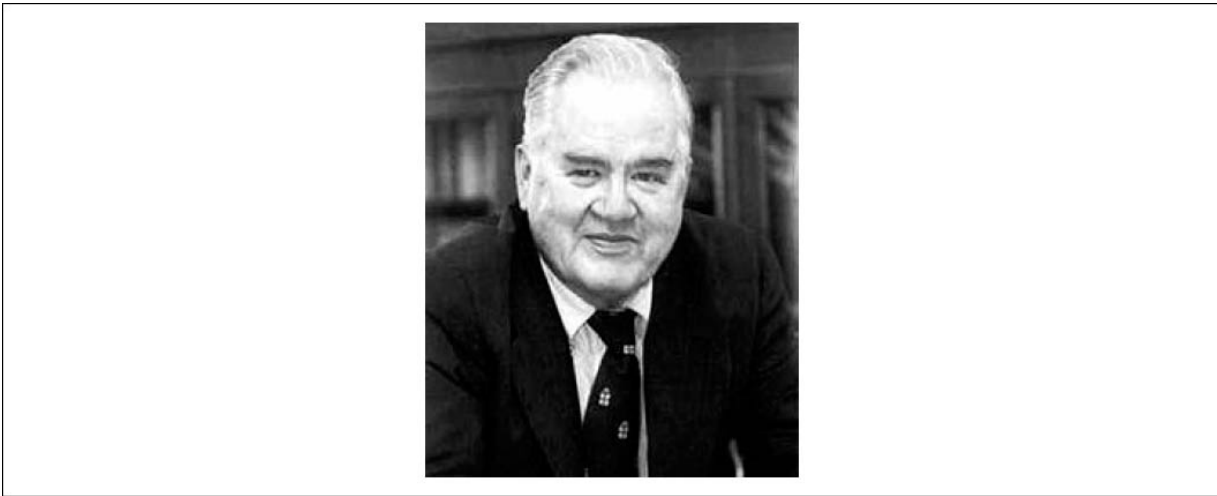


Abbildung 1-1: John Tukey, der bedeutende Statistiker, dessen vor über 50 Jahren entwickelte Ideen die Grundlage der Data Science bilden

Mit der zunehmenden Verfügbarkeit von Rechenleistung und leistungsfähigen Datenanalyseprogrammen hat sich die explorative Datenanalyse weit über ihren ursprünglichen Rahmen hinaus weiterentwickelt. Die wichtigsten Triebkräfte dieser Disziplin waren die rasche Entwicklung neuer Technologien, der Zugang zu mehr und umfangreicheren Daten und der verstärkte Einsatz der quantitativen Analyse in einer Vielzahl von Disziplinen. David Donoho, Professor für Statistik an der Stanford University und ehemaliger Student Tukeys, verfasste einen ausgezeichneten Artikel auf der Grundlage seiner Präsentation auf dem Workshop zur Hundertjahrfeier von Tukey in Princeton, New Jersey [Donoho-2015]. Donoho führt die Entwicklung der Data Science auf Tukeys Pionierarbeit in der Datenanalyse zurück.

Strukturierte Datentypen

Es gibt zahlreiche unterschiedliche Datenquellen: Sensormessungen, Ereignisse, Text, Bilder und Videos. Das *Internet der Dinge* (engl. *Internet of Things* (IoT)) produziert ständig neue Informationsfluten. Ein Großteil dieser Daten liegt unstrukturiert vor: Bilder sind nichts anderes als eine Zusammenstellung von Pixeln, wobei jedes Pixel RGB-Farbinformationen (Rot, Grün, Blau) enthält. Texte sind Folgen von Wörtern und Nicht-Wortzeichen, die oft in Abschnitte, Unterabschnitte usw. gegliedert sind. Clickstreams sind Handlungsverläufe eines Nutzers, der mit einer Anwendung oder einer Webseite interagiert. Tatsächlich besteht eine große Herausforderung der Datenwissenschaft darin, diese Flut von Rohdaten in verwertbare Informationen zu überführen. Um die in diesem Buch behandelten statistischen Konzepte in Anwendung zu bringen, müssen unstrukturierte Rohdaten zunächst aufbereitet und in eine strukturierte Form überführt werden. Eine der am häufigsten vorkommenden Formen strukturierter Daten ist eine Tabelle mit Zeilen und Spalten – so wie Daten aus einer relationalen Datenbank oder Daten, die für eine Studie erhoben wurden.

Es gibt zwei grundlegende Arten strukturierter Daten: numerische und kategoriale Daten. Numerische Daten treten in zwei Formen auf: *kontinuierlich*, wie z.B. die Windgeschwindigkeit oder die zeitliche Dauer, und *diskret*, wie z.B. die Häufigkeit des Auftretens eines Ereignisses. *Kategoriale* Daten nehmen nur einen bestimmten Satz von Werten an, wie z.B. einen TV-Bildschirmtyp (Plasma, LCD, LED usw.) oder den Namen eines Bundesstaats (Alabama, Alaska usw.). *Binäre* Daten sind ein wichtiger Spezialfall kategorialer Daten, die nur einen von zwei möglichen Werten annehmen, wie z.B. 0 oder 1, ja oder nein oder auch wahr oder falsch. Ein weiterer nützlicher kategorialer Datentyp sind *ordinalskalierte* Daten, bei denen die

Kategorien in einer Reihenfolge geordnet sind; ein Beispiel hierfür ist eine numerische Bewertung (1, 2, 3, 4 oder 5).

Warum plagen wir uns mit der Taxonomie der Datentypen herum? Es stellt sich heraus, dass für die Zwecke der Datenanalyse und der prädiktiven Modellierung der Datentyp wichtig ist, um die Art der visuellen Darstellung, der Datenanalyse oder des statistischen Modells zu bestimmen. Tatsächlich verwenden datenwissenschaftliche Softwareprogramme wie *R* und *Python* diese Datentypen, um die Rechenleistung zu optimieren. Noch wichtiger ist es, dass der Datentyp einer Variablen ausschlaggebend dafür ist, wie das Programm die Berechnungen für diese Variable handhabt.

Schlüsselbegriffe zu Datentypen

Numerisch

Daten, die auf einer numerischen Skala abgebildet sind.

Kontinuierlich

Daten, die innerhalb eines Intervalls einen beliebigen Wert annehmen können.

Synonyme

intervallskaliert, Gleitkommazahl, numerisch

Diskret

Daten, die nur ganzzahlige Werte annehmen können, wie z. B. Häufigkeiten bzw. Zählungen.

Synonyme

Ganzzahl, Zählwert

Kategorial

Daten, die nur einen bestimmten Satz von Werten annehmen können, die wiederum einen Satz von möglichen Kategorien repräsentieren.

Synonyme

Aufzählungstyp, Faktor, faktoriell, nominal

Binär

Ein Spezialfall des kategorialen Datentyps mit nur zwei möglichen Ausprägungen, z.B. 0/1, wahr/falsch.

Synonyme

dichotom, logisch, Indikatorvariable, boolesche Variable

Ordinalskaliert

Kategoriale Daten, die eine eindeutige Reihenfolge bzw. Rangordnung haben.

Synonym

geordneter Faktor

Softwareingenieure und Datenbankprogrammierer fragen sich vielleicht, warum wir überhaupt den Begriff der *kategorialen* und *ordinalskalierten* Daten für unsere Analyse benötigen. Schließlich sind Kategorien lediglich eine Sammlung von Text- (oder numerischen) Werten, und die zugrunde liegende Datenbank übernimmt automatisch die interne Darstellung. Die explizite Bestimmung von Daten als kategoriale Daten im Vergleich zu Textdaten bietet jedoch einige Vorteile:

- Die Kenntnis, dass Daten kategorial sind, kann als Signal dienen, durch das ein Softwareprogramm erkennen kann, wie sich statistische Verfahren wie die Erstellung eines Diagramms oder die Anpassung eines Modells verhalten sollen. Insbesondere ordinalskalierte Daten können als `ordered.factor` in *R* angegeben werden, wodurch eine benutzerdefinierte Ordnung in Diagrammen, Tabellen und Modellen erhalten bleibt. In *Python* unterstützt `scikit-learn` ordinalskalierte Daten mit der Methode `sklearn.preprocessing.OrdinalEncoder`.
- Das Speichern und Indizieren kann optimiert werden (wie in einer relationalen Datenbank).
- Die möglichen Werte, die eine gegebene kategoriale Variable annehmen kann, werden in dem Softwareprogramm erzwungen (wie bei einer Aufzählung).

Der dritte »Vorteil« kann zu unbeabsichtigtem bzw. unerwartetem Verhalten führen: Das Standardverhalten

von Datenimportfunktionen in *R* (z.B. `read.csv`) besteht darin, eine Textspalte automatisch in einen `factor` umzuwandeln. Bei nachfolgenden Operationen auf dieser Spalte wird davon ausgegangen, dass die einzigen zulässigen Werte für diese Spalte die ursprünglich importierten sind und die Zuweisung eines neuen Textwerts eine Warnung verursacht sowie einen Eintrag mit dem Wert `NA` (ein fehlender Wert) erzeugt. Das `pandas`-Paket in *Python* nimmt diese Umwandlung nicht automatisch vor. Sie können jedoch in der Funktion `read_csv` eine Spalte explizit als kategorial spezifizieren.

Kernideen

- Daten werden in Softwareprogrammen typischerweise in verschiedene Typen eingeteilt.
- Zu den Datentypen gehören numerische (kontinuierlich, diskret) und kategoriale (binär, ordinalskaliert).
- Die Datentypisierung dient als Signal für das Softwareprogramm, wie die Daten zu verarbeiten sind.

Weiterführende Literatur

- Datentypen können verwirrend sein, da sich Typen überschneiden und die Taxonomie in einem Softwareprogramm von der in einem anderen abweichen kann. Auf der R-Tutorial-Webseite (<https://oreil.ly/2YUoA>) können Sie die Taxonomie in *R* nachvollziehen. Die `pandas`-Dokumentation (<https://oreil.ly/UGX-4>) beschreibt die verschiedenen Datentypen in *Python* und wie sie verändert werden können.
- Datenbanken sind in ihrer Einteilung der Datentypen detaillierter und berücksichtigen Präzisionsniveaus, Datenfelder fester oder variabler Länge und mehr

(siehe den W3Schools-SQL-Leitfaden (<https://oreil.ly/cThTM>).)

Tabellarische Daten

Der typische Bezugsrahmen für eine Analyse in der Data Science ist ein *tabellarisches Datenobjekt* (engl. *Rectangular Data Object*), wie eine Tabellenkalkulation oder eine Datenbanktabelle.

»Tabellarische Daten« ist der allgemeine Begriff für eine zweidimensionale Matrix mit Zeilen für die Beobachtungen (Fälle) und Spalten für die Merkmale (Variablen); in *R* und *Python* wird dies als *Data Frame* bezeichnet. Die Daten sind zu Beginn nicht immer in dieser Form vorhanden: Unstrukturierte Daten (z.B. Text) müssen zunächst so verarbeitet und aufbereitet werden, dass sie als eine Reihe von Merkmalen in tabellarischer Struktur dargestellt werden können (siehe »[Strukturierte Datentypen](#)« auf [Seite 2](#)). Daten in relationalen Datenbanken müssen für die meisten Datenanalyse- und Modellierungsaufgaben extrahiert und in eine einzelne Tabelle überführt werden.

Schlüsselbegriffe zu tabellarischen Daten

Data Frame

Tabellarische Daten (wie ein Tabellenkalkulationsblatt) sind die grundlegende Datenstruktur für statistische und maschinelle Lernmodelle.

Merkmal

Eine Spalte innerhalb einer Tabelle wird allgemein als *Merkmal* (engl. *Feature*) bezeichnet.

Synonyme

Attribut, Eingabe, Prädiktorvariable, Prädiktor, unabhängige Variable

Ergebnis

Viele datenwissenschaftliche Projekte zielen auf die Vorhersage eines *Ergebnisses* (engl. *Outcome*) ab – oft in Form eines Ja-oder-Nein-Ergebnisses (ob beispielsweise in [Tabelle 1-1](#) eine »Auktion umkämpft war oder nicht«). Die *Merkmale* werden manchmal verwendet, um das *Ergebnis* eines statistischen Versuchs oder einer Studie vorherzusagen..

Synonyme

Ergebnisvariable, abhängige Variable, Antwortvariable, Zielgröße, Ausgabe, Responsevariable

Eintrag

Eine Zeile innerhalb einer Tabelle wird allgemein als *Eintrag* (engl. *Record*) bezeichnet.

Synonyme

Fall, Beispiel, Instanz, Beobachtung

Tabelle 1-1: Ein typisches Data-Frame-Format

Kategorie	Währung	Verkäufer-Rating	Dauer	Schluss-tag	Schluss-preis	Eröffnungs-preis	umkämpft?
Musik/Film/Spiel	USD	3249	5	Mon	0.01	0.01	0
Musik/Film/Spiel	USD	3249	5	Mon	0.01	0.01	0
Automobil	USD	3115	7	Die	0.01	0.01	0
Automobil	USD	3115	7	Die	0.01	0.01	0
Automobil	USD	3115	7	Die	0.01	0.01	0
Automobil	USD	3115	7	Die	0.01	0.01	0
Automobil	USD	3115	7	Die	0.01	0.01	1
Automobil	USD	3115	7	Die	0.01	0.01	1

In [Tabelle 1-1](#) gibt es eine Kombination aus Mess- oder Zähl-daten (z.B. Dauer und Preis) und kategorialen Daten (z.B. Kategorie und Währung). Wie bereits erwähnt, ist eine besondere Form der kategorialen Variablen eine binäre Variable (ja/nein oder 0/1), wie in der Spalte ganz rechts in [Tabelle 1-1](#) – eine Indikatorvariable, die angibt, ob eine Auktion umkämpft war (mehrere Bieter hatte) oder nicht. Diese Indikatorvariable ist zufällig auch eine

Ergebnisvariable, wenn das Modell vorhersagen soll, ob eine Auktion umkämpft sein wird oder nicht.

Data Frames und Tabellen

Klassische Datenbanktabellen haben eine oder mehrere Spalten, die als Index bezeichnet werden und im Wesentlichen eine Zeilennummer darstellen. Dies kann die Effizienz bestimmter Datenbankabfragen erheblich verbessern. In *Pythons* pandas-Bibliothek wird die grundlegende tabellarische Datenstruktur durch ein Data-Frame-Objekt umgesetzt. Standardmäßig wird automatisch ein ganzzahliger Index für ein Data-Frame-Objekt basierend auf der Reihenfolge der Zeilen erstellt. In pandas ist es auch möglich, mehrstufige bzw. hierarchische Indizes festzulegen, um die Effizienz bestimmter Operationen zu verbessern.

In *R* ist die grundlegende tabellarische Datenstruktur mittels eines `data.frame`-Objekts implementiert. Ein `data.frame` hat auch einen impliziten ganzzahligen Index, der auf der Zeilenreihenfolge basiert. Der standardmäßige `data.frame` in *R* unterstützt keine benutzerdefinierten oder mehrstufigen Indizes. Jedoch kann über das Argument `row.names` ein benutzerdefinierter Schlüssel erstellt werden. Um diesem Problem zu begegnen, werden immer häufiger zwei neuere Pakete eingesetzt: `data.table` und `dplyr`. Beide unterstützen mehrstufige Indizes und bieten erhebliche Beschleunigungen bei der Arbeit mit einem `data.frame`.

Unterschiede in der Terminologie

Die Terminologie bei tabellarischen Daten kann verwirrend sein. Statistiker und Data Scientists verwenden oftmals unterschiedliche Begriffe für ein und denselben Sachverhalt. Statistiker nutzen in einem Modell *Prädiktorvariablen*, um eine *Antwortvariable* (engl. *Response*) oder eine *abhängige Variable* vorherzusagen. Ein Datenwissenschaftler spricht von *Merkmalen* (engl. *Features*), um eine *Zielgröße* (engl. *Target*) vorherzusagen.

Ein Synonym ist besonders verwirrend: Informatiker verwenden den Begriff *Stichprobe* (engl. *Sample*) für eine einzelne Datenzeile, für einen Statistiker ist eine *Stichprobe* hingegen eine Sammlung von Datenzeilen.

Nicht tabellarische Datenstrukturen

Neben tabellarischen Daten gibt es noch andere Datenstrukturen.

Zeitreihendaten umfassen aufeinanderfolgende Messungen derselben Variablen. Sie sind das Rohmaterial für statistische Prognosemethoden und auch eine zentrale Komponente der von Geräten – dem Internet der Dinge – erzeugten Daten.

Räumliche Daten- bzw. Geodatenstrukturen, die bei der Kartierung und Standortanalyse verwendet werden, sind komplexer und vielfältiger als tabellarische Datenstrukturen. In der *Objektdarstellung* (engl. *Object Representation*) stehen ein Objekt (z.B. ein Haus) und seine räumlichen Koordinaten im Mittelpunkt der Daten. Die *Feldansicht* (engl. *Field View*) hingegen konzentriert sich auf kleine räumliche Einheiten und den Wert einer relevanten Metrik (z.B. Pixelhelligkeit).

Graphen- (oder Netzwerk-) Datenstrukturen werden verwendet, um physikalische, soziale oder abstrakte Beziehungen darzustellen. Beispielsweise kann ein Diagramm eines sozialen Netzwerks wie Facebook oder LinkedIn Verbindungen zwischen Menschen im Netzwerk darstellen. Ein Beispiel für ein physisches Netzwerk sind Vertriebszentren, die durch Straßen miteinander verbunden sind. Diagrammstrukturen sind für bestimmte Arten von Fragestellungen nützlich, wie z.B. bei der Netzwerkoptimierung und bei Empfehlungssystemen.

Jeder dieser Datentypen hat seine eigene spezifische Methodologie in der Data Science. Der Schwerpunkt dieses

Buchs liegt auf tabellarische Daten, dem grundlegenden Baustein der prädiktiven Modellierung.

Graphen in der Statistik

In der Informatik und der Informationstechnologie bezieht sich der Begriff *Graph* typischerweise auf die Darstellung von Verbindungen zwischen Entitäten und auf die zugrunde liegende Datenstruktur. In der Statistik wird der Begriff *Graph* verwendet, um sich auf eine Vielzahl von Darstellungen und Visualisierungen zu beziehen, nicht nur von Verbindungen zwischen Entitäten. Zudem bezieht er sich ausschließlich auf die Visualisierung und nicht auf die Datenstruktur.

Kernideen

- Die grundlegende Datenstruktur in der Data Science ist eine rechteckige Matrix, in der die Zeilen den Beobachtungen entsprechen und die Spalten den Variablen (Merkmalen).
- Die Terminologie kann verwirrend sein; es gibt eine Vielzahl von Synonymen, die sich aus den verschiedenen Disziplinen ergeben, die zur Data Science beitragen (Statistik, Informatik und Informationstechnologie).

Weiterführende Literatur

- Dokumentation zu Data Frames in *R* (<https://oreil.ly/NsONR>)
- Dokumentation zu Data Frames in *Python* (<https://oreil.ly/oxDKQ>)

Lagemaße

Variablen für Mess- oder Zähldaten können Tausende von unterschiedlichen Werten haben. Ein grundlegender Schritt bei der Erkundung Ihrer Daten ist die Ermittlung eines »typischen Werts« für jedes Merkmal (Variable) – ein sogenanntes Lagemaß (engl. *Estimates of Location*): eine

Schätzung darüber, wo sich die Mehrheit der Daten konzentriert (d.h. ihre zentrale Tendenz).

Schlüsselbegriffe zu Lagemaßen

Mittelwert

Die Summe aller Werte dividiert durch die Anzahl der Werte.

Synonyme

arithmetisches Mittel, Durchschnitt

Gewichteter Mittelwert

Die Summe aller Werte, die jeweils mit einem Gewicht bzw. einem Gewichtungsfaktor multipliziert werden, geteilt durch die Summe aller Gewichte.

Synonym

gewichteter Durchschnitt

Median

Der Wert, bei dem die Hälfte der Daten oberhalb und die andere Hälfte unterhalb dieses Werts liegt.

Synonym

50%-Perzentil

Perzentil

Der Wert, bei dem $P\%$ der Daten unterhalb dieses Werts liegen.

Synonym

Quantil

Gewichteter Median

Der Wert, bei dem die Summe der Gewichte der sortierten Daten exakt die Hälfte beträgt und der die Daten so einteilt, dass sie entweder oberhalb oder unterhalb diesen Werts liegen.

Getrimmter Mittelwert

Der Mittelwert aller Werte, nachdem eine vorgegebene Anzahl von Ausreißern entfernt wurde.

Synonym

gestutzter Mittelwert

Robust

Nicht sensibel gegenüber Ausreißern.

Ausreißer

Ein Datenwert, der sich stark von den übrigen Daten unterscheidet.

Synonym

Extremwert