

Ein Buch zum Mitmachen und Verstehen

Datenanalyse

von Kopf bis Fuß



Sagen Sie mit
der linearen
Regression
Ihre Gehalts-
erhöhung
voraus



Finden Sie heraus,
wer Ihre Kunden
wirklich sind



Spielen Sie zentrale
Statistik-Konzepte
direkt in Ihr Hirn



**Hier geht's um große
Zahlen, Statistik und
richtige Entscheidungen**

Verkaufen Sie mehr
Spielzeug, indem Sie
Ihr Geschäftsmodell
optimieren



Befreien Sie
sich von
Denkfehlern



Bereinigen Sie
unsaubere Daten
für eine effiziente
Auswertung

O'REILLY®

Michael Milton
Deutsche Übersetzung von Jörg Beyer

Die Informationen in diesem Buch wurden mit größter Sorgfalt erarbeitet. Dennoch können Fehler nicht vollständig ausgeschlossen werden. Verlag, Autoren und Übersetzer übernehmen keine juristische Verantwortung oder irgendeine Haftung für eventuell verbliebene Fehler und deren Folgen. D.h., wenn Sie beispielsweise ein Kernkraftwerk unter Verwendung dieses Buchs betreiben möchten, tun Sie dies auf eigene Gefahr.

Alle Warennamen werden ohne Gewährleistung der freien Verwendbarkeit benutzt und sind möglicherweise eingetragene Warenzeichen. Der Verlag richtet sich im Wesentlichen nach den Schreibweisen der Hersteller. Das Werk einschließlich aller seiner Teile ist urheberrechtlich geschützt. Alle Rechte vorbehalten einschließlich der Vervielfältigung, Übersetzung, Mikroverfilmung sowie Einspeicherung und Verarbeitung in elektronischen Systemen.

Kommentare und Fragen können Sie gerne an uns richten:

O'Reilly Verlag
Balthasarstr. 81
50670 Köln
Tel.: 0221/9731600
Fax: 0221/9731608
E-Mail: kommentar@oreilly.de

Copyright der deutschen Ausgabe:

© 2010 by O'Reilly Verlag GmbH & Co. KG



Die Originalausgabe erschien 2009 unter dem Titel
Head First Data Analysis bei O'Reilly Media, Inc.

Bibliografische Information Der Deutschen Bibliothek
Die Deutsche Bibliothek verzeichnet diese Publikation in der
Deutschen Nationalbibliografie; detaillierte bibliografische Daten
sind im Internet über <http://dnb.ddb.de> abrufbar.

Deutsche Übersetzung und Bearbeitung: Jörg Beyer, Weimar (Lahn)
Lektorat: Christine Haite, Köln
Korrektur: Sibylle Feldmann, Düsseldorf

Satz: Ulrich Borstelmann, Dortmund
Umschlaggestaltung: Karen Montgomery, Boston
Produktion: Andrea Miß, Köln
Belichtung, Druck und buchbinderische Verarbeitung: Media-Print, Paderborn

ISBN-13 978-3-89721-959-5

Dieses Buch ist auf 100% chlorfrei gebleichtem Papier gedruckt.

Bei der Produktion dieses Buchs kamen keine Daten zu Schaden.

Dem Andenken an meine Großmutter Jane Reese Gibbs gewidmet.

Der Autor von Datenanalyse von Kopf bis Fuß



Michael Milton ↗

Michael Milton hat einen wesentlichen Teil seines bisherigen Berufslebens damit verbracht, gemeinnützige Organisationen darin zu unterstützen, ihre Mittelbeschaffung zu verbessern, indem er Spenderdaten interpretiert und entsprechende strategische Maßnahmen daraus abgeleitet hat.

Er hat einen Abschluss in Philosophie des New College of Florida und einen weiteren in Religion und Ethik der Yale University. Nach jahrelangem Lesen *langweiliger*, mit unglaublich wichtigen Sachen vollgestopfter Bücher war die Lektüre der Bücher aus der *Von Kopf bis Fuß*-Reihe eine Offenbarung für ihn, und deshalb ist er besonders dankbar, die Gelegenheit bekommen zu haben, jetzt selbst *ein aufregendes* und mit unglaublich wichtigen Sachen vollgestopftes Buch zu schreiben.

Wenn er sich nicht gerade in einer Bibliothek oder Buchhandlung herumtreibt, können Sie ihm beim Laufen, Fotografieren oder Bierbrauen zusehen.

Der Übersetzer dieses Buchs

Jörg Beyer arbeitete ein paar Jahre als Schreiner, bevor er den Reboot-Knopf drückte und in Marburg Psychologie studierte.

Zur Erleichterung seines Arbeitslebens hat er sich im Lauf der Zeit allerhand angeeignet, mit dem sich große Datenmengen bewältigen und Routinearbeiten beschleunigen lassen: Datenbankentwicklung und SQL, AppleScript, Perl, XML, R ... eine lange Liste, die immer noch ein bisschen länger zu werden droht.

Seine Beziehung zur Statistik lässt sich am besten als Liebe auf den zweiten Blick beschreiben, und wenn sich heute irgendwo ein Datensatz blicken lässt, löst das einen sofortigen Beutefangreflex bei ihm aus.

Er arbeitet als Statistik-Consultant im Gesundheitswesen und als Übersetzer für wissenschaftliche und IT-Fachliteratur. Freizeit hat er gelegentlich auch, die er gerne gemeinsam mit Freunden oder allein mit sich selbst und einem guten Buch oder einer CD verbringt.

Verwandte Titel von O'Reilly

Statistik von Kopf bis Fuß

R in a Nutshell (*engl.*)

Analyzing Business Data with Excel (*engl.*)

Excel Scientific and Engineering Cookbook (*engl.*)

Access Data Analysis Cookbook (*engl.*)

Weitere Bücher aus O'Reillys Reihe Von Kopf bis Fuß

Java von Kopf bis Fuß

Objektorientierte Analyse & Design von Kopf bis Fuß (OOA & D)

HTML mit CSS & XHTML von Kopf bis Fuß

Entwurfsmuster von Kopf bis Fuß

Servlets & JSP von Kopf bis Fuß

Head First EJB (*engl.*)

Head First PMP (*engl.*)

SQL von Kopf bis Fuß

Softwareentwicklung von Kopf bis Fuß

JavaScript von Kopf bis Fuß

Ajax von Kopf bis Fuß

Head First Physics (*engl.*)

Statistik von Kopf bis Fuß

Head First Rails (*engl.*)

PHP & MySQL von Kopf bis Fuß

Head First Algebra (*engl.*)

Webdesign von Kopf bis Fuß

Netzwerke von Kopf bis Fuß

Head First Excel (*engl., in Vorber.*)

Der Inhalt (im Überblick)

	Einführung	xxv
1	Einführung in die Datenanalyse: <i>Wir zerlegen alles in seine Einzelteile</i>	1
2	Experimente: <i>Überprüfen Sie Ihre Hypothesen</i>	37
3	Optimierung: <i>Holen Sie das Äußerste raus</i>	75
4	Datenvisualisierung: <i>Aus Bildern lernen Sie was</i>	111
5	Hypothesen prüfen: <i>Sag, dass das nicht wahr ist</i>	139
6	Bayes-Statistik: <i>Bloß nicht die Bodenhaftung verlieren!</i>	169
7	Subjektive Wahrscheinlichkeiten: <i>Der Glaube an Zahlen</i>	191
8	Heuristiken: <i>Analysieren wie ein echter Mensch</i>	225
9	Histogramme: <i>Zahlen nehmen Form an</i>	251
10	Regression: <i>Vorhersagen</i>	279
11	Der Zufallsfehler: <i>Ups, daneben!</i>	315
12	Relationale Datenbanken: <i>Sind Sie beziehungsfähig?</i>	359
13	Datenbereinigung: <i>Ordnung erzwingen</i>	385
A	Was übrig bleibt: <i>Die Top Ten der Themen, die wir nicht behandelt haben</i>	417
B	R installieren: <i>Machen Sie R einsatzbereit!</i>	427
C	Excels Erweiterungen aktivieren: <i>Solver und Analyse-Funktionen</i>	431
	Index	435

Der Inhalt (jetzt ausführlich)

Einführung

Ihr Datenanalysten-Gehirn. Sie versuchen, etwas zu *lernen*, und Ihr *Hirn* tut sein Bestes, damit das Gelernte nicht hängen bleibt. Es denkt nämlich: »Wir sollten lieber ordentlich Platz für wichtigere Dinge lassen, zum Beispiel dafür, welche Tiere einem gefährlich werden könnten, oder dass es eine ganz schlechte Idee ist, nackt Snowboard zu fahren.« Tja, wie schaffen Sie es nun, Ihr Gehirn davon zu überzeugen, dass Ihr *Leben* von soliden datenanalytischen Kompetenzen abhängen könnte?

Für wen ist dieses Buch?	xxvi
Wir wissen, was Sie gerade denken	xxvii
Metakognition: Nachdenken übers Denken	xxix
Und das können SIE tun, um sich Ihr Gehirn gefügig zu machen	xxxI
Lies mich!	xxxii
Die Gutachter	xxxiv
Danksagungen	xxxv

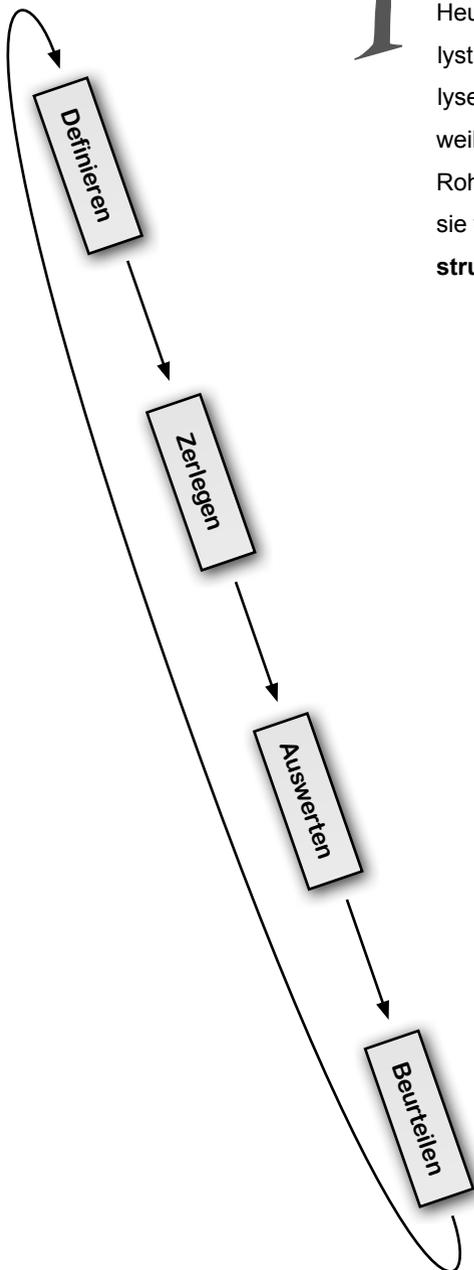
Einführung in die Datenanalyse

Wir zerlegen alles in seine Einzelteile

1

Überall sind Daten.

Heutzutage muss jeder mit Bergen von Daten fertig werden, ob er sich nun »Datenanalytist« nennt oder nicht. Diejenigen allerdings, in deren Werkzeugkasten sich Datenanalyse-Kompetenzen finden, haben einen **entscheidenden Vorsprung** vor allen anderen, weil sie wissen, was man mit all dem Zeug **machen** kann. Sie wissen, wie man aus Rohdaten Informationen gewinnt, mit denen sich **reale Prozesse** steuern lassen, und sie wissen, wie man komplexe Fragestellungen und Datenmengen so **aufbricht und strukturiert**, dass man zum Kern der Probleme im jeweiligen Geschäftsfeld vordringt.



René Sans Kosmetik braucht Ihre Hilfe	2
Der Geschäftsführer würde den Absatz gern mit einer Datenanalyse anschieben	3
Datenanalyse heißt, sorgfältig über die Befundlage nachzudenken	4
Definieren Sie das Problem	5
Ihr Auftraggeber hilft Ihnen, das Problem zu definieren	6
Feedback von René Sans für Sie	8
Brechen Sie Problemstellung und Daten in besser überschaubare Teile auf	9
Sehen Sie sich ein weiteres Mal an, was Sie haben	10
Werten Sie die Teilprobleme aus	13
Eine Analyse haben Sie erst, wenn Sie sich selbst einbringen	14
Geben Sie eine Empfehlung	15
Ihr Bericht ist fertig	16
Der Geschäftsführer mag Ihre Arbeit	17
Gerade erhalten wir Nachricht von einem Zeitungsartikel	18
Sie haben sich von den Einschätzungen des Geschäftsführers in die falsche Richtung schicken lassen	20
Ihre Annahmen und Meinungen zur Realität sind Ihr mentales Modell	21
Ihr statistisches Modell hängt von Ihrem mentalen Modell ab	22
Mentale Modelle sollten immer das einschließen, was Sie nicht wissen	25
Der Geschäftsführer informiert Sie darüber, was er nicht weiß	26
René Sans hat Ihnen gerade eine Riesenliste mit Rohdaten zugeschickt	28
Zeit, die Daten weiter aufzuboahren	31
Der Großhändler Alles für Alle bestätigt Ihnen Ihren Eindruck	32
So sind Sie vorgegangen	35
Ihre Analyse hat Ihrem Auftraggeber zu einer brillanten Entscheidung verholfen	36

Experimente

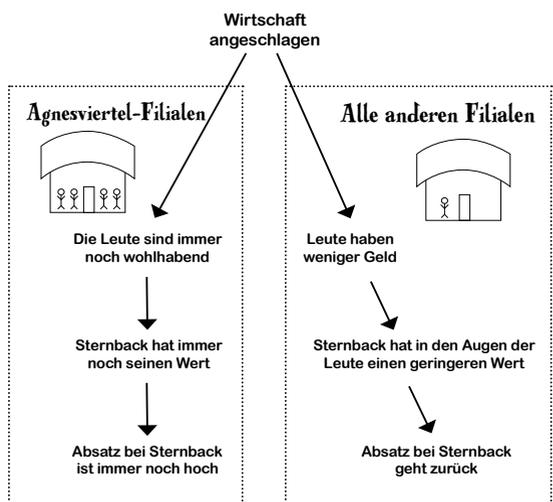
Überprüfen Sie Ihre Hypothesen

2

Können Sie belegen, was Sie glauben?

Mit einer echten **empirischen** Versuchsanordnung? Es geht nichts über ein ordentliches Experiment, um offene Fragen zu klären und zu demonstrieren, wie etwas in der Realität eigentlich zusammenhängt. Anstatt sich ausschließlich auf **anfallende Beobachtungsdaten** zu verlassen, kann ein sauber geplantes und durchgeführtes Experiment oft dabei helfen, **kausale Verbindungen** herzustellen. Mit einer überzeugenden empirischen Datengrundlage werden Ihre analytischen Urteile umso schlagkräftiger.

Kaffee-Rezession!	38
Nächste Vorstandssitzung bei Sternback in drei Monaten	39
Der Sternback-Fragebogen	41
Immer die Methode des Vergleichs benutzen	42
Vergleiche sind das A und O im Umgang mit Beobachtungsdaten	43
Könnte das Wertempfinden den Ertragsrückgang verursacht haben?	44
Denkweise eines typischen Kunden	46
Beobachtungsstudien sind voll von Störvariablen	47
Mögliche Konfundierung der Ergebnisse durch die Standortfrage	48
Bändigen Sie konfundierende Variablen durch Aufteilung der Daten in Blöcke	50
Es ist schlimmer, als wir dachten!	53
Sie müssen ein Experiment durchführen, um die beste Strategie zu finden	54
Der Sternback-Geschäftsführer hat es ziemlich eilig	55
Sternback senkt die Preise	56
Einen Monat später ...	57
Eine Kontrollgruppe verschafft Ihnen eine Baseline	58
Wie man nicht gefeuert wird (oder: Notruf 112)	61
Machen wir noch ein ^{mal ein richtiges!} Experiment	62
Einen Monat später ...	63
Auch Experimente werden von Störvariablen geplagt	64
Konfundierung durch sorgfältige Gruppenbildung vermeiden	65
Durch Randomisierung homogene Gruppen zusammenstellen	67
Ihr Experiment kann starten	71
Die Ergebnisse sind da	72
Sternback verfügt jetzt über eine empirisch überprüfte Verkaufsstrategie	73



Optimierung

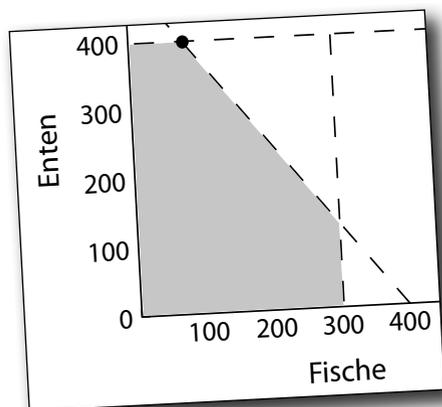
Holen Sie das Äußerste raus

3

Wir wollen immer so viel wie möglich.

Und ständig versuchen wir herauszufinden, wie wir das hinkriegen können. Wenn das, wovon wir möglichst viel wollen – Gewinn, Geld, Effizienz, Geschwindigkeit – *quantifiziert* werden kann, dann stehen die Chancen gut, dass es in der Datenanalyse Instrumente gibt, mit deren Hilfe wir uns unsere **Entscheidungsvariablen** so zurechtbiegen, dass es uns bei der **Bestimmung des optimalen Punkts** hilft, an dem wir von dem, was wir haben wollen, das meiste bekommen. In diesem Kapitel werden Sie eine dieser Methoden kennen lernen, genauso wie die leistungsfähige **Solver-Erweiterung** in Excel, die diese Methode umsetzt.

Diesmal haben Sie es mit Badetieren zu tun	76
Die kontrollierten Variablen werden durch Nebenbedingungen beschränkt	79
Die Entscheidungsvariablen sind das, was sich kontrollieren lässt	79
Sie haben ein Optimierungsproblem	80
Die Zielvorgabe bestimmen Sie mit der Zielfunktion	81
Ihre Zielfunktion	82
Produktmixe für weitere Nebenbedingungen	83
Zeichnen Sie mehrere Nebenbedingungen in dasselbe Diagramm	84
Ihre Wahlmöglichkeiten liegen ausschließlich im zulässigen Bereich	85
Ihre neue Nebenbedingung hat den zulässigen Bereich verändert	87
Ihre Tabellenkalkulation kann optimieren	90
Der Solver erledigt Ihr Optimierungsproblem im Handumdrehen	94
Der Profit geht in den Keller	97
Ihr Modell beschreibt nur, was Sie darin aufnehmen	98
Gleichen Sie Ihre Modellannahmen mit den Zielvorgaben Ihrer Analyse ab	99
Achten Sie auf negativ gekoppelte Variablen	103
Ihr neuer Fertigungsplan funktioniert reibungslos	108
Ihre Annahmen beruhen auf einer sich permanent im Fluss befindenden Realität	109



Datenvisualisierung

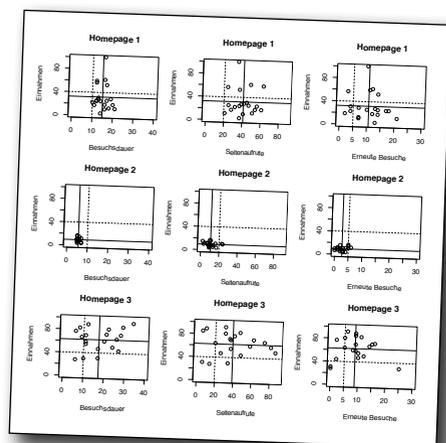
Aus Bildern lernen Sie was

4

Sie brauchen mehr als eine Tabelle voller Zahlen.

Daten sind über alle Maßen **komplex** – mit Variablen wie Sand am Meer! Und dann über Unmengen von Tabellen zu grübeln, ist nicht nur langweilig, es kann tatsächlich reine Zeitverschwendung sein. Sie können sich von einer **klaren grafischen Darstellung** komplexer multivariater Daten auf überschaubarem Raum den Wald zeigen lassen, den Sie vor lauter Bäumen glatt übersehen würden, wenn Sie die ganze Zeit nur auf Datenblätter starrten.

Bei Anziehend anders möchte man die Website optimieren	112
Die Ergebnisse sind drin, und der Informationsgestalter ist raus	113
Der vorige Informationsgestalter hat diese drei Infografiken eingereicht	114
Welche Daten stehen hinter den Grafiken?	115
Wo sind die Daten?!	116
Ein paar Gratisratschläge vom vorigen Informationsgestalter	117
Zu viele Daten sind nie Teil Ihres Problems	118
Daten hübsch aussehen zu lassen, ist ebenfalls nicht Ihr Problem	119
Auch bei der Datenvisualisierung geht es immer um die richtigen Vergleiche	120
Ihre Grafik ist jetzt schon nützlicher als die zurückgewiesenen	123
Nehmen Sie Streudiagramme, wenn Sie Zusammenhänge explorieren	124
Die besten Diagramme sind multivariat	125
Stellen Sie mehr Variablen dar, indem Sie mehrere Diagramme gemeinsam anzeigen	126
Das Diagramm ist hervorragend, der Web-Crack ist aber noch nicht zufrieden	130
Eine gute visuelle Umsetzung hilft Ihnen, über Ursachen nachzudenken	131
Die Versuchsplaner schalten sich ein	132
Die Versuchsplaner haben ihre eigenen Hypothesen	135
Ihr Auftraggeber ist zufrieden mit Ihrer Arbeit	136
Von überall laufen Bestellungen ein!	137



Hypothesen prüfen

Sag, dass das nicht wahr ist

5

Es kann ganz schön schwierig sein, die Realität zu erfassen.

Gerade wenn Sie es mit komplexen, heterogenen Daten zu tun haben, wird es verflucht kompliziert, sich auf **kommende Entwicklungen** einzustellen. Aus diesem Grund begnügen sich Datenanalysten nicht mit **scheinbar naheliegenden Einschätzungen**, die sie der Einfachheit halber für zutreffend halten: Sorgfältiges Abwägen während der Datenanalyse versetzt Sie in die Lage, vorhandene Optionen akribisch zu bewerten und alle verfügbaren Informationen im Modell zu integrieren. Sie sind im Begriff, das Wesen der **Falsifikation** kennenzulernen, eine zunächst wenig intuitive, aber sehr wirkungsvolle Methode, genau das zu erreichen.



Ich brauch was zum Anziehen ...	140
Wann soll die Produktion der neuen Handyskins starten?	141
PodPhone möchte nicht, dass man voraussehen kann, was sie als Nächstes tun	142
Das hier ist alles, was wir wissen	143
Die Hypothese von Skinner passt zu den Daten	144
Skinner ist in den Besitz eines vertraulichen Strategiepapiers gelangt	145
Variablen können positiv oder negativ gekoppelt sein	146
Zusammenhänge sind in der Realität nicht linear, sondern vernetzt	149
Hypothesen zu den Optionen bei PodPhone	150
Sie haben alles, was Sie zur Überprüfung der Hypothesen brauchen	151
Falsifikation – das A und O beim Prüfen von Hypothesen	152
Diagnostizität hilft, die Hypothese mit den wenigsten Gegenargumenten zu finden	160
Nicht jede Hypothese lässt sich ausschließen, Sie können aber festlegen, welche stärker ist als die anderen	163
Sie haben gerade eine Bildnachricht bekommen ...	164
Es ist raus!	167

Bayes-Statistik

Bloß nicht die Bodenhaftung verlieren!

Sie erfassen ständig neue Daten.

Und Sie müssen sichergehen, dass jede Ihrer Analysen alle verfügbaren Daten einbezieht, die für das Problem relevant sind. Sie haben gelernt, wie man im Umgang mit heterogenen Daten die *Falsifikation* zu Hilfe nimmt, aber wie sieht es mit **richtiggehenden Wahrscheinlichkeiten** aus? Die Antwort darauf ist eine äußerst praktische Methode, die **Bayes-Regel**, mit deren Hilfe Sie **Basisraten** berücksichtigen und zu nicht ganz selbstverständlichen Einsichten über Daten kommen, die sich permanent ändern.

Schlechte Nachrichten von Ihrem Arzt	170
Sehen wir uns die Analyse zur Testsicherheit Aussage für Aussage an	173
Wie stark ist die Leguangrippe tatsächlich verbreitet?	174
Sie haben die Falsch-Positiven berechnet	175
Alle diese Begriffe beschreiben bedingte Wahrscheinlichkeiten	176
Sie müssen alle falsch-positiven, richtig-positiven, falsch-negativen und richtig-negativen Fälle durchzählen	177
1% der Bevölkerung hat Leguangrippe	178
Ihr Risiko, an Leguangrippe erkrankt zu sein, ist im Grunde ziemlich gering	181
Denken Sie in simplen ganzen Zahlen über komplexe Wahrscheinlichkeiten nach	182
Die Bayes-Regel kümmert sich um Basisraten und bedingte Wahrscheinlichkeiten	182
Sie können die Bayes-Regel wiederholt anwenden	183
Der zweite Test ist negativ	184
Die Testsicherheit beim neuen Test ist anders	185
Neue Informationen können Ihre Basisrate verändern	186
Große Erleichterung!	189

Röchel



Subjektive Wahrscheinlichkeiten

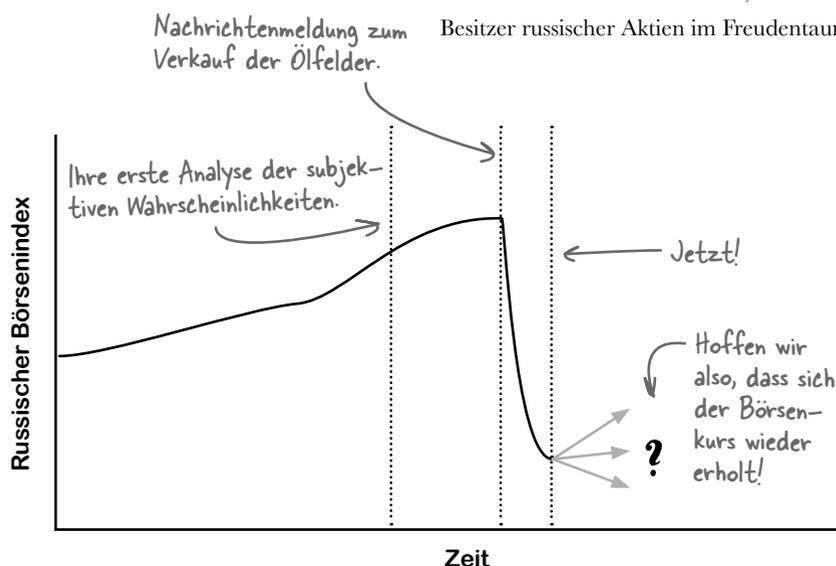
Der Glaube an Zahlen

7

Manchmal ist es ganz gut, sich Zahlen auszudenken.

Im Ernst. Aber nur solange Sie damit Ihre eigenen Überlegungen und Überzeugungen ausdrücken. Mit **subjektiven Wahrscheinlichkeiten** können Sie *schwammige Schätzungen* auf sehr direkte Art und Weise handfest und greifbar machen – Sie werden gleich sehen, wie das geht. Und ganz nebenbei erfahren Sie, wie man mit der **Standardabweichung** die Streuung in Daten einschätzt. Als Special Guest wird außerdem eine der besonders leistungsfähigen analytischen Methoden vorbeischaun, die Sie bereits kennengelernt haben.

Bei Terra Inco Invest ist man auf Ihre Hilfe angewiesen	192
Die Analysten gehen sich gegenseitig an den Hals	193
Subjektive Wahrscheinlichkeiten als Ausdruck von Expertenmeinungen	198
Vielleicht zeigen sich in den subjektiven Wahrscheinlichkeiten gar keine Meinungsverschiedenheiten	199
Die Analysten haben ihre subjektiven Wahrscheinlichkeiten eingereicht	201
Der Geschäftsführer sieht nicht, worauf Sie hinauswollen	202
Der Geschäftsführer mag Ihre Arbeit	207
Die Standardabweichung misst, wie weit Datenpunkte vom Mittelwert abweichen	208
Diese Nachricht trifft Sie aus heiterem Himmel	213
Die Bayes-Regel eignet sich gut zum Revidieren subjektiver Wahrscheinlichkeiten	217
Der Geschäftsführer weiß exakt, was er mit dieser Information machen muss	223
Besitzer russischer Aktien im Freudentaumel!	224



Heuristiken

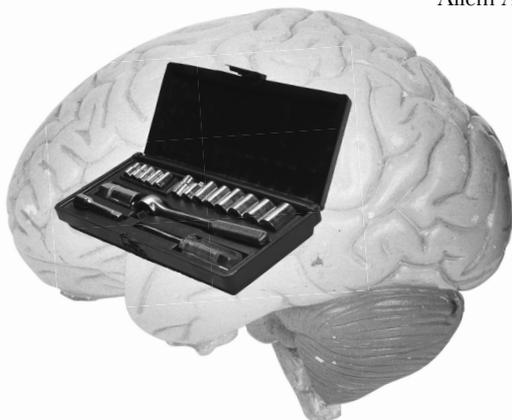
8

Analysieren wie ein echter Mensch

Die Realität konfrontiert Sie mit mehr Variablen, als Sie bewältigen können.

Es wird immer Daten geben, an die Sie nicht rankommen. Und *selbst wenn* Sie Daten zu mehr oder weniger allem haben, was Sie zu verstehen versuchen, sind passende Analysestrategien oft **schwer umsetzbar** und **zeitaufwendig**. Zum Glück ist die Art und Weise, wie Sie im normalen Leben denken, kein »rationales Berechnen des Optimums« – stattdessen werden unvollständige und unsichere Informationen mithilfe von **Näherungsregeln** verarbeitet, um schnell Entscheidungen treffen zu können. Das Faszinierende daran ist, dass diese Regeln **tatsächlich funktionieren** – deshalb sind sie wichtige (und notwendige) Hilfsmittel für den Datenanalysten.

Die Abfall-Scouts haben dem Stadtrat ihren Bericht vorgelegt	226
Die Abfall-Scouts haben das Städtchen richtig auf Vordermann gebracht	227
Die Abfall-Scouts haben die Effektivität ihrer Kampagne erfasst	228
Der Auftrag lautet, das Abfallaufkommen zu reduzieren	229
Die Höhe des Abfallaufkommens ist unmöglich zu messen	230
Stellen Sie eine schwierige Frage, und man wird eine leichtere beantworten	231
Hinter der Datenreuther Abfallfrage steht ein komplexes System	232
Ein umfassendes Abfallmessmodell lässt sich weder planen noch umsetzen	233
Heuristiken sind der Mittelweg zwischen reinem Bauchgefühl und Optimierung	236
Nehmen Sie einen schnellen und sparsamen Baum	239
Gibt es eine einfachere Möglichkeit, den Erfolg der Abfall-Scouts zu beurteilen?	240
Auch Stereotype sind Heuristiken	244
Ihre Analyse ist bereit zur Vorlage	246
Allem Anschein nach hat Ihre Analyse den Stadtrat beeindruckt	249



Histogramme

Zahlen nehmen Form an



Wie viel kann Ihnen ein Balkendiagramm mitteilen?

Es gibt schätzungsweise an die hunderttausend Möglichkeiten, **Daten grafisch darzustellen**, eine davon ist allerdings etwas Besonderes. **Histogramme**, die eine gewisse Ähnlichkeit mit Balkendiagrammen haben, sind eine superschnelle und einfache Methode, Daten zusammenzufassen. Diese leistungsfähigen kleinen Diagramme werden Sie in Kürze zur Darstellung von **Verteilung, zentraler Tendenz, Streuung und anderen Dingen** in Ihren Daten einsetzen. Egal wie umfangreich Ihre Datenreihe ist, wenn Sie ein Histogramm daraus machen, können Sie quasi »sehen«, was sich in ihrem Inneren abspielt. Und ... Sie stehen kurz davor, das alles mit einer **frei verfügbaren Software** zu machen, die so leistungsstark ist, dass Sie feuchte Augen bekommen werden.

Ihre jährliche Leistungsbeurteilung erwartet Sie	252
Nach mehr Geld zu fragen, könnte auf Verschiedenes hinauslaufen	254
Ein paar Daten mit Gehalts- und Honorarerhöhungen	255
Histogramme bilden gruppierte Häufigkeiten ab	262
Lücken zwischen Histogrammbalken entsprechen Lücken in der Werteverteilung	263
R installieren und in Betrieb nehmen	264
Laden Sie Ihre Daten in R	265
R macht schöne Histogramme	266
Legen Sie Histogramme aus Untergruppen Ihrer Daten an	271
Verhandeln zahlt sich aus	276
Was würde verhandeln für Sie bedeuten?	277



10

Regression

Vorhersagen

Machen Sie mal eine Prognose.

Die **Regressionsrechnung** ist ein unglaublich mächtiges statistisches Instrument, das – bei korrektem Einsatz – die Fähigkeit besitzt, Ihnen die **Vorhersage bestimmter Werte** zu ermöglichen; in kontrollierten experimentellen Anordnungen könnten Sie sich dadurch sogar in die Lage versetzen lassen, die Zukunft vorherzusagen. In der Wirtschaft wird sie bis zum Abwinken dazu benutzt, Modelle zur **Aufklärung** von Kunden- und Kaufverhalten zu entwickeln. Sie werden bald erfahren, dass der umsichtige Einsatz der Regressionsrechnung tatsächlich sehr einträglich sein kann.

Was werden Sie jetzt mit all dem Geld machen?	280
Zu analysieren, was man fordern sollte, könnte heftig werden	283
Sieh mal an ... ein Gehaltsschätzer!	284
Der Algorithmus muss eine Vorschrift zur Vorhersage von Gehaltserhöhungen enthalten	286
Streudiagramme, die zwei Variablen vergleichen	292
Eine Gerade könnte Ihren Kunden sagen, worauf sie abzielen sollten	294
Sagen Sie mit dem Mittelwertsgraphen Werte in allen Bändern voraus	297
Die Regressionsgerade sagt die bewilligten Gehaltserhöhungen voraus	298
Die Gerade ist nützlich, wenn Ihre Daten ausreichend hoch korrelieren	300
Für präzise Vorhersagen brauchen Sie eine Formel	304
Was passiert, wenn R ein Regressionsobjekt anlegt?	306
Regressionsgleichung und Streudiagramm arbeiten Hand in Hand	309
Die Regressionsgleichung ist der Algorithmus für den Gehaltsschätzer	310
Ihr Gehaltsschätzer hat nicht ganz wie geplant funktioniert ...	313



DER GEHALTSSCHÄTZER

Was passiert, wenn Sie für Ihre Gehaltserhöhung einen bestimmten Prozentsatz vorschlagen?
Mit dieser Gleichung finden Sie es heraus:

$$y = 2,3 + 0,7x$$

x entspricht Ihrem Vorschlag,
y ist der Satz, den Sie erwarten können.



Der Zufallsfehler

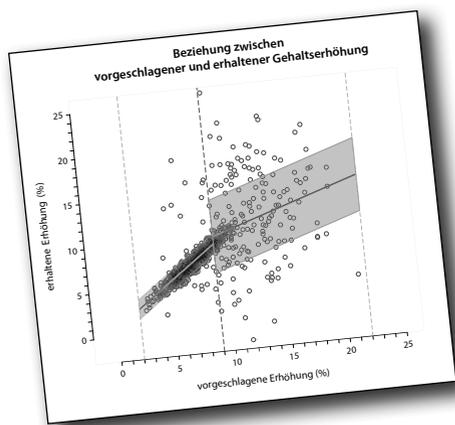
Ups, daneben!

11

Das Leben ist ein furchtbares Durcheinander.

Es sollte Sie deshalb nicht allzu sehr überraschen, dass Ihre Vorhersagen eher im Ausnahmefall exakt im Ziel landen. Wenn Sie Ihre Prognosen allerdings **inklusive Irrtumsspielraum** anbieten, wissen Ihre Kunden nicht nur über die durchschnittliche Ausprägung des vorhergesagten Werts Bescheid, sondern gleichzeitig auch, mit welchen Abweichungen normalerweise zu rechnen ist. Sobald Sie den Zufallsfehler kenntlich machen, stellen sich Ihre Vorhersagen und Überzeugungen aus einer ganz anderen Perspektive dar. Und mit den Methoden in diesem Kapitel lernen Sie darüber hinaus, wie sich der Zufallsfehler **unter Kontrolle** bringen lässt, indem man ihn so klein wie möglich macht und damit die Schwankungsbreite einer Vorhersage gering hält.

Sie haben Ihre Kunden ganz schön aufgemischt	316
Was hat Ihr Gehaltsschätzer-Algorithmus gemacht?	317
Kunden-Untergruppen	318
Der Typ, der 25% mehr gefordert hat, liegt außerhalb unseres Modells	321
Vom Umgang mit Kunden, die Vorhersagen außerhalb des Geltungsbereichs wollen	322
Der Typ, der wegen der Extrapolation gefeuert wurde, hat sich wieder beruhigt	327
Sie haben lediglich einen Teil des Grundproblems gelöst	328
Wie sehen die Daten zu diesen chaotischen Ergebnissen aus?	329
Zufallsfehler sind Abweichungen von Ihrer Vorhersage	330
Der Zufallsfehler hilft sowohl Ihnen als auch Ihren Kunden	334
Beziffern Sie den Zufallsfehler	336
Messen Sie die Residuenverteilung mit dem Standardschätzfehler	337
Das von R berechnete Modell kennt den Standardschätzfehler bereits	338
Die Zusammenfassung Ihres linearen Modells in R liefert Ihnen den Standardschätzfehler	340
Datensegmentierung hat immer den Umgang mit dem Zufallsfehler zum Ziel	346
Gute Regressionsmodelle balancieren Datenaufklärung und Vorhersagegüte aus	350
Die Modelle für die Datensegmente funktionieren besser als die ursprünglichen	352
Ihre Kunden kommen scharenweise zurück	357



Relationale Datenbanken

Sind Sie beziehungsfähig?

12

Wie organisiert man Daten, die ganz schlimm multivariat sind?

Eine Tabelle – beispielsweise eine in Ihrer Tabellenkalkulation – hat *nur zwei* Dimensionen: Zeilen und Spalten. Wenn Sie jetzt aber Daten haben, die diverse Dimensionen mehr aufweisen, sieht **das herkömmliche Tabellenformat** sehr schnell sehr alt aus. In diesem Kapitel werden Sie hautnah erleben, wann das Spreadsheet-Format von Tabellenkalkulationen es Ihnen wirklich schwer macht, multivariate Daten zu bändigen, und warum es mit einem **relationalen Datenbankmanagementsystem** einfach ist, Unmengen multivariater Datensätze abzuspeichern und wieder abzurufen.

Bei der Datenreuther Depesche möchte man die verkaufte Auflage analysieren	360
Diese Daten dokumentieren ihren Geschäftsverlauf	361
Sie müssen wissen, wie die Datentabellen zusammenhängen	362
Datenbank = Sammlung von Informationen mit klar definierten Beziehungen	365
An Ihre Information kommen Sie, wenn Sie einen Pfad entlang der Beziehungen anlegen	366
Legen Sie eine Arbeitsmappe an, die diesen Pfad abbildet	366
Ihre Zusammenfassung verknüpft Artikelanzahl und verkaufte Auflage	371
Allem Anschein nach kommt Ihr Streudiagramm gut an	374
Das war schon übel, diese ganzen Daten zusammenkopieren zu müssen	375
Relationale Datenbanken nehmen Ihnen die Verwaltung von Beziehungen ab	376
Bei der Depesche hat man nach Ihrem Diagramm eine Datenbank eingerichtet	377
Die Depesche hat Ihnen Daten mit einer SQL-Abfrage zusammengestellt	379
Endlose Auswertungsmöglichkeiten für Daten aus einer relationalen Datenbank	382
Die Person auf dem Cover – das sind SIE	383

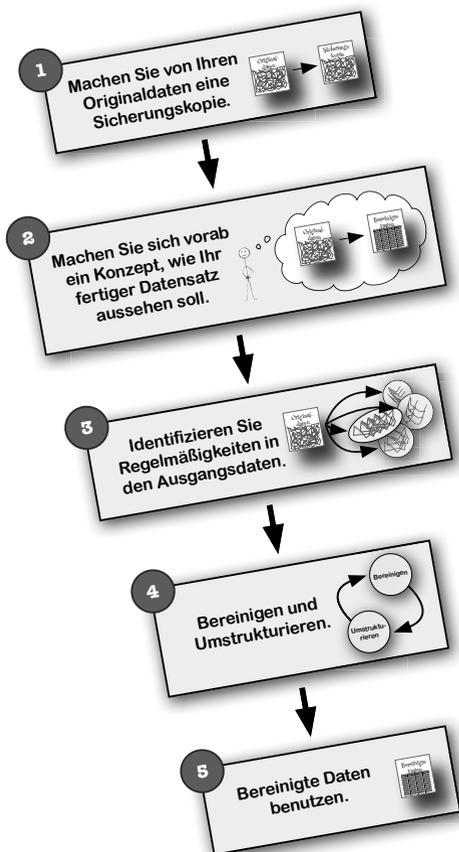


13

Datenbereinigung Ordnung erzwingen

Ihre Daten sind nutzlos ...

... solange sie unstrukturiert sind. Und eine Menge Leute, die **Daten zusammentragen**, leisten lausige Arbeit beim Anlegen einer angemessenen Struktur. Wenn Ihre Daten nicht vernünftig strukturiert sind, ist an Gruppieren und Vergleichen nicht zu denken, Formeln können Sie auch keine anwenden, und wenn Sie Pech haben, bekommen Sie nicht einmal etwas Vernünftiges zu sehen. Solche Daten könnten Sie genauso gut ignorieren, stimmt's? Genau besehen, ginge das alles besser. Wenn Sie eine **klare Vorstellung** davon haben, welche Struktur Sie brauchen, und mit ein paar **Mitteln zum Manipulieren von Textdaten** können Sie auch die größte vorstellbare Zumutung eines Datendurcheinanders so **einrichten**, das tatsächlich irgendetwas Brauchbares dabei herauskommt.



Gerade erhalten Sie den Kandidatenstamm einer aufgelösten Personalberatungsfirma	386
Die schmutzige Seite der Datenanalyse	387
O'Lymp braucht die Liste für seine Vermittler	388
Datenbereinigung hängt entscheidend von der Vorbereitung ab	392
Sobald Ihr Konzept steht, können Sie mit der eigentlichen Datenbereinigung beginnen	393
Nehmen Sie die Raute als Spaltentrenner	394
Excel hat die Daten mit dem Trenner auf die Spalten verteilt	395
Beseitigen Sie das Caret mit der WECHSELN-Funktion	399
Sie haben gerade die Vornamen bereinigt	400
Für die WECHSELN-Funktion sind die Nachnamen zu komplex	402
Komplexe Muster mit verschachtelten Textfunktionen angehen	403
Komplexe Muster mit regulären Ausdrücken in R aufbrechen	404
Der Aufruf der sub()-Funktion hat die Nachnamen repariert	406
Sie können Ihrem Auftraggeber die Daten jetzt liefern	407
Möglich, dass Sie noch nicht ganz fertig sind ...	408
Sortieren Sie Ihre Daten, dann stehen die Dubletten untereinander	409
Die Daten könnten aus einer relationalen Datenbank stammen	412
Entfernen Sie die doppelten Namenseinträge	413
Jetzt sind alle Einträge im Datensatz sauber und eindeutig	414
Die Rekrutierungsaktion bei O'Lymp ist ein Wahnsinnserefolg!	415
Zeit, Abschied zu nehmen ...	416

Was übrig bleibt

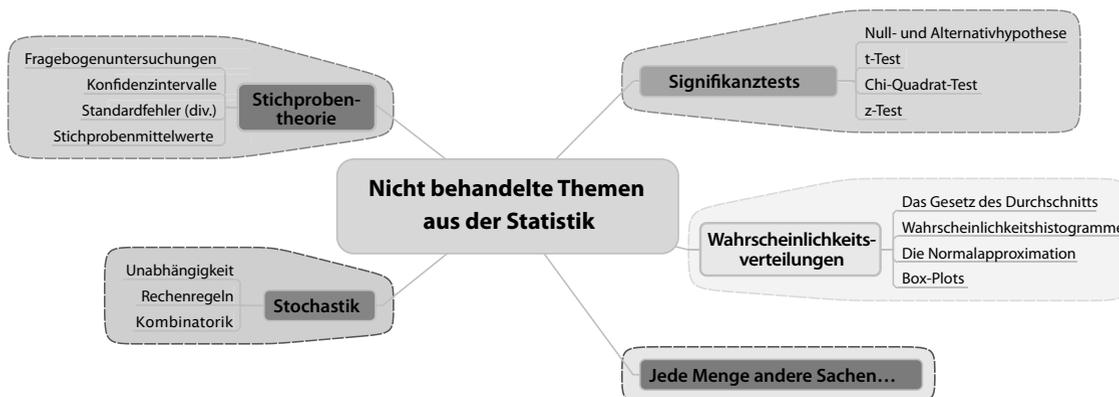
Die Top Ten der Themen, die wir nicht behandelt haben



Wir haben ein ganz schönes Stück des Weges zurückgelegt.

Allerdings ist die Datenanalyse ein weites, **sich permanent weiterentwickelndes Arbeitsfeld**, und es ist wirklich viel Lernstoff übrig geblieben. In diesem Anhang gehen wir mit Ihnen **zehn Punkte** durch, für deren Behandlung zwar im Buch kein Raum mehr war, die aber **ganz oben auf Ihre Liste** der Dinge gehören, mit denen Sie sich als Nächstes befassen sollten.

- | | |
|---|-----|
| 1. Alles, was irgendwie mit Statistik zu tun hat | 418 |
| 2. Excel-Kenntnisse | 419 |
| 3. Edward Tufte und seine Prinzipien der Informationsvisualisierung | 420 |
| 4. Pivot-Tabellen | 421 |
| 5. Das R-Ökosystem | 422 |
| 6. Nicht-lineare und multiple Regression | 423 |
| 7. Statistisches Hypothesentesten | 424 |
| 8. Zufall | 424 |
| 9. Google-Docs | 425 |
| 10. Ihre Fachkompetenz | 426 |



R installieren

Machen Sie R einsatzbereit!

B

Unter der Haube ist die geballte Leistungsfähigkeit dieses Datenakrobats enorm komplex.

R zu installieren und in Betrieb zu nehmen, ist allerdings etwas, das Sie in ein paar Minuten erledigt haben, und in diesem Anhang wollen wir Ihnen zeigen, wo und wie Sie ohne Reibungsverluste an Ihre eigene R-Installation kommen.

Bringen Sie R ans Laufen

428

The screenshot shows the R Project for Statistical Computing website. The browser address bar displays `http://www.r-project.org/`. The page features the R logo and a navigation menu on the left with categories like 'About R', 'Download Packages', 'R Project Foundation', 'Documentation', and 'Misc'. The main content area is titled 'The R Project for Statistical Computing' and displays several statistical visualizations: a PCA plot of 5 variables (Family, Examination, Education, Catholic, Agriculture) with a (1-3) 60% label; a Clustering dendrogram with 4 groups; and two histograms for Factor 1 [41%] and Factor 3 [19%]. Below the plots is a 'Getting Started' section with bullet points: 'R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To download R, please choose your preferred CRAN mirror.' and 'If you have questions about R like how to download and install the software, or what the license terms are, please read our answers to frequently asked questions before you send an email.' A 'News' section at the bottom lists: 'R version 2.10.1 has been released on 2009-12-14. The source code will first become available in this directory, and eventually via all of CRAN.', 'The first issue of The R Journal is now available', and 'useR! 2010, the R user conference, will be held at NIST, Gaithersburg, Maryland, USA, July 21-23, 2010.'

Excels Erweiterungen aktivieren



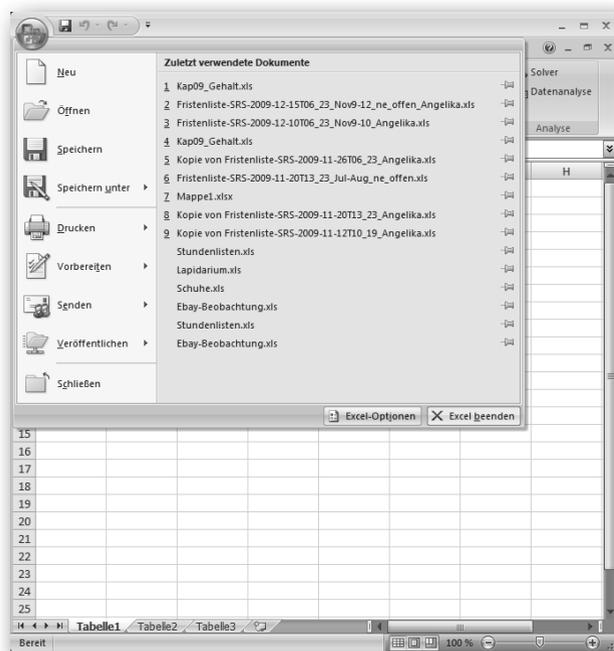
Solver und Analyse-Funktionen

Ein paar der interessantesten Funktionen von Excel werden standardmäßig nicht mit installiert.

Genau so ist es. Um die Optimierungsprobleme in Kapitel 3 und die Histogramme in Kapitel 9 bearbeiten zu können, müssen Sie den **Solver** und die **Analyse-Funktionen** aktivieren, zwei Erweiterungen, die je nach Plattform und Programmversion während der Erstinstallation nicht unbedingt in Excel aktiviert werden, es sei denn, Sie greifen ausdrücklich in die Konfiguration ein.

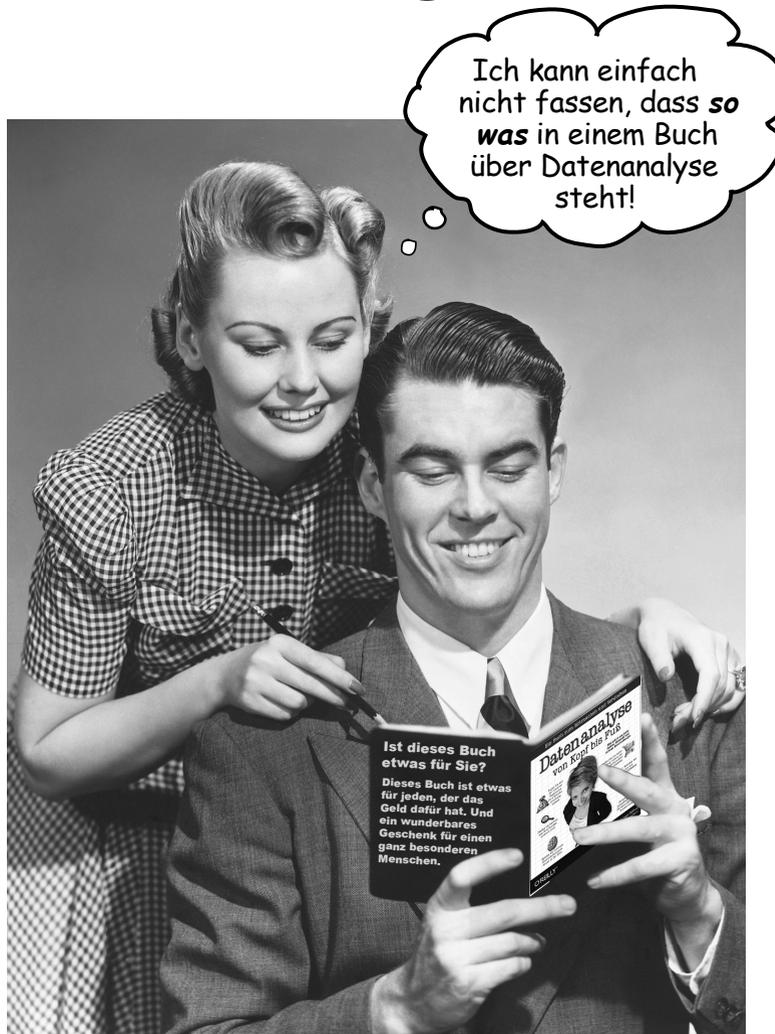
Installieren Sie die Datenanalysefunktionen in Excel

432



Wie man dieses Buch benutzt

Einführung



In diesem Abschnitt beantworten wir die brennende Frage: »Und? Warum STEHT so was in einem Buch über Datenanalyse?«

Für wen ist dieses Buch?

Wenn Sie *alle* folgenden Fragen mit »Ja« beantworten können ...

- 1 Haben Sie das Gefühl, dass in Ihren Daten ungeahnte Einsichten und Erkenntnisse vergraben liegen könnten, an die Sie aber nur mit passendem Werkzeug herankommen?
- 2 Wollen Sie lernen, verstehen und abrufen können, wie man brillante Diagramme anlegt, Hypothesen prüft, ein Regressionsmodell anpasst und chaotische Daten bereinigt?
- 3 Ziehen Sie eine anregende Unterhaltung beim Abendessen einem trockenen, langweiligen Vortrag vor?

... dann ist dieses Buch etwas für Sie.

Wer sollte eher die Finger von diesem Buch lassen?

Wenn Sie *eine* der folgenden Fragen mit »Ja« beantworten müssen ...

- 1 Sind Sie ein Topdatenanalyst mit Erfahrung, der nach einer **Quelle brandaktueller Techniken und Strategien** sucht?
- 2 Haben Sie noch nie mit Microsoft Excel oder OpenOffice Calc gearbeitet?
- 3 Haben Sie **Angst, etwas Neues auszuprobieren**? Ist Ihnen eine Wurzelkanalbehandlung lieber, als Streifen kombiniert mit Karos zu tragen? Glauben Sie, dass ein Fachbuch über Datenanalyse, in dem Kontrollbedingungen und Zielfunktionen vermenschlicht werden, nicht seriös sein kann?

... dann ist dieses Buch *nicht* das richtige für Sie.



[Anmerkung der Marketing-Abteilung: Dieses Buch ist etwas für jeden, der eine Kreditkarte besitzt oder im Laden sein Wechselgeld zählen kann.]

Wir wissen, was Sie gerade denken.

»Kann *das* wirklich ein seriöses Buch über Datenanalyse sein?«

»Was ist mit den ganzen *Abbildungen*?«

»Kann ich das auf diese Art wirklich *lernen*?«

Und wir wissen, was Ihr Gehirn gerade denkt.

Ihr Gehirn lechzt nach Neuem. Es ist ständig dabei, Ihre Umgebung abzusuchen, und es *wartet* auf etwas Ungewöhnliches. So ist es nun einmal gebaut, und es hilft Ihnen zu überleben.

Also, was macht Ihr Gehirn mit all den gewöhnlichen, normalen Routinesachen, denen Sie begegnen? Es tut alles in seiner Macht stehende, um dadurch nicht bei seiner *eigentlichen* Arbeit gestört zu werden: Dinge zu erfassen, die *wirklich* wichtig sind. Es gibt sich nicht damit ab, die langweiligen Sachen zu speichern, sondern lässt diese gar nicht erst durch den »Das-ist-offensichtlich-nicht-wichtig«-Filter.

Woher *weiß* Ihr Gehirn denn, was wichtig ist? Nehmen Sie an, Sie machen einen Tagesausflug und ein Tiger springt vor Ihnen aus dem Gebüsch: Was passiert dabei in Ihrem Kopf und Ihrem Körper?

Neuronen feuern. Gefühle werden angekurbelt. *Botenstoffe werden ausgeschüttet.*

Und so weiß Ihr Gehirn:

Das hier muss wichtig sein! Vergiss es nicht!

Aber nun stellen Sie sich vor, Sie sind zu Hause oder in einer Bibliothek. In einer sicheren, warmen, tigerfreien Zone. Sie lernen. Bereiten sich auf eine Prüfung vor. Oder Sie versuchen, irgendein schwieriges Thema zu erarbeiten, von dem Ihr Chef glaubt, Sie bräuchten dafür eine Woche oder höchstens zehn Tage.

Da ist nur ein Problem: Ihr Gehirn versucht Ihnen einen großen Gefallen zu tun. Es versucht, dafür zu sorgen, dass diese *offensichtlich* unwichtigen Inhalte nicht knappe Ressourcen blockieren. Ressourcen, die besser dafür verwendet würden, die *wirklich wichtigen* Dinge zu speichern. Wie Tiger. Wie die Gefahren des Feuers. Oder dass Sie nie, nie, niemals diese Partyfotos auf Ihre Facebook-Seite hätten stellen dürfen.

Und es gibt keine einfache Möglichkeit, Ihrem Gehirn zu sagen: »Hey, Gehirn, vielen Dank, aber egal, wie langweilig dieses Buch auch ist und wie klein der Ausschlag auf meiner emotionalen Richterskala gerade ist, aber ich *will wirklich*, dass du dir diesen Kram merkst.«



Wir stellen uns unseren Leser als einen aktiv Lernenden vor.

Also, was ist nötig, damit Sie etwas *lernen*? Erst einmal müssen Sie es aufnehmen und dann dafür sorgen, dass Sie es nicht wieder vergessen. Es geht nicht darum, Fakten in Ihren Kopf zu schieben. Nach den neuesten Forschungsergebnissen aus Kognitionswissenschaft, Neurobiologie und pädagogischer Psychologie gehört zum Lernen viel mehr als nur Text auf einer Seite. Wir wissen, was Ihr Gehirn anmacht.

Einige der Lernprinzipien dieser Buchreihe:

Bilder einsetzen. An Bilder kann man sich viel besser erinnern als an Worte allein und lernt so viel effektiver (bis zu 89 Prozent Verbesserung bei Gedächtnisabruf- und Lerntransferstudien). Außerdem werden die Dinge **dadurch verständlicher, Text in oder neben die Grafiken setzen**, auf die sie sich beziehen, anstatt darunter oder auf eine andere Seite. Die Leser werden auf den Bildinhalt bezogene Probleme dann mit *doppelt* so hoher Wahrscheinlichkeit lösen können.



Wir verwenden einen gesprächsorientierten Stil mit persönlicher Ansprache. Nach neueren Untersuchungen haben Studenten nach dem Lernen bei Tests bis zu 40 Prozent besser abgeschnitten, wenn der Inhalt den Leser direkt in der ersten Person und im lockeren Stil angesprochen hat statt in einem formalen Ton. Wir halten keinen Vortrag, sondern erzählen Geschichten. Wir benutzen eine zwanglose Sprache. Nehmen Sie sich selbst nicht zu ernst. Würden Sie einer anregenden Unterhaltung beim Abendessen mehr Aufmerksamkeit schenken oder einem Vortrag?

Wir bringen den Lernenden dazu, intensiver nachzudenken. Mit anderen Worten: Falls Sie nicht aktiv Ihre Neuronen strapazieren, passiert in Ihrem Gehirn nicht viel. Ein Leser muss motiviert, begeistert und neugierig sein und angeregt werden, Probleme zu lösen, Schlüsse zu ziehen und sich neues Wissen anzueignen. Und dafür brauchen Sie Herausforderungen, Übungen, zum Nachdenken anregende Fragen und Tätigkeiten, durch die beide Seiten des Gehirns und mehrere Sinne beansprucht werden.



Wir wollen Ihre Aufmerksamkeit – und sie behalten. Wir alle haben schon Erfahrungen dieser Art gemacht: »Ich will das wirklich lernen, aber ich kann einfach nicht über Seite 1 hinaus wach bleiben.« Ihr Gehirn passt auf, wenn Dinge ungewöhnlich, interessant, merkwürdig, auffällig, unerwartet sind. Ein neues, schwieriges, technisches Thema zu lernen, muss nicht langweilig sein. Wenn es das nicht ist, lernt Ihr Gehirn viel schneller.

Wir sprechen Gefühle an. Wir wissen, dass Ihre Fähigkeit, sich an etwas zu erinnern, wesentlich von dessen emotionalem Gehalt abhängt. Sie erinnern sich *an* das, was Sie *bewegt*. Sie erinnern sich, wenn Sie etwas *fühlen*. Nein, wir erzählen keine herzerreißenden Geschichten über einen Jungen und seinen Hund. Was wir erzählen, ruft Überraschungs-, Neugier-, Spaß- und »Was-soll-das?«-Emotionen hervor und dieses Triumphgefühl, das Sie beim Lösen eines Puzzles empfinden oder wenn Sie etwas lernen, das alle anderen schwierig finden. Oder wenn Sie merken, dass Sie etwas können, was dieser »Ich-bin-technisch-echt-vielgabter-Typ« aus Ihrem Seminar über rechnergestützte Statistik *nicht* kann.



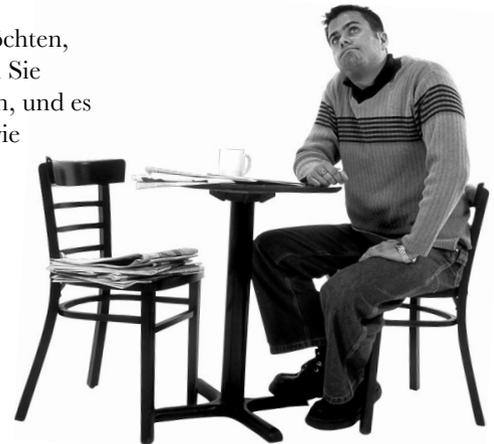
Metakognition: Nachdenken übers Denken

Wenn Sie wirklich lernen möchten, und zwar schneller und nachhaltiger, dann schenken Sie Ihrer Aufmerksamkeit Aufmerksamkeit. Denken Sie darüber nach, wie Sie denken. Lernen Sie, wie Sie lernen.

Die meisten von uns haben in ihrer Jugend keine Kurse in Metakognition oder Lerntheorie gehabt. Es wurde von uns *erwartet*, dass wir lernen, aber nur selten wurde uns auch *beigebracht*, wie man lernt.

Wir nehmen aber an, dass Sie wirklich etwas über Datenanalyse lernen möchten, wenn Sie dieses Buch in den Händen halten. Und wahrscheinlich möchten Sie nicht viel Zeit aufwenden. Und Sie wollen sich an das *erinnern*, was Sie lesen, und es anwenden können. Und deshalb müssen Sie es *verstehen*. Wenn Sie so viel wie möglich von diesem Buch profitieren wollen oder von *irgendeinem* anderen Buch oder einer anderen Lernerfahrung, übernehmen Sie Verantwortung für Ihr Gehirn. Ihr Gehirn im Zusammenhang mit **diesem** Lernstoff.

Der Trick besteht darin, Ihr Gehirn dazu zu bringen, neuen Lernstoff als etwas wirklich Wichtiges anzusehen. Als entscheidend für Ihr Wohlbefinden. So wichtig wie ein Tiger. Andernfalls stecken Sie in einem dauernden Kampf, in dem Ihr Gehirn sein Bestes gibt, um die neuen Inhalte davon abzuhalten, hängen zu bleiben.



Wie bringen Sie also Ihr Gehirn dazu, Datenanalyse für so wichtig zu halten wie einen beutehungrigen Tiger?

Da gibt es den langsamen, ermüdenden Weg oder den schnelleren, effektiveren Weg. Der langsame Weg geht über bloße Wiederholung. Natürlich ist Ihnen klar, dass Sie lernen und sich sogar an die langweiligsten Themen erinnern *können*, wenn Sie sich die gleiche Sache immer wieder einhämmern. Wenn Sie nur oft genug wiederholen, sagt Ihr Gehirn: »Er hat zwar nicht das *Gefühl*, dass das wichtig ist, aber er sieht sich dieselbe Sache *immer und immer wieder* an – also muss sie wohl wichtig sein, nehme ich an.«

Der schnellere Weg besteht darin, **alles zu tun, was die Gehirnaktivität erhöht**, vor allem *verschiedene* Arten von Gehirnaktivität. Eine wichtige Rolle dabei spielen die auf der vorhergehenden Seite erwähnten Dinge – alles Dinge, die nachweislich helfen, dass Ihr Gehirn *für* Sie arbeitet. So hat sich beispielsweise in Untersuchungen gezeigt: Wenn Wörter *in* den Abbildungen stehen, die sie beschreiben (und nicht irgendwo an anderer Stelle auf der Seite, zum Beispiel in einer Bildunterschrift oder im Text), versucht Ihr Gehirn, herauszufinden, wie die Wörter und das Bild zusammenhängen, und dadurch feuern mehr Neuronen. Und je mehr Neuronen feuern, umso größer ist die Chance, dass Ihr Gehirn *mitbekommt*: Bei dieser Sache lohnt es sich, aufzupassen, und vielleicht auch, sich daran zu erinnern.

Ein lockerer Sprachstil hilft, denn Menschen tendieren zu erhöhter Aufmerksamkeit, wenn ihnen bewusst ist, dass sie ein Gespräch führen – man erwartet dann ja von ihnen, dass sie dem Gespräch folgen und sich beteiligen. Das Erstaunliche daran ist: Es ist Ihrem Gehirn ziemlich egal, dass die »Unterhaltung« zwischen Ihnen und einem Buch stattfindet! Wenn der Schreibstil dagegen formal und trocken ist, hat Ihr Gehirn den gleichen Eindruck wie in einem Vortrag, bei dem Sie in einem Raum zusammen mit vielen anderen passiven Zuhörern sitzen. Nicht nötig, wach zu bleiben.

Aber Abbildungen und ein lockerer Sprachstil sind erst der Anfang.

Das haben WIR getan:

Wir haben **Bilder** verwendet, weil Ihr Gehirn auf visuelle Eindrücke eingestellt ist, nicht auf Text. Soweit es Ihr Gehirn betrifft, sagt ein Bild *wirklich* mehr als 103 Worte. Und dort, wo Text und Abbildungen zusammenwirken, haben wir den Text *in* die Bilder eingebettet, denn Ihr Gehirn arbeitet besser, wenn der Text *innerhalb* der Sache steht, auf die er sich bezieht, und nicht in einer Bildunterschrift oder irgendwo vergraben im Text.

Wir haben **Redundanz** eingesetzt, das heißt dasselbe auf *unterschiedliche* Art und mit verschiedenen Medientypen ausgedrückt, damit Sie es über *mehrere* Sinneskanäle aufnehmen. Das erhöht die Chance, dass die Inhalte an verschiedenen Stellen in Ihrem Gehirn abgelegt werden.

Wir haben Konzepte und Bilder in **unerwarteter** Weise eingesetzt, weil Ihr Gehirn auf Neuigkeiten programmiert ist. Und wir haben Bilder und Ideen mit zumindest *etwas emotionalem Charakter* verwendet, weil Ihr Gehirn darauf eingestellt ist, auf die Biochemie von Gefühlen zu achten. An alles, was ein *Gefühl* in Ihnen auslöst, können Sie sich mit höherer Wahrscheinlichkeit erinnern, selbst wenn dieses Gefühl nicht mehr ist als ein bisschen **Belustigung, Überraschung oder Interesse**.

Wir haben einen **umgangssprachlichen** Stil mit direkter Anrede benutzt, denn Ihr Gehirn ist von Natur aus aufmerksamer, wenn es den Eindruck hat, Sie seien in einer Unterhaltung, als wenn es davon ausgeht, dass Sie passiv einer Präsentation zuhören – sogar dann, wenn Sie *lesen*.

Wir haben mehr als 80 **Aktivitäten** für Sie vorgesehen, denn Ihr Gehirn lernt und behält von Natur aus besser, wenn Sie Dinge **tun**, als wenn Sie nur darüber *lesen*. Und wir haben die Übungen zwar anspruchsvoll, aber doch lösbar gemacht, weil das für die meisten Leser den größten Ansporn bietet und das Lernen durch selbst erarbeitete Erfolgsergebnisse direkt und aktiv unterstützt.

Wir haben **mehrere unterschiedliche Lernstile** eingesetzt, denn vielleicht bevorzugen *Sie* ein Schritt-für-Schritt-Vorgehen, während jemand anders erst einmal den groben Zusammenhang verstehen und ein Dritter einfach nur Beispieldaten sehen möchte. Aber ganz abgesehen von den jeweiligen Lernvorlieben profitiert *jeder* davon, wenn er die gleichen Inhalte in unterschiedlicher Form präsentiert bekommt.

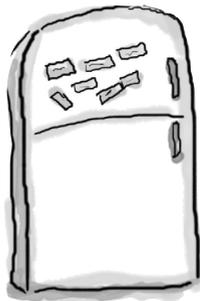
Wir liefern Inhalte für **beide Seiten Ihres Gehirns**, denn je mehr Sie von Ihrem Gehirn einsetzen, umso wahrscheinlicher werden Sie lernen und behalten, und umso länger bleiben Sie konzentriert. Wenn Sie mit einer Seite des Gehirns arbeiten, bedeutet das häufig, dass sich die andere Seite des Gehirns ausruhen kann; so können Sie über einen längeren Zeitraum produktiver lernen.

Und wir haben **Geschichten** und Übungen aufgenommen, die **mehr als einen Blickwinkel repräsentieren**, denn Ihr Gehirn lernt von Natur aus intensiver, wenn es gezwungen ist, selbst zu analysieren, zu bewerten und zu urteilen.

Wir haben **Herausforderungen** eingefügt: in Form von Übungen und indem wir **Fragen** stellen, auf die es nicht immer eine eindeutige Antwort gibt, denn Ihr Gehirn ist darauf eingestellt, zu lernen und sich zu erinnern, wenn es an etwas *arbeiten* muss. Überlegen Sie: Ihren Körper bekommen Sie ja auch nicht in Form, wenn Sie im Fitnessstudio *die Leute nur beobachten*. Aber wir haben unser Bestes getan, um dafür zu sorgen, dass Sie – wenn Sie schon hart arbeiten – an den *richtigen* Dingen arbeiten. Dass Sie **nicht eine einzige Synapse darauf verschwenden**, ein schwer verständliches Beispiel zu verarbeiten oder einen schwierigen, mit Fachbegriffen gespickten oder sonstwie überladenen Text zu analysieren.

Wir haben **Menschen** eingesetzt. In Geschichten, Beispielen, Bildern usw. – denn *Sie sind* ein Mensch. Und Ihr Gehirn schenkt *Menschen* mehr Aufmerksamkeit als *Dingen*.





Und das können SIE tun, um sich Ihr Gehirn gefügig zu machen

So, wir haben unseren Teil geleistet. Der Rest liegt bei Ihnen. Diese Tipps sind ein Anfang; hören Sie auf Ihr Gehirn und finden Sie heraus, was bei Ihnen funktioniert und was nicht. Probieren Sie neue Wege aus.

Schneiden Sie das hier aus und heften Sie es an Ihren Kühlschrank.

- 1 Immer langsam. Je mehr Sie verstehen, umso weniger müssen Sie auswendig lernen.**
Lesen Sie nicht nur. Halten Sie inne und denken Sie nach. Wenn das Buch Sie etwas fragt, springen Sie nicht einfach zur Antwort. Stellen Sie sich vor, dass Sie das wirklich jemand *fragt*. Je gründlicher Sie Ihr Gehirn zum Nachdenken zwingen, umso größer ist die Chance, dass Sie lernen und behalten.
- 2 Bearbeiten Sie die Übungen. Machen Sie sich eigene Notizen.**
Wir haben die Übungen entworfen, aber wenn wir sie auch für Sie lösen würden, wäre das, als würde ein anderer Ihr Training für Sie absolvieren. Und sehen Sie sich die Übungen *nicht einfach nur an – benutzen Sie Papier und Bleistift*. Es deutet vieles darauf hin, dass körperliche Aktivität beim Lernen den Lernerfolg erhöhen kann.
- 3 Lesen Sie die Abschnitte »Es gibt keine dummen Fragen«.**
Und zwar alle. Das sind keine Zusatzanmerkungen – *sie gehören zum Kerninhalt!* Überspringen Sie sie nicht.
- 4 Lesen Sie dies als Letztes vor dem Schlafengehen. Oder lesen Sie danach zumindest nichts Anspruchsvolles mehr.**
Ein Teil des Lernprozesses (vor allem die Übertragung ins Langzeitgedächtnis) findet erst statt, *nachdem* Sie das Buch zur Seite gelegt haben. Ihr Gehirn braucht für die weitere Verarbeitung Zeit für sich. Wenn Sie in dieser Zeit etwas Neues aufnehmen, geht ein Teil dessen, was Sie gerade gelernt haben, verloren.
- 5 Trinken Sie Wasser. Viel.**
Ihr Gehirn arbeitet am besten in einem schönen Flüssigkeitsbad. Dehydrierung, also Austrocknung (zu der es schon kommen kann, bevor Sie überhaupt Durst verspüren), beeinträchtigt die kognitiven Funktionen.
- 6 Reden Sie drüber. Laut.**
Sprechen aktiviert andere Hirnareale. Wenn Sie etwas verstehen wollen oder Ihre Chancen verbessern wollen, sich später daran zu erinnern, sprechen Sie es laut aus. Noch besser: Versuchen Sie, es jemand anderem laut zu erklären. Sie lernen dann schneller und haben vielleicht Ideen, auf die Sie beim bloßen leisen Lesen nie gekommen wären.
- 7 Hören Sie auf Ihr Gehirn.**
Achten Sie darauf, Ihr Gehirn nicht zu überlasten. Wenn Sie merken, dass Sie etwas nur noch überfliegen oder dass Sie sich das gerade erst Gelesene nicht merken können, ist es Zeit für eine Pause. Ab einem bestimmten Punkt lernen Sie nicht mehr schneller, indem Sie mehr hineinzustopfen versuchen; das kann den Lernprozess sogar beeinträchtigen.
- 8 Aber bitte mit Gefühl!**
Ihr Gehirn muss wissen, dass es *um etwas Wichtiges* geht. Lassen Sie sich in die Geschichten hineinziehen. Erfinden Sie eigene Bildunterschriften für die Fotos. Über einen schlechten Scherz zu stöhnen, ist *immer noch besser*, als gar nichts zu empfinden.
- 9 Krempeln Sie die Ärmel hoch und tun Sie was!**
Es gibt nur eine Möglichkeit, die Datenanalyse beherrschen zu lernen: die Ärmel hochkrempeln und sich die Hände schmutzig machen. Und genau das werden Sie in diesem Buch von Anfang bis Ende tun. Datenanalyse ist eine Kunst, und Übung ist die einzige Möglichkeit, gut darin zu werden. Dazu werden wir Ihnen eine Menge Gelegenheit geben: Jedes Kapitel enthält Übungen, in denen wir Ihnen Probleme zum Lösen anbieten. Überfliegen Sie sie nicht einfach – ein wesentlicher Teil des Lernprozesses findet statt, während Sie an den Übungen arbeiten. Zu jeder Übung finden Sie unsere Lösung – *trauen Sie sich, dort nachzusehen, wenn Sie stecken bleiben!* (Es passiert leicht, dass man sich bei irgendeiner Kleinigkeit festfährt.) Versuchen Sie aber immer erst selbst, die Übung zu lösen. Und sorgen Sie auf alle Fälle dafür, dass Ihre Lösung funktioniert, bevor Sie den nächsten Abschnitt in Angriff nehmen!

Lies mich!

Dies ist ein Lernerlebnis, kein Nachschlagewerk. Wir haben bewusst alles herausgestrichen, was an irgendeiner Stelle des Buchs hinderlich für den Lernprozess sein könnte. Und wenn Sie das Buch das erste Mal durcharbeiten, müssen Sie am Anfang beginnen, denn das Buch setzt voraus, dass Sie bestimmte Sachen schon gesehen und gelernt haben.

In diesem Buch geht es nicht um Software für die Datenanalyse.

Viele Bücher, in deren Titel »Datenanalyse« auftaucht, arbeiten lediglich die Liste aller Excel-Funktionen, die irgendetwas mit Datenanalyse zu tun haben, von oben nach unten ab und zeigen zu jeder ein paar Beispiele. In der *Datenanalyse von Kopf bis Fuß* beschreiben wir dagegen, *wie man als Datenanalyst arbeitet*. Sicher erfahren Sie in diesem Buch auch allerhand über Ihre Softwarewerkzeuge – allerdings nur als Mittel zum Zweck, um zu lernen, wie man gute Datenanalysen macht.

Wir gehen davon aus, dass Sie wissen, wie man mit einfachen Tabellenkalkulationsformeln arbeitet.

Haben Sie schon mal die SUMMENfunktion Ihrer Tabellenkalkulation benutzt? Falls nicht, sollten Sie sich vielleicht ein wenig intensiver in Ihre Tabellenkalkulation einarbeiten, bevor Sie dieses Buch in Angriff nehmen. Auch wenn in vielen Kapiteln nicht von Ihnen erwartet wird, eine Tabellenkalkulation zu benutzen, wird in einigen anderen vorausgesetzt, dass Sie wissen, wie man Formeln einrichtet. Sind Sie beispielsweise mit der SUMMENfunktion vertraut, bedeutet das für Sie einen wesentlichen Vorteil.

In diesem Buch geht es um mehr als nur Statistik.

Das Buch enthält einiges an Statistik, und als Datenanalyst sollten Sie so viel darüber lernen, wie Sie können. Sobald Sie also die *Datenanalyse von Kopf bis Fuß* durchgearbeitet haben, ist es in jedem Fall eine gute Idee, gleich noch *Statistik von Kopf bis Fuß* zu lesen. Allerdings umfasst das Arbeitsfeld »Datenanalyse« neben der Statistik eine Reihe anderer Bereiche, und die diversen nicht-statistischen Gebiete, die für dieses Buch ausgewählt wurden, konzentrieren sich auf die essenziellen, praktischen Aspekte von *Datenanalysen unter Realbedingungen*.

Die Übungen sind NICHT optional.

Die Übungen und sonstigen Aktivitäten sind keine Zugaben, sondern Grundbestandteil des Buchs. Einige davon helfen beim Einprägen, andere beim Verständnis und wieder andere bei der Anwendung des Gelernten. **Überspringen Sie die Übungen nicht.**

Zu den Kopfnuss-Übungen gibt es keine Lösungen.

Für manche gibt es keine richtige Lösung, bei anderen wiederum ist es ein Teil der Lernerfahrung, dass Sie selbst entscheiden, ob und wann Ihre Antworten richtig sind. Bei einigen Kopfnuss-Übungen finden Sie Hinweise, die Sie in die richtige Richtung lenken.

Die Redundanz ist beabsichtigt und wichtig.

Eine der Besonderheiten eines Buchs dieser Reihe ist: Wir wollen, dass Sie *wirklich* verstehen. Und wenn Sie mit dem Buch fertig sind, sollen Sie sich an das Gelernte erinnern. Bei den meisten Nachschlagewerken besteht das Ziel nicht im Behalten und Erinnern. Aber in *diesem* Buch geht es ums *Lernen*, und deshalb werden manche Ideen und Begriffe *mehr als ein Mal* besprochen.

Das Buch ist nach der letzten Seite noch nicht zu Ende.

Wir finden es gut, wenn Sie auf Begleitseiten zu einem Buch auf nützliches und unterhaltsames Extramaterial zugreifen können. Mehr zur Datenanalyse finden Sie unter den folgenden URLs:

http://examples.oreilly.de/german_examples/hfdataanalysisger/

<http://www.headfirstlabs.com/books/hfida/>.

Die Gutachter

Eric Heilman



Tony Rose



Bill Mietelski



Fachgutachter:

Eric Heilman ist Phi-Beta-Kappa-Absolvent der Walsh School of Foreign Service an der Georgetown University mit einem Abschluss in International Economics. Während seiner Studentenzeit in Washington D.C. hat er im Weißen Haus im State Department und im National Economic Council gearbeitet. Seine Abschlussarbeit in Wirtschaftswissenschaften legte er an der University of Chicago ab. Zurzeit unterrichtet er Statistik und Mathematik an der Georgetown Preparatory School in Bethesda, MD.

Bill Mietelski ist Softwareingenieur und dreimaliger Fachgutachter für Titel der *Von Kopf bis Fuß*-Reihe. Er kann es kaum erwarten, seine Golf-Statistik einer Datenanalyse zu unterziehen, was ihm hoffentlich dabei helfen wird, auf dem Platz abzuräumen.

Anthony Rose arbeitet seit annähernd zehn Jahren im Bereich Datenanalyse und ist zurzeit Geschäftsführer von Support Analytics, einem Consulting-Unternehmen für Datenanalyse und Datenvisualisierung. Anthony hat einen Abschluss als Betriebswirt mit Schwerpunkt Management und Finanzwesen, woher auch seine Leidenschaft für die Datenanalyse stammt. Wenn er nicht gerade arbeitet, kann man ihn auf dem Golfplatz in Columbia, Maryland, treffen, oder er ist in ein gutes Buch vertieft, lässt sich einen edlen Wein schmecken oder genießt einfach die Zeit, die er mit seinen Töchtern und seiner wundervollen Frau verbringt.

Danksagungen

Mein Lektor:

Brian Sawyer war als Lektor unglaublich. Mit Brian zu arbeiten, hat etwas davon, von einem professionellen Tänzer geführt zu werden. Permanent passiert eine Menge, ohne dass du es wirklich verstehst, aber du machst eine großartige Figur und hast ein fantastisches Erlebnis. Uns war eine aufregende Zusammenarbeit vergönnt, und seine Unterstützung, seine Rückmeldungen und seine Vorschläge waren unbezahlbar.

Das O'Reilly-Team:

Brett McLaughlin kannte die Vision dieses Projekts seit dessen Anfängen, begleitete es durch die schwierigen Phasen und war insgesamt eine kontinuierliche Unterstützung. Bretts unerbittliche Konzentration auf das, was **Sie** mit einem Buch der *Von Kopf bis Fuß*-Reihe erleben, war eine Inspiration. Er ist der Mann mit dem Durchblick.

Karen Shaner hat sich um die logistische Unterstützung des Projekts gekümmert und um diverse begeisterte Rückmeldungen an dem einen oder anderen kalten Morgen in Cambridge. Von **Brittany Smith** stammen ein paar der coolen Grafikelemente, die wir wieder und wieder benutzt haben.

Wirklich clevere Menschen, deren Einflüsse sich in diesem Buch wiederfinden:

Zwar sind viele der großen Konzepte, die in diesem Buch vermittelt werden, für das Thema Datenanalyse eher unkonventionell, nur wenige davon stammen aber im eigentlichen Sinn von mir. Ich habe mich mächtig aus den Werken der folgenden intellektuellen Vordenker (um nicht zu sagen, Superstars) bedient: Dietrich Dörner, Gerd Gigerenzer, Richards Heuer und Edward Tufte. Lesen Sie sie *alle*! Die Idee des Anti-Kompetenzprofils stammt aus Nassim Talebs »Der schwarze Schwan« (sollte es zu einem zweiten Band der Datenanalyse kommen, können Sie mit mehr seiner Ideen rechnen). **Richards Heuer** war freundlicherweise bereit, sich schriftlich mit mir über dieses Buch auszutauschen, er hat eine Reihe nützlicher Anregungen beigesteuert.

Freunde und Kollegen:

Für **Lou Barrs** intellektuelle, moralische, logistische und gestalterische Unterstützung bei diesem Buch möchte ich mich sehr bedanken. **Vezen Wu** hat mich in die Theorie relationaler Datenbanken eingeführt. **Aron Edidin** hat in meiner Studentenzeit einen umwerfend tollen Einführungskurs über die »Analyse nachrichtendienstlicher Informationen« gesponsort. Von meinen Poker-Freunden – **Paul, Brewster, Matt, Jon** und **Jason** – habe ich eine gründliche Ausbildung darin erhalten, wie man zwischen heuristischem und Optimierungsansatz einen Mittelweg wählt.

Leute, die für mich unersetzlich sind:

Die **Fachgutachter** haben exzellente Arbeit geleistet, haufenweise Fehler aufgespürt, eine Menge guter Vorschläge gemacht, und überhaupt waren sie sehr hilfreich.

Während ich an diesem Buch geschrieben habe, war mir mein Freund **Blair Christian** eine immense Hilfe. Blair ist Statistiker und ein sehr analytisch denkender Mensch. Seinen Einfluss finden Sie auf jeder einzelnen Seite. Danke für alles, Blair.

Meine Familie, **Michael Sr., Elizabeth, Sara, Gary** und **Marie** waren eine unglaubliche Hilfe. Und darüber hinaus bin ich dankbar für die unentwegte Unterstützung durch meine Frau **Julia**, die mir alles bedeutet. Danke euch allen!



Brian Sawyer



Brett McLaughlin



Blair und Niko Christian



Julia Burch

1 Einführung in die Datenanalyse

Wir zerlegen alles in seine Einzelteile



Überall sind Daten.

Heutzutage muss jeder mit Bergen von Daten fertig werden, ob er sich nun »Datenanalyst« nennt oder nicht. Diejenigen allerdings, die über Datenanalyse-Kompetenzen verfügen, haben einen **entscheidenden Vorsprung** vor allen anderen, weil sie wissen, was man mit all dem Zeug **machen** kann. Sie wissen, wie man aus Rohdaten Informationen gewinnt, mit denen sich **reale Prozesse steuern** lassen, und sie wissen, wie man komplexe Fragestellungen und Datenmengen so **aufschlüsselt und strukturiert**, dass man zum Kern der Probleme im jeweiligen Geschäftsfeld vordringt.

René Sans Kosmetik braucht Ihre Hilfe

Ihr erster Tag als Datenanalyst, und gerade haben Sie diese Umsatzzahlen zur Bewertung zugeschickt bekommen. Die Daten beschreiben die Verkaufszahlen von *Morgentau Plus*, der führenden Feuchtigkeitslotion von *René Sans Kosmetik*, Ihrem ersten Auftraggeber.

Was ist in den letzten sechs Monaten mit dem Umsatz passiert?

Wie ist das Verhältnis zwischen Bruttoumsatz und Absatzziel des Auftraggebers?

	September	Oktober	November	Dezember	Januar	Februar
Bruttoumsatz	5.280.000 €	5.501.000 €	5.469.000 €	5.480.000 €	5.533.000 €	5.554.000 €
Absatzziel	5.280.000 €	5.500.000 €	5.729.000 €	5.968.000 €	6.217.000 €	6.476.000 €
Werbekosten	1.056.000 €	950.400 €	739.200 €	528.000 €	316.800 €	316.800 €
Kosten für soziale Netzwerke	0 €	105.600 €	316.800 €	528.000 €	739.200 €	739.200 €
Grundpreis (100 ml)	2,00 €	2,00 €	2,00 €	1,90 €	1,90 €	1,90 €

Erkennen Sie ein Muster in der Kostenentwicklung bei René Sans Kosmetik?

Was läuft Ihrer Ansicht nach da ab bei den Grundpreisen? Warum sinken die?

Sehen Sie sich die Daten an. Es ist okay, wenn Sie nicht über alles Bescheid wissen – **nehmen Sie sich einfach Zeit** und sehen Sie sich die Daten genau an.

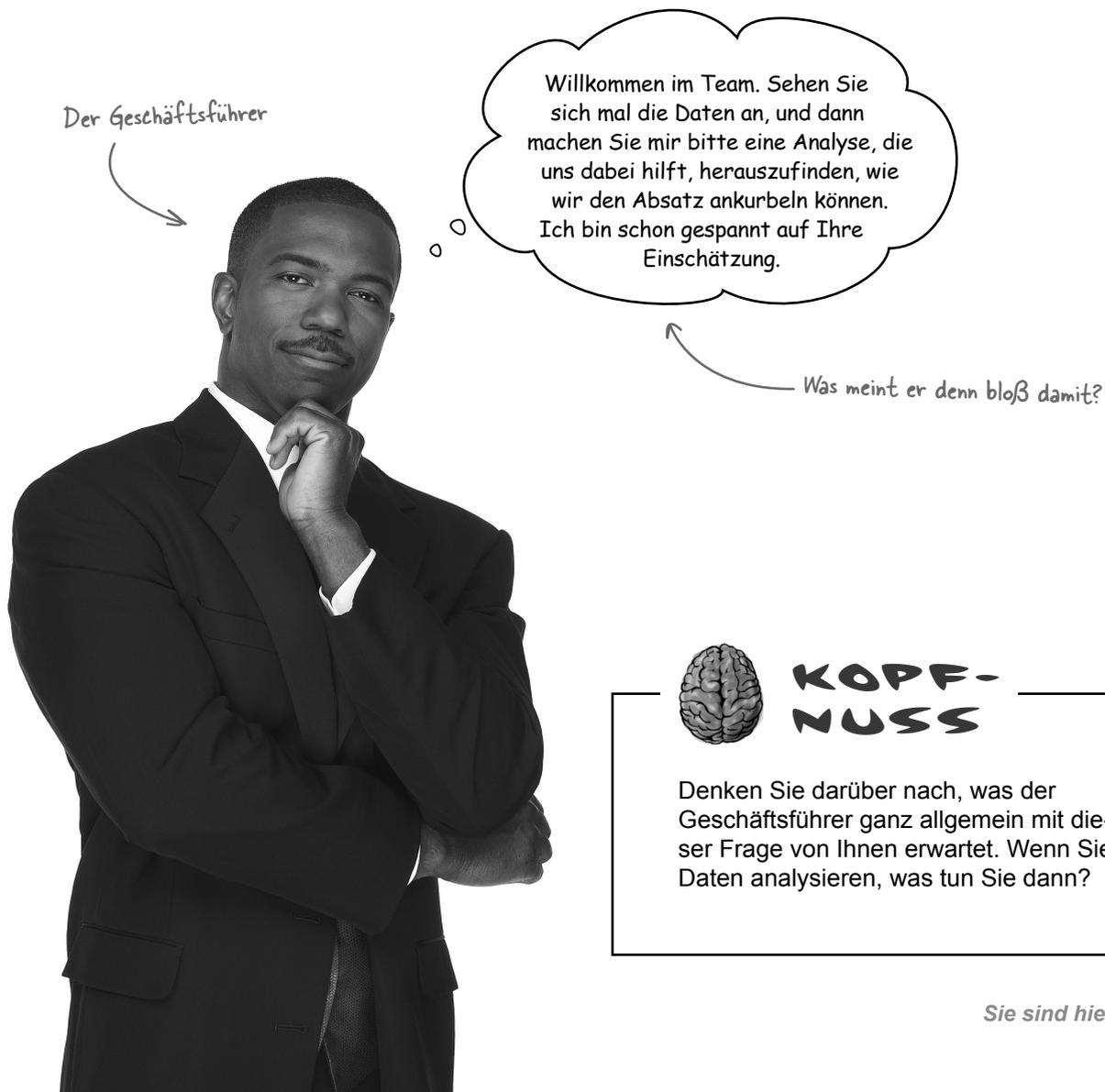
Was sehen Sie? Wie viel teilt Ihnen diese Tabelle über die Geschäftssituation von *René Sans Kosmetik* mit? Und über ihre Feuchtigkeitslotion *Morgentau Plus*?

Ein guter Datenanalyst will immer die Daten sehen.

Der Geschäftsführer würde den Absatz gern mit einer Datenanalyse anschieben

Er möchte, dass Sie ihm »eine Analyse machen«.

Ein irgendwie *vage* formulierter Auftrag, finden Sie nicht auch? Er klingt einfach, aber ist die Aufgabe wirklich so eindeutig? Klar, er will mehr Absatz. Sicher, er denkt sich, dass irgendwas in den Daten dabei helfen wird, dieses Ziel zu erreichen. Bloß was? Und ... wie?



KOPF- NUSS

Denken Sie darüber nach, was der Geschäftsführer ganz allgemein mit dieser Frage von Ihnen erwartet. Wenn Sie Daten analysieren, was tun Sie dann?

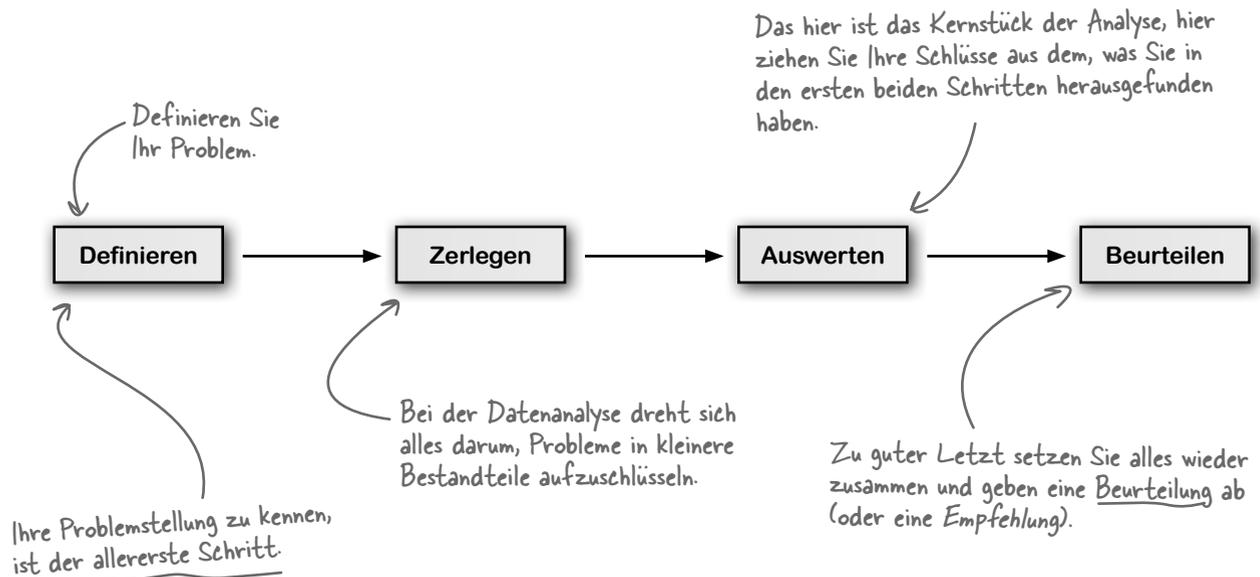
Datenanalyse heißt, sorgfältig über die Befundlage nachzudenken

Der Ausdruck *Datenanalyse* deckt eine Menge unterschiedlicher Tätigkeiten und verschiedene fachliche Kompetenzen (oder *Skills*) ab. Wenn Ihnen jemand erzählt, er sei Datenanalyst, sagt Ihnen das noch nicht viel darüber, was er nun *ganz genau* kann oder macht.

Jeder gute Datenanalyst, unabhängig von der jeweiligen Kompetenz und Ausrichtung, hält sich bei der Arbeit allerdings an den **immer gleichen grundlegenden Ablauf** und stützt sich beim Nachdenken über Problemstellungen immer sorgfältig auf empirische Belege.

Sie können sicher sein, dass diese Person Excel beherrscht – das war's dann aber auch schon, mehr wissen Sie vorab nicht!

In formaleren Zusammenhängen spricht man eher von einer *Fragestellung*, die Sie beantworten müssen – darin besteht Ihr Auftrag!



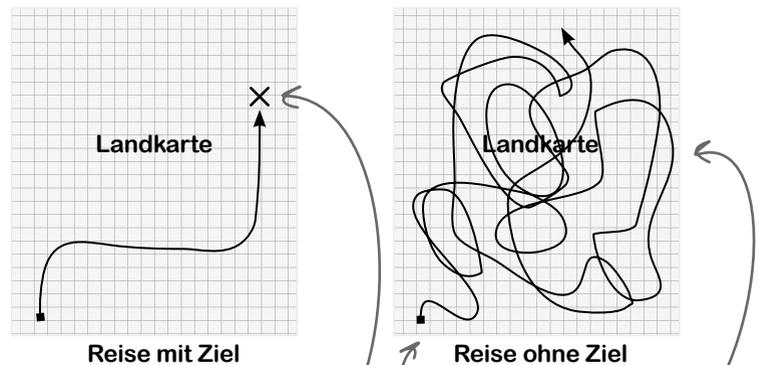
In jedem Kapitel dieses Buchs gehen Sie diese Schritte immer und immer wieder durch; Sie werden das sehr schnell verinnerlichen.

In letzter Konsequenz ist jede Datenanalyse dazu da, zu **besseren Urteilen** zu kommen, und Sie werden im Folgenden erfahren, wie man fundiertere Entscheidungen trifft, indem man in einem Meer von Daten nach Erkenntnissen sucht.

Definieren Sie das Problem

Eine Datenanalyse ohne **explizit** definierte Problemstellung oder Zielsetzung ist dasselbe, wie mit dem Auto zu einer Spritztour aufzubrechen, ohne sich für ein Ziel entschieden zu haben.

Sicher, Sie könnten auf ein paar interessante Orte stoßen, und manchmal *wollen* Sie einfach nur in der Hoffnung unterwegs sein, irgendetwas Aufregendes zu entdecken, **aber wer garantiert Ihnen, dass Sie auf diese Weise überhaupt etwas finden?**



Das hier ist ein gigantischer Analysebericht.

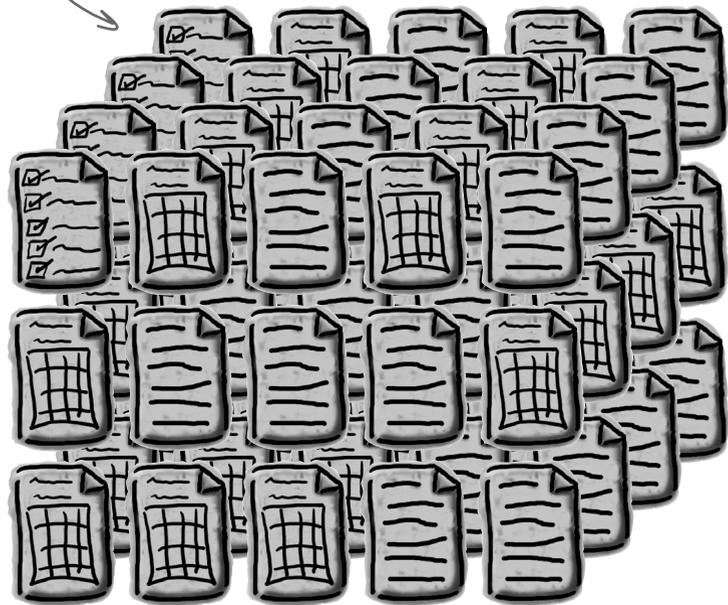
Sehen Sie die Gemeinsamkeiten?

Wer weiß schon, wo Sie das hinführt?

Haben Sie je einen Analysereport mit einer gefühlten Länge von **hunderttausend Seiten** gesehen, voll mit Unmengen an Flussgrafiken und Diagrammen?

Immer mal wieder braucht ein Datenanalyst einen Berg Papier oder eine stundenlange Präsentation, um zum Punkt zu kommen. In so einem Fall hat der Analyst aber oft sein Problem **nicht ausreichend eingegrenzt**. Er schüttet Sie mit Informationen zu und entzieht sich damit seiner Verpflichtung, eine **Empfehlung zur Lösung des Problems** auszusprechen.

Manchmal ist die Situation aber noch schlimmer: Das Problem selbst ist **alles andere als klar definiert**, und der Analyst möchte nicht, dass Sie merken, wie er einfach nur *die Datengrundlage* wiederkaut.



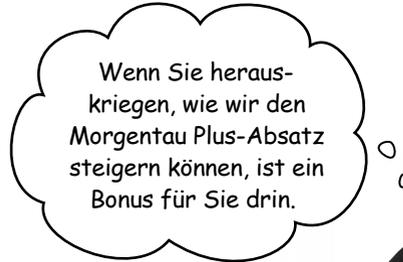
Wie definiert man das Problem?

Ihr Auftraggeber hilft Ihnen, das Problem zu definieren

Er ist derjenige, dem Ihre Analyse helfen soll. Ihr Auftraggeber könnte Ihr Vorgesetzter sein oder der Geschäftsführer des Unternehmens, für das Sie arbeiten, unter Umständen sind Sie es sogar selbst.

Ihr Auftraggeber ist *derjenige*, der auf der Grundlage Ihrer Analyse eine Entscheidung treffen wird. **Sie müssen so viele Informationen wie möglich von ihm bekommen**, um die Fragestellung definieren zu können.

In unserem jetzigen Fall möchte der Geschäftsführer von René Sans Kosmetik mehr Absatz. Das ist aber nur eine sehr weit gefasste Wunschvorstellung. Sie müssen *genauer* verstehen, was er eigentlich meint, um eine Analyse entwerfen zu können, die sein Problem löst.



Das ist Ihr Auftraggeber, der Mensch, für den Sie arbeiten.

Es ist ausgesprochen nützlich, Ihren Auftraggeber so gut wie möglich kennenzulernen.

Punkt für Punkt

Ihr Auftraggeber könnte:

- seine Daten gut oder schlecht kennen
- seine Fragestellung oder seine Zielsetzung gut oder schlecht kennen bzw. formulieren
- seine Branche gut oder schlecht kennen
- zielstrebig oder unschlüssig sein
- deutlich oder unklar sein
- intuitiv oder analytisch vorgehen

Der Geschäftsführer von René Sans Kosmetik.

Behalten Sie während dieses Kapitels die Orientierungshilfe am Seitenende im Auge, sie zeigt Ihnen, wo Sie gerade sind.

Je besser Sie Ihren Auftraggeber kennen, desto wahrscheinlicher ist es, dass ihn Ihre Analyse unterstützen kann.

Definieren

Zerlegen

Auswerten

Beurteilen

Es gibt keine Dummen Fragen

F: Ich mag es eigentlich, in meinen Daten herumzustöbern. Meinen Sie denn, dass ich immer ein klares Ziel im Kopf haben muss, bevor ich mir meine Daten überhaupt nur ansehe?

A: Nur um sich mal Ihre Daten anzusehen, müssen Sie keine Fragestellung im Kopf haben. Behalten Sie aber im Gedächtnis, dass *reines Datensichten* selbst noch keine Datenanalyse ist. Datenanalyse hat immer damit zu tun, Probleme zu identifizieren und zu beantworten.

F: Ich habe mal was von »exploratorischer Datenanalyse« gehört, wobei man Daten nach Ideen für weitergehende Analysen absucht. Bei dieser Art Datenanalyse gibt es aber keine Fragestellungen!

A: Natürlich gibt es die. Die Fragestellung besteht bei einer exploratorischen Datenanalyse (oder *Datenexploration*) darin, Hypothesen zu finden, die eine Überprüfung wert sind. Das ist eindeutig ein sehr konkretes Problem, das man da zu lösen hat.

F: Gut. Erzählen Sie mir mehr über Auftraggeber, die sich nicht ganz im Klaren darüber sind, was ihr Problem ist. Brauchen solche Leute überhaupt einen Datenanalysten?

A: Aber sicher!

F: Klingt für mich eher nach der Sorte, die professionelle Lebenshilfe braucht.

A: In der Tat hilft ein (auch sozial) kompetenter Datenanalyst seinem Auftraggeber dabei, sein Problem zu durchdenken; er sitzt nicht einfach da und wartet darauf, dass ihm der Auftraggeber sagt, was er tun soll. Ihre Kunden werden es sehr schätzen, wenn Sie ihnen zeigen können, dass sie Probleme haben, von denen sie noch nicht einmal wussten, dass sie sie haben.

F: Für mich klingt das albern. Wer will denn mehr Probleme?

A: Leute, die einen Datenanalysten anheuern, haben erkannt, dass jemand mit analytischen Kompetenzen in der Lage ist, ihren Geschäftserfolg zu steigern. Manche ergreifen Probleme als *Chance*, und Datenanalysten, die ihren Auftraggebern zeigen, wie man Spielräume nutzt, verschaffen ihnen einen Wettbewerbsvorteil.



Spitzen Sie Ihren Bleistift

Das Ausgangsproblem ist hier, dass wir den Absatz steigern sollen. Welche Fragen würden Sie dem Geschäftsführer stellen, um besser zu verstehen, was er sich genau vorstellt? Denken Sie sich mindestens fünf aus.

1

2

3

4

5

Feedback von René Sans für Sie

Diese E-Mail hier kam gerade als Antwort auf Ihre Fragen rein. Haufenweise Information ...

Hier haben wir ein paar Beispiel-fragen, um den Geschäftsführer dazu zu bringen, die Ziele für Ihre Analyse zu definieren.

Fragen Sie immer: »Wie viel?« Quantifizieren Sie Ihre Ziele und Ihre Einschätzungen.

Schätzen Sie ab, worüber sich Ihr Auftraggeber Gedanken macht. Er wird sich ganz sicher Sorgen wegen seiner Mitbewerber machen.

Sie bemerken etwas in den Zahlen, das Ihnen seltsam vorkommt? Fragen Sie nach!

Ihre Fragen könnten anders aussehen.

Von: Geschäftsleitung, René Sans Kosmetik
An: Datenanalyse von Kopf bis Fuß
Betreff: Re: Problemstellung definieren

Um wie viel wollen Sie den Absatz steigern?

Ich muss wieder mit unserem Absatzziel gleichziehen, das finden Sie in der Tabelle. Unsere gesamte Budgetierung baut darauf auf; wenn wir das Absatzziel verpassen, werden wir Schwierigkeiten bekommen.

Wie wollen wir das Ihrer Meinung nach angehen?

Nun, das herauszufinden ist Ihre Sache. Die Strategie wird aber wohl beinhalten, die Leute dazu zu bringen, mehr zu kaufen, und mit »Leute« meine ich weibliche Teenager (Alter 11–15 Jahre). Irgendwie werden Sie die Verkäufe doch wohl mit der einen oder anderen Marketingmaßnahme ankurbeln können. Sie sind derjenige für die Daten – finden Sie's selbst heraus!

Welcher Zuwachs in den Verkaufszahlen wäre Ihrer Ansicht nach angemessen? Sind die bestehenden Absatzzielzahlen realistisch?

Diese Tween-Mädel haben's reichlich. Geld fürs Babysitten, Eltern und so weiter. Ich glaube nicht, dass es irgendeine obere Grenze dafür gibt, wie viel Morgentau Plus wir an sie verkaufen können.

Wie sieht der Absatz bei Ihren Mitbewerbern aus?

Dazu habe ich keine harten Zahlen, mein Eindruck ist aber, dass die uns abhängen werden. Ich schätze, man ist uns in Bezug auf den Bruttoerlös bei Feuchtigkeitslotionen um 50 bis 100% voraus.

Was hat es mit den Budgets für Werbung und soziale Netzwerke auf sich?

Wir versuchen was Neues. Das Gesamtbudget (also für beides zusammen) beträgt 20% vom Erlös im ersten Monat. Ursprünglich ging das alles komplett in die traditionelle Werbung, aber wir sind dabei, es in die sozialen Netzwerke zu verlagern. Ich bekomme Gänsehaut, wenn ich daran denke, was passieren würde, wenn wir die herkömmliche Werbung auf dem gleichen Niveau hielten.

Definieren

Zerlegen

Auswerten

Beurteilen

Brechen Sie Problemstellung und Daten in besser überschaubare Teile auf

Der nächste Schritt in einer Datenanalyse besteht darin, alles, was Sie mithilfe Ihres Auftraggebers über das Problem in Erfahrung gebracht haben, sowie Ihre Daten zu nehmen und so zu zerlegen, dass genau das Maß an **Feinauflösung** entsteht, das Ihrer Analyse am meisten nützt.



Zerlegen Sie das Problem in kleinere Teilprobleme

Sie müssen Ihr Problem in **handliche, lösbare Blöcke** zerlegen. Häufig wird Ihre Fragestellung *diffus* sein, etwa so:

»Wie lässt sich unser Absatz steigern?«

- »Was erwarten unsere besten Kunden von uns?«
- »Welche Art von Werbung ist am aussichtsreichsten?«
- »Wie bewähren sich unsere Werbemaßnahmen?«

Die komplette Fragestellung lässt sich nicht direkt beantworten. Indem Sie aber die kleineren Teilfragen beantworten, die Sie durch die *Analyse des Gesamtproblems* formulieren konnten, finden Sie die Lösung für das Gesamtproblem.

Lösen Sie die Gesamtfragestellung, indem Sie die Teilfragen beantworten.

Teilen Sie die Daten in kleinere Blöcke auf

Gleiches Spiel mit den Daten. Man wird Ihnen die präzise quantifizierte Informationen, die Sie benötigen, nicht frei Haus liefern; Sie müssen die wesentlichen Bestandteile selbst isolieren.

Wenn die Daten, die Sie erhalten haben, als **Zusammenfassung oder Aufstellung** vorliegen, wie diejenigen, die Sie von René Sans Kosmetik bekommen haben, werden Sie wissen wollen, welche einzelnen Bestandteile am wichtigsten für Sie sind.

Liegen Ihnen dagegen **Rohdaten** vor, werden Sie die Einzeldaten zusammenfassen müssen, damit Sie etwas Verwertbares erhalten.

	September	Oktober	November	Dezember	Januar	Februar
Bruttumsatz	5.280.000 €	5.501.000 €	5.469.000 €	5.480.000 €	5.533.000 €	5.554.000 €
Absatzziel	5.280.000 €	5.500.000 €	5.729.000 €	5.968.000 €	6.217.000 €	6.476.000 €
Werbekosten	1.056.000 €	950.400 €	739.200 €	528.000 €	316.800 €	316.800 €
Kosten für soziale Netzwerke	0 €	105.600 €	316.800 €	528.000 €	739.200 €	739.200 €
Grundpreis (100 ml)	2,00 €	2,00 €	2,00 €	1,90 €	1,90 €	1,90 €

Bruttumsatz Dezember 5.480.000 €
 ... versus Grundpreis Dezember 1,90 €

Das hier könnten zwei Teilbereiche sein, die Sie im Auge behalten müssen.

In Kürze mehr zu diesen Schlagworten!

Versuchen wir es mal mit dem Zerlegen ...

Sehen Sie sich ein weiteres Mal an, was Sie haben

Beginnen wir mit den Daten. Sie haben hier eine Aufstellung der Umsatzzahlen von René Sans Kosmetik, und am Anfang versuchen Sie am besten, die wichtigsten Anteile dadurch zu isolieren, dass Sie nach Ansatzpunkten für **aussagekräftige Vergleiche** suchen.

Brechen Sie eine Datenaufstellung auf, indem Sie nach **interessanten Vergleichsmöglichkeiten** suchen.

In welchem Verhältnis stehen im November Brutto- und Zielumsatz?

Wie verhält sich der Bruttoumsatz vom Januar im Vergleich zu dem vom Februar?

	September	Oktober	November	Dezember	Januar	Februar
Bruttoumsatz	5.280.000 €	5.501.000 €	5.469.000 €	5.480.000 €	5.533.000 €	5.554.000 €
Absatzziel	5.280.000 €	5.500.000 €	5.729.000 €	5.968.000 €	6.217.000 €	6.476.000 €
Werbekosten	1.056.000 €	950.400 €	739.200 €	528.000 €	316.800 €	316.800 €
Kosten für soziale Netzwerke	0 €	105.600 €	316.800 €	528.000 €	739.200 €	739.200 €
Grundpreis (100 ml)	2,00 €	2,00 €	2,00 €	1,90 €	1,90 €	1,90 €

Wie ändern sich über die Zeit die Kosten für Werbung und soziale Netzwerke im Verhältnis zueinander?

Geht die Reduzierung des Grundpreises in irgendeiner Weise mit Änderungen im Bruttoumsatz einher?

Clevere Vergleiche sind essenzieller Bestandteil jeder Datenanalyse, Sie werden das im gesamten Buch immer wieder machen.

Im vorliegenden Fall möchten Sie sich einen **klaren Eindruck** davon verschaffen, wie das Geschäft mit Morgentau Plus läuft, indem Sie die Kennwerte der Aufstellung miteinander vergleichen.



Sie haben die *Problemstellung* definiert: **herausfinden, wie der Absatz gesteigert werden kann**. Aber deren Formulierung sagt Ihnen wenig darüber, wie man sich Ihre *Problemlösung* vorstellt. Sie mussten also allerhand nützliche Anmerkungen aus dem Geschäftsführer selbst herauskitzeln.

Seine Kommentare enthalten eine **Reihe wichtiger Grundüberzeugungen** darüber, wie die Kosmetikbranche funktioniert. Es ist zu hoffen, dass der Geschäftsführer mit seinen Einschätzungen recht hat, sie werden **das Rückgrat** Ihrer Analyse bilden! Welche sind die wichtigsten Argumente des Geschäftsführers?

Dieser Kommentar selbst birgt eine Art Datensatz in sich. Welche Teile sind am wichtigsten?

Was bringt am meisten?

Hier haben wir eine der »Wie«-Fragen.

Von: Geschäftsleitung, René Sans Kosmetik
 An: Datenanalyse von Kopf bis Fuß
 Betreff: Re: Problemstellung definieren

Um wie viel wollen Sie den Absatz steigern?

Ich muss wieder mit unserem Absatzziel gleichziehen, das finden Sie in der Tabelle. Unsere gesamte Budgetierung baut darauf auf; wenn wir das Absatzziel verpassen, werden wir Schwierigkeiten bekommen.

Wie wollen wir das Ihrer Meinung nach angehen?

Nun, das herauszufinden ist Ihre Sache. Die Strategie wird aber wohl beinhalten, die Leute dazu zu bringen, mehr zu kaufen, und mit »Leute« meine ich weibliche Teenager (Alter 11–15 Jahre). Irgendwie werden Sie die Verkäufe doch wohl mit der einen oder anderen Marketingmaßnahme ankurbeln können. Sie sind derjenige für die Daten – finden Sie's selbst heraus!

Welcher Zuwachs in den Verkaufszahlen wäre Ihrer Ansicht nach angemessen? Sind die bestehenden Absatzzielzahlen realistisch?

Diese Tween-Mädels haben's reichlich. Geld fürs Babysitten, Eltern und so weiter. Ich glaube nicht, dass es irgendeine obere Grenze dafür gibt, wie viel Morgentau Plus wir an sie verkaufen können.

Wie sieht der Absatz bei Ihren Mitbewerbern aus?

Dazu habe ich keine harten Zahlen, mein Eindruck ist aber, dass die uns abhängen werden. Ich schätze, man ist uns in Bezug auf den Bruttoerlös bei Feuchtigkeitslotionen um 50 bis 100% voraus.

Was hat es mit den Budgets für Werbung und soziale Netzwerke auf sich?

Wir versuchen was Neues. Das Gesamtbudget (also für beides zusammen) beträgt 20% vom Erlös im ersten Monat. Ursprünglich ging das alles komplett in die traditionelle Werbung, aber wir sind dabei, es in die sozialen Netzwerke zu verlagern. Ich bekomme Gänsehaut, wenn ich daran denke, was passieren würde, wenn wir die herkömmliche Werbung auf dem gleichen Niveau hielten.

Spitzen Sie Ihren Bleistift



Fassen Sie die Annahmen Ihres Auftraggebers und Ihre eigenen Überlegungen zu den Daten, die Sie für die Analyse bekommen haben, zusammen. **Zerlegen** Sie sowohl die E-Mail oben als auch die Daten in handlichere Teile, die die vorliegende Situation beschreiben.

Die Einschätzung Ihres Auftraggebers:

Ihre Gedanken zu den Daten:

- 1
- 2
- 3
- 4

- 1
- 2
- 3
- 4



Spitzen Sie Ihren Bleistift Lösung

Sie haben gerade eine Bestandsaufnahme Ihrer Einschätzungen und derjenigen Ihres Auftraggebers zur Situation gemacht. Was haben Sie herausgefunden?

Die Einschätzung Ihres Auftraggebers:

Ihre eigenen Antworten könnten sich etwas von diesen hier unterscheiden.

- 1 Morgentau Plus-Käufer sind Tweens (also weibliche Teenager zwischen 11 und 15 Jahren). Im Großen und Ganzen ist das die einzige Käufergruppe.
- 2 René Sans Kosmetik versucht, Kosten für Werbung auf soziale Netzwerke umzuverteilen, aber der Erfolg dieser Bemühungen ist zurzeit unklar.
- 3 Eine Grenze für die potenziellen Zuwächse der Verkaufszahlen bei den Tweens wird nicht gesehen.
- 4 Die Mitbewerber von René Sans Kosmetik sind extrem gefährlich.

Gut ... so etwas gehört heutzutage dazu.

Es könnte wichtig sein, sich das zu merken.

Ihre Gedanken zu den Daten:

Großes Problem!

- 1 Die Verkaufszahlen sind weit von den Absatzzielen entfernt, seit November driften sie auseinander.
- 2 Verglichen mit Januar ist der Absatz im Februar geringfügig höher, wirkt aber irgendwie flau.
- 3 Den Werbeetat zu kürzen, könnte René Sans Kosmetik daran gehindert haben, ihre Ziele zu erreichen.
- 4 Preissenkungen scheinen nicht dabei geholfen zu haben, die Verkaufszahlen mit den Absatzzielen Schritt halten zu lassen.

Was sollten sie als Nächstes tun?

Sie haben Ihre Problemstellung in kleinere, besser handhabbare Teilstücke zerlegt.

Jetzt sind wir so weit, diese Stücke detaillierter auszuwerten...



Werten Sie die Teilprobleme aus

Jetzt kommt der Teil, der richtig Spaß macht. Sie wissen, wonach Sie *suchen* müssen, und Sie wissen, welche Teile der Daten Sie *hinführen* werden. Sehen Sie sich jetzt diese Teilmformationen sehr genau und konzentriert an und machen Sie sich Ihr *eigenes* Bild.



Genau wie beim Zerlegen liegt der Schlüssel beim Auswerten im **Vergleichen** der isolierten Teilstücke.

Was fällt Ihnen auf, wenn Sie diese Einzelinformationen miteinander vergleichen?

Wählen Sie je zwei Informationen aus und lesen Sie sie zusammen durch.

Was fällt Ihnen auf?

Beobachtungen zum Problem:

Morgentau Plus-Käufer sind Tweens* (also weibliche Teenager zwischen 11 und 15 Jahren). Im Großen und Ganzen ist das die einzige Käufergruppe.

René Sans Kosmetik versucht, Kosten für Werbung auf soziale Netzwerke umzuverteilen, aber der Erfolg dieser Bemühungen ist zurzeit unklar.

Eine Grenze für die potenziellen Zuwächse der Verkaufszahlen bei den Tweens wird nicht gesehen.

Die Mitbewerber von René Sans Kosmetik sind extrem gefährlich.

Nutzen Sie Ihr Vorstellungsvermögen!

Beobachtungen in den Daten:

Die Verkaufszahlen sind weit von den Absatzzielen entfernt.

Verglichen mit Januar ist der Absatz im Februar geringfügig höher, wirkt aber irgendwie flau.

Den Werbeetat zu kürzen, könnte René Sans Kosmetik daran gehindert haben, ihre Ziele zu erreichen.

Preissenkungen scheinen nicht dabei geholfen zu haben, die Verkaufszahlen mit den Absatzzielen Schritt halten zu lassen.

Sie haben fast alle nötigen Bestandteile zusammen, aber ein wichtiges Stück fehlt noch ...

*Ach ja, »Tweens« – das ist Marketing-Jargon und steht für »Pre-Teens« oder »Be-tweens« – »Too old for toys, too young for boys«. Aber schlagen Sie bitte nicht mich, ich bin nur der Übersetzer...