

---

# Angewandte Statistik mit R



---

Reiner Hellbrück

# Angewandte Statistik mit R

Eine Einführung für Ökonomen und  
Sozialwissenschaftler

3. Auflage



Springer Gabler

Reiner Hellbrück  
Fakultät Wirtschaftswissenschaften  
Hochschule für angewandte Wissenschaften  
Würzburg-Schweinfurt, Deutschland

ISBN 978-3-658-12861-6      ISBN 978-3-658-12862-3 (eBook)  
DOI 10.1007/978-3-658-12862-3

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

Springer Gabler

© Springer Fachmedien Wiesbaden 2009, 2011, 2016

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

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Der Verlag, die Autoren und die Herausgeber gehen davon aus, dass die Angaben und Informationen in diesem Werk zum Zeitpunkt der Veröffentlichung vollständig und korrekt sind. Weder der Verlag noch die Autoren oder die Herausgeber übernehmen, ausdrücklich oder implizit, Gewähr für den Inhalt des Werkes, etwaige Fehler oder Äußerungen.

Gedruckt auf säurefreiem und chlorfrei gebleichtem Papier

Springer Gabler ist Teil von Springer Nature

Die eingetragene Gesellschaft ist Springer Fachmedien Wiesbaden GmbH

# Vorwort zur dritten Auflage

Die Einleitung und Kapitel 2 wurden überarbeitet; bei letzterem wurde der Abschnitt über LimeSurvey und Moodle entfernt. Wegen seiner Bedeutung im Wirtschaftsleben ist die Berechnung des geometrischen Mittels aufgenommen worden. Da das Paket 'QRMLib' für R, Version 3.2.3, nicht verfügbar ist, wurde in Kapitel 3 eine entsprechende Anpassung vorgenommen.

In Kapitel 4 wurde der Abschnitt 'Einfache und zusammengesetzte Hypothesen' eingefügt. Die einseitigen Tests wurden auf den Standardfall, dass das Gleichheitszeichen in der Nullhypothese steht, umgestellt. Leicht geändert wurde die Gliederung des Kapitels 6 und zusätzliche Erläuterungen in Abschnitt 6.5 wurden eingefügt. Erläuterungen zur Kovarianz wurden in Kapitel 7 vorgenommen. Die Einleitung des Kapitels 8 beginnt nun mit einem Praxisbeispiel, in der Einleitung des Kapitels 12 wurde ebenfalls die praktische Bedeutung der Regressionsanalyse etwas hervorgehoben.

In Anhang C wird nun zusätzlich erläutert, wie in einzelnen Paketen vorhandene Daten eingelesen werden können und wie man Updates einzelner Pakete erhalten kann. Um Wechselkursdaten herunterzuladen gibt es Oanda, hierzu wird nun auch eine Alternative, der Service 'TrueFX', vorgestellt. Die Befehle dieses Kapitels wurden zudem auf Funktion geprüft.

In der jüngsten Zeit hat Herr Florian Schubert, M.Sc., die Veranstaltung in parallelen Gruppen zu mir gehalten. Diskussionen, die sich hierbei mit ihm ergaben, sind in die neue Auflage eingeflossen. Vielen Dank, das machte Laune. Und auch Frau Hasenbalg, Lektorin beim Verlag Springer Gabler, möchte ich gerne für die gute Zusammenarbeit danken. Alle verbliebenen Fehler gehen, wie immer, zu meinen Lasten.

Würzburg, im März 2016

Reiner Hellbrück

# Vorwort zur ersten Auflage

Dieses Buch entstand im Zuge der Neustrukturierung meiner Statistikveranstaltungen an der FH Würzburg-Schweinfurt. Die fortschreitende Digitalisierung macht auch vor der Statistik nicht halt und so entstand der Wunsch, die Veranstaltungen neu auszurichten. Wie an Fachhochschulen üblich, liegt der Schwerpunkt auf der Anwendung. Aus diesem Grund sind im allgemeinen nach einer kurzen Darstellung des nötigen Hintergrundwissens Beispiele angefügt. Hierbei kommt die Statistiksoftware R zum Einsatz.

R wird sehr selektiv eingesetzt, allen Anwendungen ist zuvor ein Grundlagenkapitel vorgeschaltet, in dem Maßzahlen, Teststatistiken, Hypothesen und dergleichen vorgestellt werden. Leser, die einen schnellen Überblick über die Software wünschen, seien auf den Anhang C verwiesen, wo die wichtigsten Befehle dargestellt werden. Um Mißverständnissen vorzubeugen, sei ausdrücklich darauf hingewiesen, dass es sich hier um keine Einführung in das Programmpaket R handelt: Methoden und Anwendung, gestützt mit Software, stehen gleichberechtigt nebeneinander. Das Erlernen der Software ergibt sich als nützlicher Nebeneffekt.

Die anfänglichen Rechnungen erfolgten mit Version R-2.5, dann mit neueren. Die Software ist auf vielen verschiedenen Betriebssystemen lauffähig. Im vorliegenden Fall wurde Windows XP und Suse-Linux genutzt. Bei Linux wurden teilweise Rechnungen mit Hilfe einer Shell, (auch Konsole oder Befehlsfenster genannt) größtenteils aber mit Emacs-ess durchgeführt. Emacs ist ein Text-Editor, der üblicherweise mit jeder Linux-Distribution ausgeliefert wird. Das Kürzel 'ess' steht für 'emacs speaks statistics' und will heißen, dass das Zusatzwerkzeug 'Emacs-ess' als Benutzeroberfläche (als 'frontend') für Statistiksoftware eingesetzt werden kann. Hierüber ist es möglich, mit einer einheitlichen Benutzeroberfläche verschiedene Statistikprogramme, darunter auch 'SPSS' und 'Stata', anzusprechen. Eigene Versuche in dieser Richtung wurden von dem Autor bislang nicht unternommen.

Der Einstieg ist sehr einfach gehalten, um dem Studenten während der ersten Wochen genügend Zeit zu lassen, die neue Software auf seinem eigenen Rechner zu installieren und kennenzulernen. Erfahrungsgemäß stellen sich bereits bei dem Einlesen der Daten die ersten Probleme ein. Dies rührt aus der Verwendung unterschiedlicher Parameter, die zur Trennung von Zeichen bei Textdateien verwendet werden. Desweiteren gibt es üblicherweise Probleme durch die Verwendung unterschiedlicher Betriebssysteme. R ist primär für Linuxsysteme geschrieben. Hier gelten jedoch etwas andere Konventionen bei der Angabe von Pfaden: statt des '\', wie in Windowssystemen üblich, wird das Zeichen '/' verwendet.

Werden die Befehle nicht direkt in dem Befehlsfenster von R geschrieben, sondern in einem Textverarbeitungsprogramm, so kann es nach Kopieren der Befehle in das Befehlsfenster leicht zu Fehlermeldungen kommen. Ursache ist dann häufig die automatische Ersetzung der Anführungszeichen in typographische Anführungszeichen innerhalb des Textverarbeitungsprogramms. Deshalb

wird empfohlen, zum Schreiben oder Bearbeiten von Befehlen eine Software zu verwenden, die solche automatischen Ersetzungen nicht vornimmt, oder dass solche Funktionen ausgeschaltet werden.

Der deskriptiven Statistik ist vergleichsweise wenig Raum gewidmet, der Schwerpunkt liegt auf der schließenden Statistik und multivariaten Verfahren, bei denen seitens Ökonomen (speziell meiner Kolleginnen und Kollegen) Nachfrage besteht. Der Text kann, je nach Belieben, unterschiedlich verwendet werden. Einerseits besteht die Möglichkeit, die Theorie weitestgehend in den Hintergrund zu drängen, um sich ausschließlich auf die Anwendung zu konzentrieren: die Kapitel 6 und 7 zur Wahrscheinlichkeitstheorie können dann übersprungen werden. Dies bietet sich an, wenn eine Veranstaltung zur Wahrscheinlichkeitstheorie vorgeschaltet ist.

Andererseits ist es möglich, Inhalte anwendungsnahe zu präsentieren, und bei Bedarf nötiges Wissen in Wahrscheinlichkeitstheorie einzuflechten. Dann bietet es sich an, die Kapitel in der angegebenen Folge zu besprechen. Da die Kapitel 10 und 11, ohne statistische Tests auskommen, können sie auch zur Veranschaulichung multivariater Verfahren vorgezogen werden.

Bei einigen Lehrbüchern hat sich zwischenzeitlich die Unart eingeschlichen, während des laufenden Textes nicht zu zitieren. Es scheint, als habe ein sehr bekanntes Lehrbuch der Mikroökonomie, diese Entwicklung eingeleitet. Dem Autor des Lehrbuches verbrannte das Manuskript mitsamt der Zitate. Aus den verbliebenen Resten wurde es fast gänzlich ohne Zitate fertiggestellt. Hierdurch wird dem Studenten der Eindruck vermittelt, als brauche man nicht zu zitieren. Diesem Zeitgeist wird hier nicht gefolgt. Es wird angegeben, woher der Autor seine Weisheiten hat.

Dank schulde ich vielen, insbesondere meinem akademischen Lehrer Prof. Dr. Volker Steinmetz, der es außerordentlich gut verstand, theoretische Statistik und Ökonometrie zu vermitteln. Herr Prof. Dr. Rudolf Richter bot bereits in den 80-er Jahren PC-gestützte ökonometrische Auswertungen an, damals ein Novum. Beide Ansätze werden hier miteinander verknüpft. Danken möchte ich an dieser Stelle auch meinem wissenschaftlichen Mitarbeiter Manuel Hertel, für die gute Zusammenarbeit und die Entlastung durch seine Übungsstunden, die er mit großer Umsicht anbietet. Schließlich möchte ich bei meinen Söhnen, David und Simon um Nachsicht bitten, für die Zeit, die ich in meinem Arbeitszimmer den PC blockiert habe. Meine Frau genöß die Zeit, während ich 'aufgeräumt' war, ebenso wie ich.

Würzburg, im Juni 2009: Reiner Hellbrück

# Vorwort zur zweiten Auflage

Kleinere Veränderungen sind vorgenommen worden. So wurde die Bedeutung der Messbarkeit besser herausgearbeitet. In Kapitel 3 ist die logarithmische Skala hinzugefügt worden und bei der Regression wird die Thematik der Kointegration angesprochen. In Kapitel C wurde eine weitere Möglichkeit zur Installation zusätzlicher R -Pakete in Unixsystemen eingefügt. Zudem wurden einige Internetadressen und Befehle aktualisiert sowie Schreibfehler der 1. Auflage korrigiert. Ein herzliches Dankeschön geht an zwei meiner Studenten, Herrn Daniel Back und Herrn Felix Kreß, die mich freundlicherweise auf Druckfehler hingewiesen haben. Alle verbliebenen Fehler gehen selbstverständlich zu meinen Lasten.

Freundlichst wird darauf hingewiesen, dass die verwendeten Daten von der Homepage des Verlages heruntergeladen werden können. Hierbei handelt es sich um eine \*.zip-Datei. Die enthaltenen Dateien müssen zuerst entpackt werden, damit R darauf zugreifen kann. Insbesondere für Dozenten finden sich zusätzliche Materialien; so werden beispielsweise alle Abbildungen zur Verfügung gestellt.

Gerne bin ich bereit, eine Befragung mit LimeSurvey zu ermöglichen. Das Programm ist zwar kostenfrei, doch seine Installation auf einem Server, die Nutzung und die Einrichtung von Nutzungsrechten verursachen Kosten. Aktuelle Konditionen erhalten Sie auf Anfrage. Senden Sie bei Interesse eine E-Mail an

[reiner.hellbrueck@fhws.de](mailto:reiner.hellbrueck@fhws.de).

Um alle Funktionen des Programms verfügbar zu haben, ist es notwendig, zumindest R 2.11 zu installieren. Zudem kann es notwendig sein, Pakete zu aktualisieren; ansonsten kann es zu Fehlermeldungen kommen. So ist die logarithmische Skalierung beispielsweise in älteren Distributionen nicht enthalten.

Herr Christian Schuld hat mich freundlicherweise bei der Beschaffung von Literatur unterstützt. Text, Layout, Stichwortverzeichnis, Glossar und Literaturverzeichnis wurden, wie an Fachhochschulen meist der Fall, selbst erstellt.  $\LaTeX$  hat hier wertvolle Dienste geleistet.

Würzburg, im Oktober 2010: Reiner Hellbrück



# Inhaltsverzeichnis

Vorwort zur dritten Auflage	v
Vorwort zur ersten Auflage	vi
Vorwort zur zweiten Auflage	viii
Abbildungsverzeichnis	xvii
Tabellenverzeichnis	xix
<b>1 Einleitung</b>	<b>1</b>
1.1 Gegenstand . . . . .	1
1.2 Aufbau . . . . .	5
<b>2 Datenerhebung - ganz praktisch</b>	<b>9</b>
2.1 Einleitung . . . . .	9
2.2 Statistikpaket R . . . . .	10
2.3 Erhebungsplan . . . . .	11
2.3.1 Grundlagen . . . . .	11
2.3.2 Beispiel . . . . .	12
2.4 Ziehen einer Stichprobe . . . . .	13
2.4.1 Grundlagen . . . . .	13
2.4.2 Beispiel . . . . .	14
2.5 Rohdaten auslesen . . . . .	14
2.5.1 Grundlagen . . . . .	14
2.5.2 Beispiel . . . . .	15
2.6 Daten in Statistikprogramm einlesen . . . . .	15
2.6.1 Grundlagen . . . . .	15
2.6.2 Beispiel . . . . .	17
2.7 Plausibilitätsprüfung . . . . .	19
2.7.1 Grundlagen . . . . .	19
2.7.2 Beispiel 1 . . . . .	19
2.7.3 Einfache Datensätze . . . . .	20
2.7.4 Beispiel 2 . . . . .	22
2.7.5 Komplexe Datensätze . . . . .	24

2.7.6	Beispiel 3 . . . . .	26
2.8	Abschließende Bemerkungen . . . . .	27
2.9	Kontrollfragen . . . . .	28
2.10	Aufgaben . . . . .	29
<b>3</b>	<b>Datenaufbereitung</b>	<b>31</b>
3.1	Einleitung . . . . .	31
3.2	Graphische Methoden . . . . .	32
3.2.1	Grundlagen . . . . .	32
3.2.2	Beispiele . . . . .	34
3.3	Absolute Häufigkeitsverteilung . . . . .	38
3.3.1	Grundlagen . . . . .	38
3.3.2	Beispiel 1 . . . . .	39
3.3.3	Maßzahlen . . . . .	39
3.3.4	Beispiel 2 . . . . .	40
3.4	Relative Häufigkeitsverteilung . . . . .	42
3.4.1	Grundlagen . . . . .	42
3.4.2	Beispiel 1 . . . . .	42
3.4.3	Maßzahlen . . . . .	44
3.4.4	Beispiel 2 . . . . .	45
3.5	Verteilungsfunktion und Quantile . . . . .	45
3.5.1	Verteilungsfunktion . . . . .	45
3.5.2	Quantile . . . . .	46
3.5.3	Verteilungsfunktion und Quantile . . . . .	49
3.6	Histogramme . . . . .	50
3.6.1	Absolute Häufigkeit . . . . .	50
3.6.2	Durchschnittliche Häufigkeitsdichte . . . . .	52
3.7	Kontingenztafel . . . . .	54
3.7.1	Gemeinsame Verteilung . . . . .	54
3.7.2	Randverteilungen . . . . .	55
3.7.3	Bedingte Verteilung und statistische Unabhängigkeit . . . . .	57
3.8	Lorenz-Kurve . . . . .	57
3.8.1	Grundlagen . . . . .	57
3.8.2	Beispiel . . . . .	58
3.8.3	Gini-Koeffizienten . . . . .	61
3.9	Abschließende Bemerkungen . . . . .	63
3.10	Kontrollfragen . . . . .	64
3.11	Aufgaben . . . . .	65
3.A	Nützliches zu Maßzahlen* . . . . .	68
3.B	Logarithmische Skala* . . . . .	68
<b>4</b>	<b>Statistisches Testen</b>	<b>71</b>
4.1	Einleitung . . . . .	71
4.2	Binomialverteilung . . . . .	72
4.2.1	Grundlagen . . . . .	72
4.2.2	Beispiel . . . . .	74

4.3	Test . . . . .	75
4.3.1	Zweiseitige Fragestellung . . . . .	75
4.3.2	Einseitige Fragestellung - Version 1 . . . . .	78
4.3.3	Einseitige Fragestellung - Version 2 . . . . .	80
4.3.4	Fehler 1. Art . . . . .	81
4.3.5	Beispiel . . . . .	81
4.4	Einfache und zusammengesetzte Hypothesen* . . . . .	84
4.4.1	Einfache Hypothesen . . . . .	84
4.4.2	Zusammengesetzte Hypothesen . . . . .	84
4.5	Abschließende Bemerkungen . . . . .	88
4.6	Kontrollfragen . . . . .	90
4.7	Aufgaben . . . . .	91
4.A	Wirkungsanalyse* . . . . .	93
4.A.1	Grundlagen . . . . .	93
4.A.2	Test . . . . .	95
4.A.3	Beispiel . . . . .	96
4.A.4	Abschließende Bemerkungen . . . . .	97
<b>5</b>	<b>Chi-Quadrat Tests</b>	<b>99</b>
5.1	Einleitung . . . . .	99
5.2	Unabhängigkeitstest . . . . .	100
5.2.1	Grundlagen . . . . .	100
5.2.2	Beispiel . . . . .	102
5.3	Anpassungstest . . . . .	105
5.3.1	Grundlagen . . . . .	105
5.3.2	Beispiel . . . . .	106
5.4	Homogenitätstest . . . . .	107
5.4.1	Grundlagen . . . . .	107
5.4.2	Beispiel . . . . .	109
5.5	Abschließende Bemerkungen . . . . .	111
5.6	Kontrollfragen . . . . .	111
5.7	Aufgaben . . . . .	112
<b>6</b>	<b>Wahrscheinlichkeitsräume</b>	<b>115</b>
6.1	Einleitung . . . . .	115
6.2	Definitionsmenge . . . . .	116
6.3	Wahrscheinlichkeitsraum der Grundgesamtheit . . . . .	118
6.3.1	Begriff . . . . .	118
6.3.2	Laplacescher Wahrscheinlichkeitsraum . . . . .	119
6.4	Wahrscheinlichkeitsraum der Stichprobe . . . . .	121
6.4.1	Begriff . . . . .	121
6.4.2	Grundgesamtheit und Stichprobe . . . . .	123
6.5	Wichtige Zusammenhänge und Begriffe . . . . .	124
6.5.1	Rechenregeln . . . . .	124
6.5.2	Bedingte Wahrscheinlichkeit . . . . .	125
6.5.3	Stochastische Unabhängigkeit . . . . .	126

6.5.4	Multiplikationssatz . . . . .	126
6.5.5	Satz von der totalen Wahrscheinlichkeit . . . . .	127
6.5.6	Satz von Bayes . . . . .	127
6.5.7	Diskreter Wahrscheinlichkeitsraum . . . . .	129
6.6	Abschließende Bemerkungen . . . . .	130
6.7	Kontrollfragen . . . . .	131
6.8	Aufgaben . . . . .	131
<b>7</b>	<b>Abbildungen von Ergebnisräumen</b>	<b>135</b>
7.1	Einleitung . . . . .	135
7.2	Messbarkeit und Zufallsvariable . . . . .	136
7.2.1	Messbarkeit . . . . .	136
7.2.2	Zufallsvariablen . . . . .	137
7.3	Verteilungsfunktion und Dichte . . . . .	138
7.3.1	Verteilungsfunktion . . . . .	138
7.3.2	Dichte . . . . .	140
7.4	Maßzahlen . . . . .	141
7.4.1	Erwartungswert . . . . .	141
7.4.2	Kovarianz, Varianz und Standardabweichung . . . . .	142
7.4.3	Standardisierung . . . . .	143
7.5	Abschließende Bemerkungen . . . . .	144
7.6	Kontrollfragen . . . . .	144
7.7	Aufgaben . . . . .	145
<b>8</b>	<b>Einfache Korrelationsanalyse</b>	<b>149</b>
8.1	Einleitung . . . . .	149
8.2	Korrelation . . . . .	151
8.2.1	Wahrscheinlichkeitstheorie . . . . .	151
8.2.2	Empirische Korrelation (Bravais-Pearson) . . . . .	151
8.2.3	Berechnung bei Wertepaaren . . . . .	152
8.2.4	Beispiele . . . . .	153
8.3	Tests bei kardinalen Merkmalen . . . . .	155
8.3.1	Stetige normalverteilte Zufallsvariablen . . . . .	155
8.3.2	Stetige nicht-normalverteilte Zufallsvariablen . . . . .	160
8.4	Test bei ordinalen Merkmalen: Bell-Doksum Test . . . . .	166
8.4.1	Test . . . . .	166
8.4.2	Beispiel . . . . .	167
8.5	Abschließende Bemerkungen . . . . .	171
8.6	Kontrollfragen . . . . .	172
8.7	Aufgaben . . . . .	172
8.A	Weitere Tests* . . . . .	174

<b>9</b>	<b>Multivariate Korrelationsanalyse*</b>	<b>177</b>
9.1	Einleitung . . . . .	177
9.2	Vergleich zweier Korrelationen . . . . .	178
9.2.1	Grundlagen . . . . .	178
9.2.2	Beispiel . . . . .	179
9.3	Partielle Korrelation . . . . .	180
9.3.1	Grundlagen . . . . .	180
9.3.2	Beispiel 1 . . . . .	181
9.3.3	Test . . . . .	181
9.3.4	Beispiel 2 . . . . .	182
9.4	Zusammenhang zwischen mehreren Merkmalen . . . . .	182
9.4.1	Grundlagen . . . . .	182
9.4.2	Beispiel . . . . .	184
9.5	Globaltest . . . . .	185
9.5.1	Test . . . . .	185
9.5.2	Beispiel . . . . .	185
9.6	Multiple Vergleiche . . . . .	186
9.6.1	Test . . . . .	186
9.6.2	Beispiel . . . . .	188
9.7	Multiple Korrelation . . . . .	191
9.7.1	Grundlagen . . . . .	191
9.7.2	Beispiel 1 . . . . .	191
9.7.3	Test . . . . .	192
9.7.4	Beispiel 2 . . . . .	193
9.8	Kanonische Korrelation . . . . .	194
9.8.1	Grundlagen . . . . .	194
9.8.2	Beispiel 1 . . . . .	195
9.8.3	Test . . . . .	196
9.8.4	Beispiel 2 . . . . .	197
9.9	Abschließende Bemerkungen . . . . .	198
9.10	Kontrollfragen . . . . .	199
9.11	Aufgaben . . . . .	200
<b>10</b>	<b>Daten- und Distanzmatrix</b>	<b>201</b>
10.1	Einleitung . . . . .	201
10.2	Distanzmatrizen . . . . .	203
10.2.1	Definition und Eigenschaften . . . . .	203
10.2.2	Skalierung . . . . .	204
10.3	Kardinale Merkmale . . . . .	204
10.3.1	Intervall- und Verhältnisskala . . . . .	204
10.3.2	Manhattan-Distanz . . . . .	206
10.4	Ordinale Merkmale . . . . .	210
10.4.1	Grundlagen . . . . .	210
10.4.2	Beispiel . . . . .	211
10.5	Nominale Merkmale . . . . .	214
10.5.1	Grundlagen . . . . .	214

10.5.2	Beispiel . . . . .	214
10.6	Binäre Merkmale . . . . .	215
10.6.1	Grundlagen . . . . .	215
10.6.2	Beispiel . . . . .	217
10.7	Abschließende Bemerkungen . . . . .	218
10.8	Kontrollfragen . . . . .	219
10.9	Aufgaben . . . . .	220
<b>11</b>	<b>Clusteranalyse</b>	<b>223</b>
11.1	Einleitung . . . . .	223
11.2	Klassifikation . . . . .	226
11.2.1	Klassifikationstypen . . . . .	226
11.2.2	Konstruktionsverfahren . . . . .	227
11.3	PAM . . . . .	228
11.3.1	Grundlagen . . . . .	228
11.3.2	Beispiel 1 . . . . .	228
11.3.3	Bestimmung der Medoiden* . . . . .	232
11.3.4	Beispiel 2 . . . . .	236
11.3.5	Isolierte Cluster . . . . .	236
11.3.6	Beispiel 3 . . . . .	237
11.3.7	Überprüfung der Klassenbildung . . . . .	239
11.3.8	Beispiel 4 . . . . .	240
11.3.9	Bestimmung der Klassenzahl . . . . .	241
11.3.10	Beispiel 5 . . . . .	241
11.4	FANNY . . . . .	241
11.4.1	Grundlagen . . . . .	241
11.4.2	Beispiel 1 . . . . .	243
11.4.3	Partition und Überdeckung . . . . .	244
11.4.4	Beispiel 2 . . . . .	245
11.4.5	Überprüfung der Klassenbildung und Klassenanzahl . . . . .	248
11.4.6	Beispiel 3 . . . . .	249
11.5	MONA . . . . .	249
11.5.1	Grundlagen . . . . .	249
11.5.2	Beispiel 1 . . . . .	250
11.5.3	Assoziationsmaß . . . . .	252
11.5.4	Beispiel 2 . . . . .	253
11.5.5	Missings . . . . .	256
11.5.6	Beispiel 3 . . . . .	256
11.6	Abschließende Bemerkungen . . . . .	257
11.7	Kontrollfragen . . . . .	258
11.8	Aufgaben . . . . .	259

<b>12 Einfache Regression</b>	<b>261</b>
12.1 Einleitung . . . . .	261
12.2 Einfaches klassisches Regressionsmodell . . . . .	262
12.2.1 Grundlagen . . . . .	262
12.2.2 Beispiel . . . . .	265
12.3 Regressionsfunktion . . . . .	267
12.3.1 Grundlagen . . . . .	267
12.3.2 Beispiel . . . . .	268
12.4 Prognose . . . . .	270
12.4.1 Grundlagen . . . . .	270
12.4.2 Beispiel 1 . . . . .	270
12.4.3 Problem . . . . .	271
12.4.4 Beispiel 2 . . . . .	271
12.5 Bestimmtheitsmaß . . . . .	273
12.5.1 Grundlagen . . . . .	273
12.5.2 Beispiel . . . . .	275
12.6 Vollständiges Modell . . . . .	277
12.7 Tests . . . . .	278
12.7.1 Grundlagen . . . . .	278
12.7.2 Beispiel . . . . .	280
12.8 Abschließende Bemerkungen . . . . .	282
12.9 Kontrollfragen . . . . .	283
12.10 Aufgaben . . . . .	284
12.A Beweis der Streuungserlegungsformel*	286
12.B Erwartungswerte der KQ-Koeffizienten*	287
12.C Standardisierung* . . . . .	288
12.C.1 Erwartungswert . . . . .	288
12.C.2 Varianz . . . . .	288
12.D Partielle Korrelation* . . . . .	290
<b>A Theoretische Verteilungen</b>	<b>293</b>
A.1 Einleitung . . . . .	293
A.2 Diskrete Verteilungen . . . . .	294
A.2.1 Gleichverteilung* . . . . .	294
A.2.2 Bernoulli- und Binomialverteilung . . . . .	295
A.2.3 Hypergeometrische Verteilung* . . . . .	296
A.2.4 Poisson-Verteilung* . . . . .	298
A.2.5 Geometrische Verteilung* . . . . .	300
A.3 Stetige Verteilungen . . . . .	301
A.3.1 Rechteckverteilung . . . . .	301
A.3.2 Exponentialverteilung* . . . . .	303
A.3.3 Normalverteilung . . . . .	306
A.3.4 Chi-Quadrat-Verteilung . . . . .	310
A.3.5 t-Verteilung . . . . .	310
A.3.6 F-Verteilung . . . . .	311

<b>B</b>	<b>Matrizenrechnung</b>	<b>315</b>
B.1	Einleitung . . . . .	315
B.2	Matrizen . . . . .	316
B.2.1	Definition . . . . .	316
B.2.2	Vektoren . . . . .	316
B.2.3	Typen . . . . .	317
B.3	Verknüpfungen . . . . .	318
B.3.1	Gleichheitsrelation . . . . .	318
B.3.2	Addition . . . . .	319
B.3.3	Skalare Multiplikation . . . . .	320
B.3.4	Produkt zweier Matrizen . . . . .	322
B.3.5	Multiplikation von Vektoren . . . . .	324
B.4	Unabhängigkeit, Rang, Determinante, Inverse . . . . .	325
B.4.1	Lineare Unabhängigkeit . . . . .	325
B.4.2	Rang . . . . .	325
B.4.3	Determinante . . . . .	327
B.4.4	Inverse . . . . .	328
B.5	Eigenwerte, Eigenvektoren und Spur . . . . .	331
B.5.1	Definitionen . . . . .	331
B.5.2	Rechenregel . . . . .	331
B.5.3	Beispiele . . . . .	332
<b>C</b>	<b>Befehle in R</b>	<b>333</b>
C.1	Einleitung . . . . .	333
C.2	Grundlagen . . . . .	334
C.3	Daten einlesen, Objekte speichern und laden . . . . .	337
C.4	Dateneigenschaften . . . . .	340
C.5	Manipulation eingelesener Datensätze . . . . .	341
C.6	Graphik . . . . .	343
C.7	Suchen und Finden . . . . .	344
C.8	Besonderheiten in Windows . . . . .	346
C.9	Fehlermeldungen . . . . .	347
	<b>Anmerkungen und Lösungen</b>	<b>349</b>
	<b>Glossar</b>	<b>361</b>
	<b>Literaturverzeichnis</b>	<b>365</b>
	<b>Stichwortverzeichnis</b>	<b>367</b>



# Abbildungsverzeichnis

2.1	Rohdaten in Tabellenkalkulationsprogramm einlesen . . . . .	16
2.2	Anwendung empirische versus korrigierte Varianz . . . . .	22
3.1	Einfaches Liniendiagramm . . . . .	35
3.2	Liniendiagramm bei komplexen Datensätzen . . . . .	37
3.3	Kreisdiagramm . . . . .	38
3.4	Absolute Häufigkeitsverteilung . . . . .	41
3.5	Balkendiagramm . . . . .	42
3.6	Relative Häufigkeitsverteilung . . . . .	44
3.7	Empirische Verteilungsfunktion . . . . .	46
3.8	Berechnung der Quantile mit Option Typ 7 . . . . .	49
3.9	Histogramm mit absoluten Häufigkeiten . . . . .	51
3.10	Histogramm mit durchschnittlicher Häufigkeitsdichte . . . . .	53
3.11	Lorenzkurve . . . . .	60
3.12	Lorenzkurve: Konzentration auf ein Merkmal . . . . .	62
3.13	Umsatzentwicklung bei arithmetischer Skalierung . . . . .	69
3.14	Umsatzentwicklung bei halblogarithmischer Skalierung . . . . .	70
4.1	Binomialverteilung . . . . .	76
4.2	Hypothesentest: zweiseitige Fragestellung . . . . .	78
4.3	Hypothesentest: einseitige Fragestellung - Version 1 . . . . .	79
4.4	Hypothesentest: einseitige Fragestellung - Version 2 . . . . .	80
4.5	Verteilungsfunktionen der Binomialverteilung . . . . .	85
5.1	Annahme und Verwerfungsbereich . . . . .	102
6.1	Venn-Diagramme . . . . .	125
6.2	Veranschaulichung des Satzes von der totalen Wahrscheinlichkeit . . . . .	127
6.3	Baumdiagramm . . . . .	128
8.1	Streudiagramme (= Scatterplots) . . . . .	154
8.2	Veranschaulichung des Tests auf Korrelation . . . . .	159
8.3	Fishers z-Transformation . . . . .	175
10.1	Illustration der Manhattan-Distanz . . . . .	207

11.1	Verfahren . . . . .	225
11.2	Silhouette des 'output3' . . . . .	238
11.3	Silhouette des 'output8' . . . . .	240
11.4	Clusterbildung mit MONA . . . . .	252
12.1	Einkommen in Abhängigkeit des Alters . . . . .	266
12.2	KQ-Schätzung einer Cobb-Douglas Produktionsfunktion . . . . .	269
12.3	Translationsinvarianz des Bestimmtheitsmaßes . . . . .	276
A.1	Hypergeometrische Verteilung . . . . .	297
A.2	Poisson-Verteilung . . . . .	299
A.3	Verteilungsfunktion der Poisson-Verteilung . . . . .	300
A.4	Geometrische Verteilung . . . . .	301
A.5	Verteilungsfunktion der Geometrischen-Verteilung . . . . .	302
A.6	Rechteckverteilung . . . . .	303
A.7	Verteilungsfunktion der Rechteckverteilung . . . . .	304
A.8	Exponentialverteilung . . . . .	305
A.9	Verteilungsfunktion der Exponentialverteilung . . . . .	305
A.10	Standardnormalverteilung . . . . .	308
A.11	Verteilungsfunktion der Standardnormalverteilung . . . . .	308
A.12	Dichtefunktion der Chi-Quadrat-Verteilung . . . . .	309
A.13	Verteilungsfunktion der Chi-Quadrat-Verteilung . . . . .	309
A.14	Dichtefunktion der t-Verteilung . . . . .	312
A.15	Verteilungsfunktion der t-Verteilung . . . . .	312
A.16	Dichtefunktion der F-Verteilung . . . . .	314
A.17	Verteilungsfunktion der F-Verteilung . . . . .	314

# Tabellenverzeichnis

2.1	Daten YX . . . . .	17
2.2	Daten Einkommen Alter Ausbildungsjahre . . . . .	19
2.3	Daten2 . . . . .	24
3.1	Arbeitslose in Deutschland . . . . .	33
3.2	Vier mal drei Kontingenztabelle . . . . .	54
3.3	Randverteilung . . . . .	55
3.4	1. Schritt zur Erstellung einer Lorenz-Kurve . . . . .	58
3.5	2. Schritt zur Erstellung einer Lorenz-Kurve . . . . .	59
3.6	Umsatzentwicklung, Quelle: Daten frei erfunden . . . . .	68
5.1	Illustration zur Berechnung theoretischer Häufigkeiten . . . . .	101
5.2	Rohdaten . . . . .	103
5.3	Kontingenztabelle mit absoluten Häufigkeiten . . . . .	103
5.4	Eingabe x . . . . .	104
5.5	Eingabe y . . . . .	105
5.6	Kontingenztabelle mit bedingter Verteilung . . . . .	107
8.1	Beispiel: Umsatz - Bruttowertschöpfung . . . . .	157
8.2	Ränge . . . . .	165
9.1	Umsatz und Entfernung . . . . .	179
10.1	Kontingenztabelle bei binären Merkmalen . . . . .	216
11.1	Datenmatrix zur Bildung von zwei Partitionen . . . . .	228
11.2	Binäre Datenmatrix . . . . .	250