

Thomas Cleff

Deskriptive Statistik und Explorative Datenanalyse

Eine computergestützte Einführung
mit Excel, SPSS und STATA

3. Auflage



Springer Gabler

Deskriptive Statistik und Explorative Datenanalyse

Thomas Cleff

Deskriptive Statistik und Explorative Datenanalyse

Eine computergestützte Einführung mit
Excel, SPSS und STATA

3., überarbeitete und erweiterte Auflage

 Springer

Thomas Cleff
Pforzheim, Deutschland

ISBN 978-3-8349-4747-5
DOI 10.1007/978-3-8349-4748-2

ISBN 978-3-8349-4748-2 (eBook)

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

Gabler Verlag

Die 1. und 2. Auflage erschienen unter dem Titel „Deskriptive Statistik und moderne Datenanalyse“.

© Springer Fachmedien Wiesbaden 2008, 2011, 2015

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 Fachmedien Wiesbaden GmbH ist Teil der Fachverlagsgruppe Springer Science+Business Media (www.springer.com)

Vorwort zur dritten überarbeiteten und ergänzten Auflage

Mit großer Freude habe ich zur Kenntnis genommen, dass sich die zweite Auflage des Lehrbuches *Deskriptive Statistik und moderne Datenanalyse* einer so großen Nachfrage erfreut hat, dass sie beim Verlag nunmehr vergriffen ist. Freundlicherweise hat sich der Springer Gabler Verlag zur Ausgabe einer dritten – überarbeiteten und erweiterten – Auflage bereit erklärt, wofür ich mich bei der verantwortlichen Lektorin Frau Irene Buttkus herzlich bedanke. Ermöglicht es mir doch, das Buch um interessante Themenfelder der Explorativen Statistik zu erweitern. Neben den durch Software-Updates nötig gewordenen Anpassungen habe ich dem Buch zwei einführende Kapitel der Faktorenanalyse und der Clusteranalyse angefügt. Ich hoffe, ich kann somit nicht nur ein abgerundetes Programm der Deskriptiven Statistik unterbreiten, sondern dem Leser auch erste Einblicke in multivariate Verfahren der strukturentdeckenden (explorativen) Statistik ermöglichen. Entsprechend wurde der Titel des vorliegenden Lehrbuches in *Deskriptive Statistik und Explorative Datenanalyse* umbenannt.

Das Lehrbuch möchte den Studierenden der Volks- und Betriebswirtschaftslehre sowie Praktikern in Unternehmen die Grundlagen, Techniken und Anwendungsmöglichkeiten der Deskriptiven Statistik und der explorativen Datenanalyse näher bringen. Es geht zum einen auf die deskriptiven Basismethoden der univariaten und bivariaten Verfahren ein. Die Inhalte reichen von der Erhebung und Skalierung, über die univariate Analyse quantitativer Daten, bis zur Analyse bivariater Zusammenhänge. Zudem wird dem Leser ein erster Einblick in multivariate Verfahren wie der multivariaten Regression, der Faktorenanalyse und der Clusteranalyse ermöglicht. Alle Themen werden mit Hilfe von computerbasierten Berechnungen auf betriebswirtschaftliche Beispiele angewendet. Die Themengebiete decken alle wichtigen Aspekte einer Hochschulveranstaltung zur Deskriptiven Statistik ab bzw. gehen in Teilen sogar darüber hinaus.

Bei der Abfassung des Buches war ich stets bemüht, auch demjenigen einen Einblick in die Denkweise deskriptiver statistischer Verfahren zu ermöglichen, der ansonsten Schwierigkeiten mit der formalen oder methodischen Herangehensweise eines traditionellen Statistikbuches hat. An vielen Stellen habe ich versucht, auf überflüssige Formeln zu verzichten oder zunächst eine intuitive Herangehensweise an ein Thema zu wählen, bevor eine Formel abgeleitet bzw. angegeben wird. Es dürfte dennoch jeder verstehen, dass ein Buch über Statistik und Datenanalyse niemals ohne Formeln auskommen kann und es

auch nicht sollte. Da wo die Alltagssprache in ihrer Präzision versagt, ist und bleibt eine Formel letztlich die präziseste Form der sprachlichen Formulierung dessen, was methodisch ausgedrückt werden soll. Zur Vertiefung habe ich jedem Kapitel Übungsaufgaben nebst Lösungen angefügt, die ein effizientes Selbststudium erleichtern sollen.

Letztlich ermöglicht vor allem die allgemeine Verfügbarkeit von Computerprogrammen eine neue didaktische Herangehensweise an die Statistik. Jeder Studierende hat heute Zugriff auf Standardprogramme wie Excel oder auf Statistikpakete wie SPSS oder Stata. Dieses Lehrbuch beschränkt sich deshalb nicht nur auf die Darstellung der statistischen Verfahren, sondern erweitert den Blick auf dessen Anwendung mit Hilfe der Computerprogramme Excel 2010, SPSS (Version 22) und Stata (Version 13). Hierfür sind auf der Homepage des Verlages – neben anderen Zusatzmaterialien – Datensätze zur Verfügung gestellt. Mit ihnen können die Beispiel- und Übungsaufgaben durchgerechnet werden. Die Datensätze und die allgemeinen Zusatzmaterialien zu diesem Lehrbuch sind auf www.springer-gabler.de in der Rubrik „Zusätzliche Informationen“ zu finden. In derselben Rubrik befinden sich auch zusätzliche Materialien für Dozenten.

Ich danke an dieser Stelle allen Fachkollegen für die kritische Durchsicht des Manuskripts und für ihre wertvollen Hinweise. Verbleibende Fehler und Unzulänglichkeiten gehen selbstverständlich weiterhin zu meinen Lasten. Abschließend wäre dieses Buch niemals ohne die Unterstützung meiner Familie möglich gewesen. Ihr gilt mein ganz besonderer Dank.

Ich hoffe auch in Zukunft auf Anregungen und Verbesserungsvorschläge (z. B. an meine Emailadresse thomas.cleff@hs-pforzheim.de), denn gemäß einer chinesischen Weisheit sind nur mit den Augen der anderen die eigenen Fehler gut zu sehen.

Pforzheim, im Januar 2015

Thomas Cleff

Inhaltsverzeichnis

1	Statistik und empirische Forschung	1
1.1	Statistik lügt?	1
1.2	Zwei Arten von Statistik	3
1.3	Statistik als Erkenntnisprozess	5
1.4	Phasen empirischer Forschung	7
1.4.1	Von der Erkundung zur Theorie	7
1.4.2	Von der Theorie zum Modell	8
1.4.3	Vom Modell zur „Business Intelligence“	12
	Literatur	14
2	Vom Zahlenwust zum Datensatz	15
2.1	Möglichkeiten der Datenbeschaffung	15
2.2	Die Entscheidung für ein Skalenniveau	18
2.3	Datenerfassung mit dem Computer: Skalierung und Kodierung	22
2.4	Fehlende Werte oder Missing Values	23
2.5	Ausreißer und offensichtlich falsche Werte	26
2.6	Übungsaufgaben zum Abschnitt	27
	Literatur	28
3	Vom Datensatz zur Information	29
3.1	Erste Auswertungsschritte und grafische Darstellungen	29
3.2	Lageparameter als Informationsreduktion	36
3.2.1	Modus oder Modalwert	37
3.2.2	Der Mittelwert	37
3.2.3	Geometrisches Mittel	42
3.2.4	Harmonisches Mittel	44
3.2.5	Der Median	47
3.2.6	Quartile und Quantile	50
3.3	Boxplot – Erster Einblick in die Verteilung	51
3.4	Streuungsparameter	54
3.4.1	Die Standardabweichung und die Varianz	55

3.4.2	Der Variationskoeffizient	57
3.5	Schiefe und Kurtosis	59
3.6	Robustheit von Parametern	62
3.7	Konzentrationsmaße	63
3.8	Berechnung univariater Parameter mit dem Computer	66
3.8.1	Berechnung univariater Parameter mit SPSS	66
3.8.2	Berechnung univariater Parameter mit Stata	67
3.8.3	Berechnung univariater Parameter mit Excel 2010	68
3.9	Übungsaufgaben zum Abschnitt	69
	Literatur	72
4	Bivariate Zusammenhänge	73
4.1	Bivariate Skalenniveau-Kombinationen	73
4.2	Zusammenhang zweier nominaler Variablen	74
4.2.1	Kontingenztafeln	74
4.2.2	Die Chi-Quadrat Berechnung	75
4.2.3	Der Phi-Koeffizient	80
4.2.4	Der Kontingenzkoeffizient	83
4.2.5	Cramers V	84
4.2.6	Nominale Zusammenhänge mit SPSS	85
4.2.7	Nominale Zusammenhänge mit Stata	89
4.2.8	Nominale Zusammenhänge mit Excel	90
4.2.9	Übungsaufgaben zum Abschnitt	91
4.3	Zusammenhang zweier metrischer Variablen	94
4.3.1	Das Streudiagramm	94
4.3.2	Der Korrelationskoeffizient nach Bravais-Pearson	98
4.4	Zusammenhang ordinalskalierten Variablen	101
4.4.1	Die Rangkorrelation nach Spearman (Rho)	102
4.4.2	Kendalls Tau (τ)	109
4.5	Zusammenhangsmaße zweier Variablen mit unterschiedlichem Skalenniveau	115
4.5.1	Zusammenhang nominaler und metrischer Variablen	115
4.5.2	Zusammenhang nominaler und ordinaler Variablen	117
4.5.3	Zusammenhang ordinaler und metrischer Variablen	118
4.6	Korrelationsrechnung mit dem Computer	119
4.6.1	Korrelationsrechnung mit SPSS	120
4.6.2	Korrelationsrechnung mit Stata	121
4.6.3	Korrelationsrechnung mit Excel	122
4.7	Scheinkorrelationen	123
4.7.1	Partielle Korrelation	126
4.7.2	Partielle Korrelation mit SPSS	128
4.7.3	Partielle Korrelation mit Stata	128

4.7.4	Partielle Korrelation mit Excel	129
4.8	Übungsaufgaben zum Abschnitt	130
	Literatur	133
5	Regressionsanalyse	135
5.1	Erste Schritte einer Regressionsanalyse	135
5.2	Koeffizienten der bivariaten Regression	138
5.3	Multivariate Regressionskoeffizienten	143
5.4	Die Anpassungsgüte der Regression	144
5.5	Regressionsrechnung mit dem Computer	146
5.5.1	Regressionsrechnung mit Excel	146
5.5.2	Regressionsrechnung mit SPSS und Stata	148
5.6	Anpassungsgüte multivariater Regressionen	149
5.7	Regression mit unabhängiger Dummy-Variable	150
5.8	Hebelwirkungen von Beobachtungen	153
5.9	Nichtlineare Regressionen	154
5.10	Ansätze einer Regressionsdiagnostik	158
5.11	Übungsaufgaben zum Abschnitt	164
	Literatur	170
6	Zeitreihen- und Indexrechnung	171
6.1	Preisindizes	172
6.2	Mengenindizes	180
6.3	Wertindizes (Umsatzindizes)	182
6.4	Deflationierung von Zeitreihen	182
6.5	Umbasierung und Verkettung von Indizes	184
6.6	Übungsaufgaben zum Abschnitt	185
	Literatur	187
7	Clusteranalyse	189
7.1	Hierarchische Clusteranalyse	190
7.2	Die Clusterzentrenanalyse	206
7.3	Clusteranalyse mit dem Computer	208
7.3.1	Clusteranalyse mit SPSS	209
7.3.2	Clusteranalyse mit Stata	209
7.4	Übungsaufgaben zur Clusteranalyse	211
	Literatur	214
8	Faktorenanalyse	217
8.1	Faktorenanalyse: Grundlagen, Vorgehensweise und Interpretation	217
8.2	Faktorenanalyse mit dem Computer	229
8.2.1	Faktorenanalyse mit SPSS	229
8.2.2	Faktorenanalyse mit Stata	231

8.3 Übungsaufgaben zur Faktorenanalyse	232
Literatur	234
9 Lösungen der Übungsaufgaben	235
Formelsammlung	253
Sachverzeichnis	261

Abbildungsverzeichnis

Abb. 1.1	Von den Daten über die Information zum Wissen	5
Abb. 1.2	Preis-Absatz-Funktion für eine sensitive Zahnpasta	6
Abb. 1.3	Phasen empirischer Forschung	7
Abb. 1.4	Systematisierung von Modellen	10
Abb. 1.5	Was heißt schon sicher	11
Abb. 1.6	Intelligence Cycle	13
Abb. 2.1	Fragebogen Kundenbefragung Einzelhandel	18
Abb. 2.2	Merkmalsträger/Merkmale/Merkmalausprägung/Skalenniveau	19
Abb. 2.3	Kodierungsplan	23
Abb. 3.1	Dateneditor: Ansicht der eingegebenen Fragebögen	30
Abb. 3.2	Häufigkeitstabelle der Angebotsvielfalt	31
Abb. 3.3	Säulendiagramm bzw. Häufigkeitsverteilung für die Variable Angebot	31
Abb. 3.4	Verteilungsfunktion für die Variable Angebot	32
Abb. 3.5	Unterschiedliche Darstellung gleicher Sachverhalte (1) ...	32
Abb. 3.6	Unterschiedliche Darstellung gleicher Sachverhalte (2) ...	33
Abb. 3.7	Klassierung der Daten durch ein Histogramm	34
Abb. 3.8	Verletzung der Flächentreue und Verteilungsfunktion	35
Abb. 3.9	Notendurchschnitt zweier Klausuren	37
Abb. 3.10	Mittelwert als ausgeglichene Balkenwaage	38
Abb. 3.11	Mittelwert und getrimmter Mittelwert anhand des Zoobeispiels	39
Abb. 3.12	Berechnung des Mittelwerts aus klassierten Daten	40
Abb. 3.13	Geometrisches Mittel: Ein Beispiel	43
Abb. 3.14	Der Median als zentraler Wert unklassierter Daten	48
Abb. 3.15	Der Median als zentraler Wert klassierter Daten	49
Abb. 3.16	Berechnung von Quantilen bei fünf Gewichten	51
Abb. 3.17	Boxplot der Variablen Absatz pro Woche	52
Abb. 3.18	Interpretationen unterschiedlicher Formen eines Boxplots	53
Abb. 3.19	Variationskoeffizient	58
Abb. 3.20	Schiefe	59
Abb. 3.21	Idee des Dritten Zentralen Moments	60
Abb. 3.22	Kurtosis einer Verteilung	61

Abb. 3.23	Robustheit von Parametern	62
Abb. 3.24	Konzentrationsmaße	64
Abb. 3.25	Lorenzkurve	65
Abb. 3.26	Univariate Parameter mit SPSS	67
Abb. 3.27	Univariate Parameter mit Stata	68
Abb. 3.28	Univariate Parameter mit Excel	69
Abb. 3.29	KFZ Produktion in GB	70
Abb. 4.1	Kontingenztabelle (Kreuztabelle)	74
Abb. 4.3	Kontingenztabellen (Kreuztabellen): Geschlecht und Kauf [2. Teil]	76
Abb. 4.4	Berechnung erwarteter Häufigkeiten in Kontingenztabellen	78
Abb. 4.5	Chi-Quadrat-Werte bei unterschiedlicher Anzahl von Beobachtungen	81
Abb. 4.6	Phi bei einer unterschiedlichen Anzahl von Zeilen und Spalten	82
Abb. 4.7	Kontingenzkoeffizient bei unterschiedlicher Zeilen- und Spaltenzahl	84
Abb. 4.8	Kreuztabellen und nominale Zusammenhänge mit SPSS Titanic	87
Abb. 4.9	Von den Rohdaten zur computerberechneten Kreuztabelle (Titanic)	88
Abb. 4.10	Computerausdruck Chi-Quadrat und nominale Zusammenhangsmaße	88
Abb. 4.11	Kreuztabellen und nominale Zusammenhangsmaße mit Stata (Titanic)	89
Abb. 4.12	Kreuztabellen und nominale Zusammenhangsmaße mit Excel (Titanic)	90
Abb. 4.13	Streudiagramm	95
Abb. 4.14	Unterschiedliche Aspekte bei Streudiagrammen	96
Abb. 4.15	Unterschiedliche Darstellung gleicher Sachverhalte (3) ...	97
Abb. 4.16	Zusammenhang der Körpergrößen bei Hochzeiten	99
Abb. 4.17	Vier-Quadranten-Schema	99
Abb. 4.18	Produkt-Moment-Korrelation bei Ausreißern	103
Abb. 4.19	Fragebogenteil zum Design einer Weinflasche	103
Abb. 4.20	Nichtlinearer Zusammenhang zweier Variablen	103
Abb. 4.21	Daten zur Befragung zum Design einer Weinflasche	104
Abb. 4.22	Rangplätze zur Befragung zum Design einer Weinflasche	106
Abb. 4.23	Kendalls Tau bei einem perfekten positiven monotonen Zusammenhang	109
Abb. 4.24	Kendalls Tau bei fehlendem monotonen Zusammenhang	111
Abb. 4.25	Kendalls Tau bei Rangbindungen	112
Abb. 4.26	Kendalls Tau-b aus einer Kontingenztabelle	114
Abb. 4.27	Punktbiseriale Korrelation	116
Abb. 4.28	Zusammenhang zwischen einer ordinalen und metrischen Variablen	119
Abb. 4.29	Korrelationsrechnung mit SPSS	120
Abb. 4.30	Korrelationsrechnung mit Stata (Beispiel: Kendalls Tau)	121
Abb. 4.31	Korrelation nach Spearman mit Excel	123
Abb. 4.32	Gründe für Scheinkorrelationen	125
Abb. 4.33	Superbenzin und Marktanteil: Ein Beispiel für eine Scheinkorrelation	127
Abb. 4.34	Partielle Korrelation mit SPSS – Beispiel Superbenzin	128
Abb. 4.35	Partielle Korrelation mit Stata (Superbenzin)	129
Abb. 4.36	Partielle Korrelation mit Excel (Superbenzin)	130

Abb. 5.1	Prognose der Nachfrage mit Hilfe der Äquivalenzmethode	137
Abb. 5.2	Prognose der Nachfrage mit Hilfe der Abbildungsgröße	137
Abb. 5.3	Berechnung von Residuen	140
Abb. 5.4	Ausgleichsgeraden mit Bedingung „minimale Summe der Abweichungen“	140
Abb. 5.5	Die Idee multivariater Analysen	145
Abb. 5.6	Regression mit Excel und SPSS	147
Abb. 5.7	Regressionsoutput der Funktion Regression bei SPSS	148
Abb. 5.8	Regressionsoutput mit einer Dummy-Variablen	151
Abb. 5.9	Grafische Wirkung einer Dummy-Variablen	152
Abb. 5.10	Leverage Effekt	153
Abb. 5.11	In den Variablen nichtlineare Verläufe	155
Abb. 5.12	Beispiel eines in den Variablen nichtlinearen Verlaufs (1)	156
Abb. 5.13	Beispiel eines in den Variablen nichtlinearen Verlaufs (2)	157
Abb. 5.14	Autokorrelierte und nicht-autokorrelierte Verläufe der Fehlerterme	159
Abb. 5.15	Homoskedastizität und Heteroskedastizität	160
Abb. 5.16	Lösung bei perfekter Multikollinearität	161
Abb. 5.17	Lösung bei nicht perfekter Multikollinearität	163
Abb. 5.18	Getränkgröße	170
Abb. 6.1	Preisentwicklung von Dieselkraftstoff	172
Abb. 6.2	Preisentwicklung von Kraftstoffen	174
Abb. 6.3	Beispiel für Lohnentwicklung in zwei Unternehmen	183
Abb. 7.1	Bierdatensatz	191
Abb. 7.2	Distanzberechnung 1	192
Abb. 7.3	Distanzmessung 2	194
Abb. 7.4	Distanzmatrix	196
Abb. 7.5	Abfolge der Fusionsschritte	198
Abb. 7.6	Zuordnungsübersicht	198
Abb. 7.7	Fusionierungsalgorithmen (Linkage-Verfahren)	200
Abb. 7.8	Dendrogramm	202
Abb. 7.9	Screeplot zur Identifizierung von sprunghaften Heterogenitätszuwächsen	203
Abb. 7.10	Bewertung der F-Werte für die Clusterlösungen 2 bis 5	203
Abb. 7.11	Fehlklassifizierung im Vergleich mit Ergebnissen der Diskriminanzanalyse	204
Abb. 7.12	Interpretation der Cluster	205
Abb. 7.13	Anfangspartition der Clusterzentrenanalyse	207
Abb. 7.14	Hierarchische Clusteranalyse mit SPSS	208
Abb. 7.15	Clusterzentrenanalyse mit SPSS	210
Abb. 7.16	Clusteranalyse mit Stata	211
Abb. 7.17	Zuordnungsübersicht	212
Abb. 7.18	Dendrogramm	213

Abb. 7.19	Streudiagramm Persönliche Zufriedenheit und Einkommen	214
Abb. 8.1	Attribute zur Beschreibung von Zahnpastaeigenschaften	218
Abb. 8.2	Screeplot für das Zahnpastabeispiel	225
Abb. 8.3	Varimax Rotation für das Zahnpastabeispiel	227
Abb. 8.4	Faktorenanalyse mit SPSS	230
Abb. 8.5	Faktorenanalyse mit Stata	231
Abb. 9.1	Säulendiagramm und Histogramm	238
Abb. 9.2	Streudiagramm	241
Abb. 9.3	Clusteranalyse Persönliche Zufriedenheit und Einkommen (1)	250
Abb. 9.4	Clusteranalyse Persönliche Zufriedenheit und Einkommen (2)	251

Tabellenverzeichnis

Tab. 2.1	Amtliche Statistiken nationaler Institutionen	16
Tab. 2.2	Nichtamtliche Statistiken nationaler Institutionen	16
Tab. 2.3	Statistiken internationaler Institutionen	16
Tab. 3.1	Beispiel für die Mittelwertberechnung aus klassierten Daten	41
Tab. 3.2	Harmonisches Mittel	45
Tab. 3.3	Absatzanteile nach Altersklassen für Windelbenutzer	47
Tab. 3.4	Absatz von Fahrzeugen	71
Tab. 3.5	Automobilpreise	71
Tab. 4.1	Zusammenhangsmaße und Skalenniveaus	74
Tab. 4.2	Blödzeitung	132
Tab. 6.1	Durchschnittswerte für Diesel- und Ottokraftstoffe in Deutschland	173
Tab. 6.2	Verkettung von Indizes für Vorwärts- und Rückwärtsrechnung	185
Tab. 6.3	Preis-/Mengenentwicklungen	186
Tab. 6.4	Preis und Wertindex	186
Tab. 7.1	Distanz- und Ähnlichkeitsmaße in Abhängigkeit vom Skalenniveau	195
Tab. 7.2	Clusterzentren der endgültigen Lösung	213
Tab. 7.3	Cluster Zugehörigkeit	214
Tab. 8.1	Korrelationsmatrix der Faktorenanalyse	219
Tab. 8.2	Inverse der Korrelationsmatrix	219
Tab. 8.3	Bewertungsintervalle des Kaiser-Meyer-Olkin-Kriteriums	220
Tab. 8.4	Prüfung der Korrelationsmatrix durch KMO und Bartlett's Test	220
Tab. 8.5	Anti-Image-Korrelationsmatrix	221
Tab. 8.6	Eigenwerte und erklärte Gesamtvarianz für die Zahnpastaeigenschaften	223
Tab. 8.7	Reproduzierte Korrelation und Residuen zur Ursprungsmatrix	224
Tab. 8.8	Unrotierte und rotierte Faktormatrix der Zahnpastaeigenschaften	226
Tab. 8.9	Koeffizientenmatrix der Faktorscores anhand des Zahnpastabeispiels	228
Tab. 8.10	KMO und Bartlett's Test	232
Tab. 8.11	Anti-Image-Matrizen	232
Tab. 8.12	Kommunalitäten	233
Tab. 8.13	Erklärte Gesamtvarianz	233
Tab. 8.14	Rotierte Faktormatrix	233

Tab. 9.1	Preis und Absatz nach Ländern	241
Tab. 9.2	Preis-/Mengenentwicklungen (Lösung)	247
Tab. 9.3	Preis und Wertindex (Lösung)	249

1.1 Statistik lügt?

Ich glaube keiner Statistik, die ich nicht selbst gefälscht habe.

Mit Statistik kann man alles beweisen.

Diese und sicherlich noch viele ähnliche Aussagen finden sich im täglichen Leben, wenn es darum geht, das Zahlenwerk eines Gegenübers zu diskreditieren. So wird die Aussage „Es gibt drei Arten von Lügen: Lügen, verdammte Lügen und Statistiken“ gerne jenem englischen Parlamentarier des 19. Jahrhunderts zugeordnet, den man mit statistisch aufbereiteten Zahlen in die Enge getrieben hatte. Letztlich unterstellt diese Aussage, dass Statistik bzw. deren angewandte Methoden eine besonders hinterhältige Form der Lüge darstellen. Bestätigung finden die Kritiker nicht selten dann, wenn durch ein Gutachten und ein entsprechendes Gegengutachten auf statistischem Wege zwei entgegengesetzte Thesen abgeleitet werden. Wofür also Statistik, wenn anscheinend „jedes Ding zwei Seiten hat“, wenn bewiesen werden kann, was man gerne hätte und die Statistik scheinbar zu einem manipulativen Instrument der Person wird, welche die Statistik erstellt.

Obwohl solche Aussagen gerne kopfnickend, schmunzelnd oder sogar zustimmend aufgenommen werden – dies vor allem von denen, die von statistischen Verfahren eher rudimentäre bis gar keine Kenntnis besitzen – scheint gerade die Statistik eine der zentralen Methoden zu sein, mit denen Aussagen belegt werden. Man schlage an einem beliebigen Tage eine Tages- oder Wochenzeitung auf und man trifft auf Tabellen, Diagramme, Zahlen und Fakten. Kein Monat vergeht ohne Politbarometer, Geschäftsklimaindex, Konjunkturprognosen, Konsumentenindex, etc. Viele Anleger vertrauen bei ihrer Geldanlage den Entwicklungsprognosen der Aktien im DAX und hoffen auf die Erfüllung der Prognosen der Finanzmarktökonometer.

Wieso scheint hier nun die eben noch gescholtene Statistik einen unwiderstehlichen Zauber, eine Magie der Präzision der Zahlen auszustrahlen? Wie kommt es, dass der oben

beschriebene Superlativ von Lügen – Statistiken – auf einmal zur Grundlage der Planung von Privatpersonen und Unternehmen wird? Swoboda (1971, S. 16) nennt für diese Unentschlossenheit gegenüber statistischen Verfahren zwei wesentliche Gründe:

- „Erstens die *mangelnde Kenntnis* statistischer Aufgaben, Methoden und Möglichkeiten, und
- zweitens der Umstand, dass vieles für Statistik gehalten wird, was lediglich *Pseudostatistik* ist“.

Insbesondere der erste Punkt ist seit den 70er Jahren des letzten Jahrhunderts noch wichtiger als zuvor. Jedem, der die vier Grundrechenarten beherrscht, wird die Kompetenz zugetraut, Statistiken zu erstellen. Im Zeitalter von Standardsoftware, in dem prinzipiell ein Mausklick genügt, um eine Tabelle, eine Grafik oder sogar eine Regression zu erzeugen, wird dem Laien der Schritt zu komplizierten Anwendungen leicht gemacht. Nicht selten werden dabei Annahmen verletzt, Sachverhalte bewusst – also manipulativ – oder unbewusst verkürzt dargestellt. Zudem werden sorgsam ausgearbeitete Statistiken von Lesern und Zweitverwertern unachtsam oder falsch interpretiert und weitergegeben. Dabei ist es nicht nur „die Presse“, die hier in die Falle der statistischen Methoden gerät, sondern auch in mancher wissenschaftlichen Abhandlung findet sich ähnliche Unzulänglichkeit, die Swoboda als Pseudostatistik bezeichnet. Hier liegt der eigentliche Grund dafür, dass Statistik einerseits Hilfsmittel und andererseits „Lüge“ sein kann. Die bewusst oder unbewusst falsche Anwendung statistischer Methoden sowie die bewusst oder unbewusst falsche Interpretation der Ergebnisse dieser Verfahren.

Krämer (2005, S. 10) fasst die Gründe für „falsche“ Statistiken so zusammen. „Einige [Statistiken] sind bewusst manipuliert, andere nur unpassend ausgesucht. In einigen sind schon die reinen Zahlen falsch, in anderen sind die Zahlen nur irreführend dargestellt. Dann wieder werden Äpfel mit Birnen zusammengeworfen, Fragen suggestiv gestellt, Trends fahrlässig fortgeschrieben, Raten, Quoten oder Mittelwerte kunstwidrig berechnet, Wahrscheinlichkeiten vergewaltigt oder Stichproben verzerrt.“ Im vorliegenden Buch werden wir eine Reihe solcher Beispiele für falsche Interpretationen oder für Manipulationsversuche kennen lernen. Und somit wäre das Ziel dieses Buches klar umrissen: Die bereits in Goethes Gesprächen mit Eckermann betonte Notwendigkeit, quantitative Verfahren zu verstehen („das aber weiß ich, dass die Zahlen uns belehren“), sie zu durchschauen und selbst anwenden zu können, ist in einer Welt, in der uns täglich Daten, Zahlen, Trends und Statistiken umgeben, unumgänglich geworden. Statistische Modelle und Methoden sind entsprechend zu einem wichtigen Instrument in der betriebswirtschaftlichen Problemanalyse, der Entscheidungsfindung und der Unternehmensplanung geworden. Vor diesem Hintergrund sollen nicht nur die wichtigsten Methoden und deren Möglichkeiten vermittelt, sondern ebenfalls der Sinn für Irrtumsquellen und Manipulationsversuche geschärft werden.

Bis hierher könnte man nun der Auffassung sein, dass für die Anwendung der Statistik der gesunde Menschenverstand ausreicht und die Mathematik bzw. formale Darstellun-

gen in Form von Modellen keine Rolle spielen. Derjenige, der jemals in den Genuss einer gängigen Statistikvorlesung gekommen ist, wird diese Meinung wohl kaum teilen. Selbstverständlich kommt auch dieses Lehrbuch nicht ohne Formeln aus. Wie könnte es auch, wenn schon in alltäglichen Fällen eine qualitative Beschreibung nicht ausreicht: Auf die studentische Frage, wie denn die Durchfallquote in der Statistik-Klausur sei, würde sich kein Student mit der Aussage *ganz ok* zufrieden geben. Vielmehr erwartet er hier eine Aussage wie beispielsweise *10 Prozent*, was wiederum nur rechnerisch – also mit einer Formel – zu ermitteln ist.

Es kann also auch in diesem Buch nicht auf ein Mindestmaß an formaler Darstellung verzichtet werden. Dennoch wird jeder bemühte Leser, der die Grundlagen der Analysis beherrscht, dieses Buch verstehen können.

1.2 Zwei Arten von Statistik

Was kennzeichnet nun aber eine Statistik oder Datenanalyse, die Irrtumsquellen und Manipulationsversuche möglichst ausschließt? Hierzu müssen wir uns zunächst darüber verständigen, was überhaupt die Aufgaben von Statistik bzw. von Datenanalyse sind.

Historisch gesehen, gehen die Methoden der Statistik weit vor Christi Geburt zurück. Schon im sechsten Jahrhundert vor Christi sah die Verfassung des Königs Servius Tullius eine periodische Erfassung aller Bürger vor. Vielen dürfte zudem folgende Geschichte bekannt sein: „Es begab sich aber zu der Zeit, dass ein Gebot von dem Kaiser Augustus ausging, dass alle Welt geschätzt würde. Und diese Schätzung war die allererste und geschah zu der Zeit, da Quirinius Statthalter in Syrien war. Und jedermann ging, dass er sich schätzen ließe, ein jeder in seine Stadt.“¹ (Lukas 2,1 ff.)

Politiker hatten also seit jeher das Interesse, die Leistungsfähigkeit der Bevölkerung bemessen zu können. Dies allerdings nicht uneigennützig, sondern mit dem Ziel, die Bevölkerung anhand dieser Leistungsfähigkeit besteuern zu können. Aus Sicht des Staatsapparates erfolgte die Sammlung von Daten mit dem Ziel der Gewinnung von Informationen über den eigenen Staat. Noch im heutigen statistischen Jahrbuch finden sich die Wurzeln dieser Interpretation von Statistik als *Staatsbeschreibung*: Abschnitte über Geographie und Klima, Bevölkerung, Familien und Lebensformen füllen die ersten Seiten des Statistischen Jahrbuches der Bundesrepublik Deutschland (Statistisches Bundesamt 2013).

Bei allen frühzeitlichen Statistiken handelt es sich um Vollerhebungen in dem Sinne, dass buchstäblich jede Person, jedes Tier, jedes Objekt gezählt wurden. Bis zum Beginn des 20. Jahrhunderts stand die Beschäftigung mit entsprechend großen Fallzahlen

¹ Im Jahr 6/7 n. Chr. wurde Judäa (mit Idumäa und Samaritanen) römische Prokuratur. Die Textstelle bezieht sich wahrscheinlich auf die unter Quirinius durchgeführte Volkszählung, bei der die Bewohner des Landes und ihr Besitz für die Erhebung von Steuern registriert wurden. Die Bibel verwendet für diesen Registrierungsprozess den Begriff „geschätzt“. Es könnte aber auch sein, dass sich diese Textstelle auf eine erste Erfassung um 8/7 v. Chr. bezieht.

im Vordergrund des Interesses. Diese Periode stellt den Ausgangspunkt der sogenannten deskriptiven (beschreibenden) Statistik dar.

Die **Deskriptive Statistik** beinhaltet somit alle Verfahren, mit denen sich durch die Beschreibung von Daten einer Grundgesamtheit (*engl.*: population) Informationen gewinnen lassen. Zu diesen Methoden bzw. Verfahren gehören unter anderem die Erstellung von Grafiken, Tabellen und die Berechnung von deskriptiven Kennzahlen bzw. Parametern.

Erst nach Beginn des 20. Jahrhunderts entwickelte sich die uns heute eher geläufige *Induktive (Schließende) Datenanalyse*, die versucht, aus Stichproben Schlüsse auf die *Gesamtheit* zu ziehen. Dominierend bei dieser Entwicklung waren unter anderem die Wissenschaftler Jacob Bernoulli (1654–1705), Abraham de Moivre (1667–1754), Thomas Bayes (um 1702–1761), Pierre-Simon Laplace (1749–1827), Carl Friedrich Gauß (1777–1855), Pafnuti Lwowitsch Tschebyschow (1821–1894), Francis Galton (1822–1911), Ronald A. Fisher (1890–1962) und William Sealy Gosset (1876–1937), auf die eine Vielzahl der heute bekannten induktiven Verfahren zurückgeht. Diesen Erkenntnissen ist es zu verdanken, dass heute nicht jede Person einer Grundgesamtheit, sondern nur eine Stichprobe (*engl.*: sample) von Personen befragt werden muss.

Erst viel später entwickelte sich die **Induktive (Schließende) Statistik**, die versucht, mit Hilfe von Stichproben Schlüsse auf die Gesamtpopulation zu ziehen. Das „Ziehen von Schlüssen“ aus einer Stichprobe führte letztlich auch zur Etablierung des Begriffes der Schließenden Statistik, der – wie auch der Begriff der Inferenzstatistik – häufig synonym zur Induktiven Statistik verwendet wird. Dominierend bei dieser Entwicklung waren unter anderem die Wissenschaftler Jacob Bernoulli (1654–1705), Abraham de Moivre (1667–1754), Thomas Bayes (um 1702–1761), Pierre-Simon Laplace (1749–1827), Carl Friedrich Gauß (1777–1855), Pafnuti Lwowitsch Tschebyschow² (1821–1894), Francis Galton (1822–1911), Ronald A. Fisher (1890–1962) und William Sealy Gosset (1876–1937), auf die eine Vielzahl der heute bekannten induktiven Verfahren zurückgeht. Diesen Erkenntnissen ist es zu verdanken, dass heute nicht jede Person einer Grundgesamtheit, sondern nur eine Stichprobe (*engl.*: sample) von Personen befragt werden muss. Dies erweist sich insbesondere dann als vorteilhaft, wenn Vollerhebungen zu teuer kämen bzw. zu lange dauern würden, oder die Erhebung mit einer Zerstörung der Untersuchungselemente einhergehen würde (z. B. bei bestimmten Formen der Materialprüfung wie z. B. auch Weinproben). Es wäre für Unternehmen sicherlich nicht finanzierbar, alle potenziellen Kunden darüber zu befragen, wie ein neues Produkt auszusehen hat. Es wird vielmehr mit einer entsprechend zusammengestellten Stichprobe gearbeitet. Auch die Wahlforscher könnten kaum alle Wahlberechtigten befragen.

Für den Auswertungsprozess bedeutet dies, dass das zu ermittelnde Wissen nun eben nicht mehr auf Daten einer Vollerhebung basiert, sondern auf besonders ausgewählten Daten einer Stichprobe. Entsprechend sind die zu ziehenden Schlüsse in Bezug auf die Grundgesamtheit auch mit einer Unsicherheit belegt. Das ist der Preis der Herangehensweise der Induktiven Statistik. Deskriptive und Induktive Statistik bilden somit eine wis-

² Früher auch als Tschebyschew, Tschebyschow oder Tschebyschew transkribiert.



Abb. 1.1 Von den Daten über die Information zum Wissen

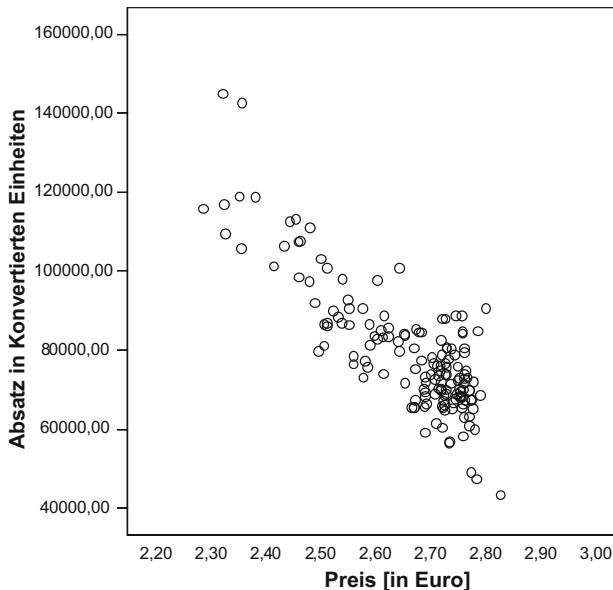
senschaftliche Disziplin für die Wirtschafts-, Sozial und Naturwissenschaften. Sie umfasst die Methoden zur Beschreibung und Analyse von Massenerscheinungen mit Hilfe von Zahlen und Daten. Auswertungsziel ist das Treffen von Aussagen bezüglich der Eigenschaften der Untersuchungseinheiten auf Basis einer Totalerhebung oder einer Stichprobe. Statistik ist eine Zusammenfassung von Methoden, welche es uns erlaubt, „vernünftige“ Entscheidungen im Falle von Unsicherheit zu treffen, und ist somit wichtigste Grundlage der *Entscheidungstheorie*.

Damit wären die beiden Hauptziele der Statistik abgrenzbar: Die Deskriptive Statistik beschränkt sich auf die zusammenfassende Darstellung von Daten und verarbeitet diese zu Informationen. Wenn diese Informationen mit Hilfe von Auswertungsverfahren der Induktiven Statistik analysiert werden, entsteht verallgemeinerbares Wissen, das politisches oder strategisches Handeln beeinflussen kann. Abbildung 1.1 stellt diesen Zusammenhang nochmals schematisch dar.

1.3 Statistik als Erkenntnisprozess

Die fundamentale Bedeutung der Statistik für den Erkenntnisprozess – oder anders ausgedrückt: für die Generierung neuen Wissens – darf nicht unterschätzt werden. Der Erkenntnisprozess in Wissenschaft und Praxis durchläuft nämlich in der Regel genau die beiden Stufen der Deskription und Induktion. Dies soll anhand eines kleinen Praxisbeispiels erläutert werden:

Ein Marktforscher aus dem Bereich der Zahnpflege stellt sich die Frage über den Zusammenhang zwischen dem Preis und dem Umfang der Abverkäufe einer bestimmten Zahnpasta. Zunächst wird er versuchen, sich durch Verdichtung von Einzelinformationen selbst ein Bild von der bestehenden Realität zu machen. So könnte er beispielsweise die Abverkäufe und die Preise der Zahnpasta pro Kalenderwoche innerhalb der letzten drei Jahre grafisch analysieren (vgl. Abb. 1.2). Wie immer bei der Datengewinnung, werden einzelne Verkaufsmärkte ihre Verkaufszahlen nicht regelmäßig melden, sodass keine Vollerhebung, sondern lediglich eine Teilerhebung vorliegt. Er stellt fest, dass bei hohen Preisen der Abverkauf zugunsten anderer Zahnpastaprodukte zurückgeht und bei niedrigen Preisen der Abverkauf entsprechend anzieht. Dieser deskriptiv ermittelte Zusammenhang entspricht nicht nur einer individuell gewonnenen Einsicht, sondern auch den Erwartungen aus der mikroökonomischen Theorie der Preis-Absatz-Funktion. In jedem Fall sind es die Methoden der Deskriptiven Statistik, mit deren Hilfe sich individuelle

**Lesehilfe:**

In der Grafik sind für drei Jahre à 52 Wochen die durchschnittlichen Preise sowie die dazugehörige Abverkaufsmenge in normierter Packungsgröße abgebildet. Jeder Punkt stellt somit eine Kombination aus Preis und Abverkaufsmenge einer bestimmten Kalenderwoche dar.

Abb. 1.2 Preis-Absatz-Funktion für eine sensitive Zahnpasta

Erkenntnisse aus Einzelinformationen gewinnen lassen und sich bestehende Erwartungen oder Theorien anhand der Verdichtung von Einzelfällen anschaulich machen lassen.

Der Forscher wird sich an dieser Stelle die Frage stellen, ob sich die aus der Teilerhebung gewonnenen Erkenntnisse – die er zudem theoretisch vorher schon vermutet hatte – für die Grundgesamtheit verallgemeinern lassen. Verallgemeinernde Informationen der Deskriptiven Statistik sind nämlich zunächst spekulativ. Mit Hilfe der Verfahren der Induktiven Statistik lässt sich aber das Risiko in Form einer Fehlerwahrscheinlichkeit bei der Übertragung der Ergebnisse der Deskriptiven Statistik auf die Grundgesamtheit bemessen. Der Forscher muss selbst entscheiden, ob er das Risiko einer Übertragung als zu hoch empfindet und die Erkenntnisse als ungesichert qualifiziert und vice versa.

Selbst wenn alle Verkaufsstellen ihre Verkaufszahl gemeldet hätten und somit eine Vollerhebung vorläge, könnte er sich die Frage stellen, ob dieser Zusammenhang zwischen Preis und Absatz *ceteris paribus* auch zukünftig noch gilt. Werte für die Zukunft liegen nämlich auf keinen Fall vor, sodass aus der Vergangenheit auf die Zukunft geschlossen werden müsste. Nur auf diese Weise lassen sich *Theorien, Annahmen und Erwartungen verifizieren* und nur so lässt sich Information in verallgemeinerbares Wissen (für das Unternehmen) transformieren.

Deskriptive und Induktive Statistik erfüllen im Forschungsprozess somit unterschiedliche Aufgaben, sodass eine differenzierte Betrachtung dieser beiden Bereiche als sinnvoll erachtet werden kann und in der Lehre häufig auch in verschiedenen Veranstaltungsteilen abgehandelt werden.

1.4 Phasen empirischer Forschung

Das obige Beispiel verdeutlicht zudem, dass der Ablauf eines Erkenntnisprozesses bestimmte Stufen durchläuft, die in Abb. 1.3 als Phasen empirischer Forschung schematisch dargestellt sind. In der *Erkundungsphase* geht es zunächst darum, sich selbst ein Bild über mögliche Zusammenhänge zu verschaffen, um diese danach in der *Theoriephase* zu einem konsistenten Modell zu verknüpfen.

1.4.1 Von der Erkundung zur Theorie

Obwohl der „Praktiker“ den Begriff der *Theorie* nur ungern verwendet, ihn im Gegenteil eher meidet, da er sonst als „weltfremd, unzugänglich, unrealistisch“ gelten könnte, so steht dieser Begriff zunächst am Anfang eines jeden Erkenntnisfortschritts. Die Herkunft

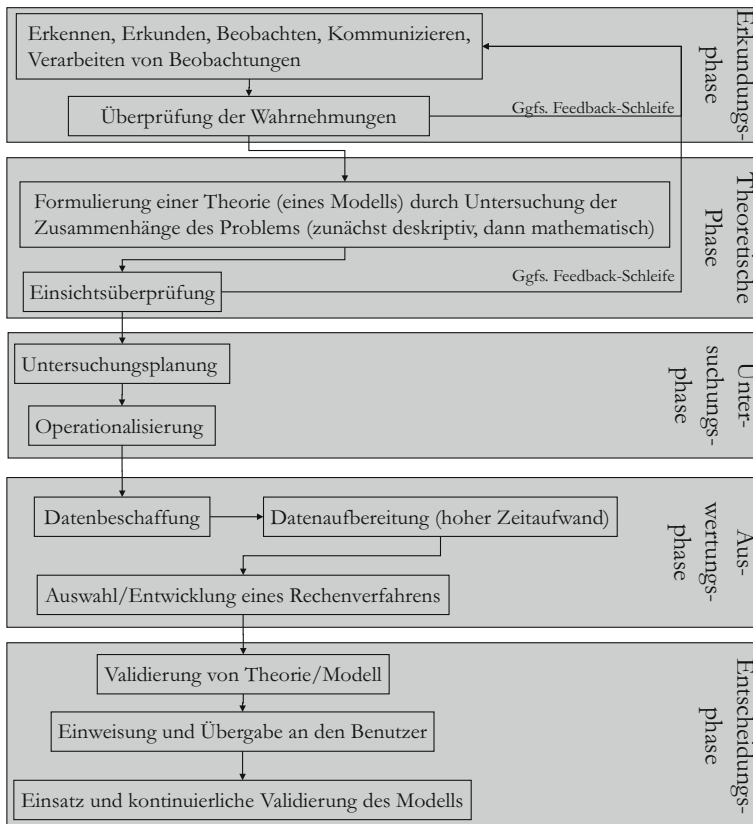


Abb. 1.3 Phasen empirischer Forschung

des Wortes Theorie leitet sich vom griechischen Wort *theorema* ab, welches mit *anschauen, betrachten, untersuchen* übersetzt werden kann. Theorie ist somit die Erkenntnis von Systemen, die zunächst eine spekulative Annäherung an einen Sachverhalt darstellt (Crow 2005, S. 14). Bereits hieraus lässt sich also schließen, dass die Aufstellung einer Theorie auf der Beobachtung und Verknüpfung von Einzelereignissen beruht, die ohne Überprüfung nicht als allgemeingültig gelten kann. Eine *erfahrungswissenschaftliche Theorie* verknüpft die Einzelereignisse der Realität, sodass bei bestimmten Anwendungsbedingungen von Tatbeständen auf Ursachen geschlossen werden kann. Kern einer jeden Theorie ist somit die Aufstellung eines einheitlichen Begriffsapparates – oder auch sprachlichen Systems – aus dem sich *gesetzmäßige Ursache-Wirkungsbeziehungen* ableiten lassen. Für unser Zahnpaste-Beispiel bedeutet dies, dass der Forscher sich zunächst einmal Gedanken darüber zu machen hat, welche Ursachen (Faktoren) auf den Absatz seines Produktes wirken. „Aus dem Bauch“ fallen dem Forscher sicherlich die wichtigsten Ursachen ein: der Preis des eigenen Produktes, der Preis der Konkurrenzprodukte, Werbemaßnahmen der Eigen- und Fremdprodukte, die Marktsegmentierung hin zu Spezialzahnpaste (Zahnweiß, empfindliche Zähne etc.).

Neben diesen Aspekten spielen in der Regel auch Ursachen eine Rolle, die dem Nichtkenner einer Branche verborgen bleiben. In Abb. 1.3 sind sowohl in der Erkundungsphase als auch in der Phase der Theoriebildung *Feedback-Schleifen* eingefügt, in denen eigene Wahrnehmungen und Einsichten von einem selbst oder von Dritten überprüft werden sollten. Eine quantitative Studie erfordert deshalb immer auch ein Höchstmaß an *kommunikativer Kompetenz*. Kontaktaufnahme zu den Branchenkennern – wie z. B. Produktmanagern –, die dem Forscher auch zunächst verborgene Ereignisse und Einflüsse erklären können, gehört deshalb zur Aufgabe einer jeden ordentlichen quantitativen Studie. Dies gilt selbstverständlich auch für Studien aus anderen Funktionsbereichen des Unternehmens: In der Beschaffungsforschung sind Einkäufer zu fragen, in der Produktionsforschung die Ingenieure und Meister, in der Finanzmarktforschung die Analysten des Bereiches, etc. Diese Kommunikation verbessert nicht nur das Verständnis des Zusammenspiels von Ursachen und Wirkung für den Forscher, sondern sie verhindert letztlich auch die Pein, in der Endpräsentation von diesen Personen erst auf wichtige fehlende Einflüsse hingewiesen werden zu müssen.

1.4.2 Von der Theorie zum Modell

Nachdem die theoretischen Zusammenhänge festgestellt worden sind, beginnt die Modellbildung. Nicht selten werden die Begriffe Theorie und Modell synonym verwendet, obwohl sich der Begriff Theorie streng genommen auf die Beschreibung der Realität mit Hilfe der Sprache bezieht. Fasst man mathematische Formalisierung auch als Sprache mit eigener Grammatik und Semiotik auf, so könnte eine Theorie auch mathematisch formal gebildet werden. In der Praxis verwendet man an dieser Stelle aber eher den Begriff des Modells, bei dem Theorien auf bestimmte Tatbestände angewendet werden.

Man bedient sich des Kunstgriffs des Modells, um durch Kombination verschiedenster theoretischer Überlegungen zu einer näherungsweise Vorstellung von der Wirklichkeit zu kommen. Durch *Abstraktion* und *Vereinfachung* wird versucht, das Realproblem möglichst strukturgleich als Formalproblem in einem *Modell* abzubilden. Unter Struktur wird dabei die relevante Gesamtheit der Eigenschaften und Relationen des Ausschnitts aus der Wirklichkeit verstanden. Schematisch scheint die Bewältigung der betriebs- und volkswirtschaftlichen Komplexität damit gelöst: Man hat lediglich alle Daten bezüglich eines Untersuchungsobjektes zu sammeln, diese statistisch auszuwerten und adäquat zu kommunizieren, um eine rationale Entscheidung zum Wohle des Betriebes oder der Volkswirtschaft fällen zu können. In der Praxis kommt man allerdings ziemlich schnell zu dem Schluss, dass eine detaillierte umfassende Beschreibung der (betrieblichen) Wirklichkeit und damit auch des Entscheidungsprozesses mit all ihren Ursachen und Wirkungszusammenhängen kaum möglich ist. Die (betriebliche) Realität ist viel zu komplex, als dass wir sie in ihrer Fülle in allen Einzelheiten erfassen könnten. Völlig strukturgleich – oder wie man es auch nennt: *isomorph* – kann die Abbildung der Wirklichkeit niemals sein. Diese Aufgabe kann kein Modell erfüllen, sodass Modelle in aller Regel reduziert – oder auch: *homomorph* – sind.

Die Realitätsnähe eines Modells – und damit der Prozess der zunehmenden *Modellverfeinerung* – hat also Grenzen. Sie liegen dort, wo das Modell seine Durchschaubarkeit verliert. Das Modell muss handhabbar bleiben und es müssen mithin die für den jeweiligen Erkenntniszweck wesentlichen Eigenschaften und Relationen des Problems wiedergegeben werden. Modelle sind also durch Abstraktion gewonnene gedankliche Hilfsmittel zur übersichtlichen Darstellung von unanschaulichen Objekten und komplexen Vorgängen (Bonhoeffer 1948, S. 3 ff.). Das Modell ist lediglich eine *Approximation der Wirklichkeit* bzw. eine *Komplexitätsreduktion*. Für die Darstellung der Teilzusammenhänge stehen verschiedene Formen und Mittel der Abbildung zur Verfügung: Die anschaulichste Form stellt das *physische* oder *ikonische* Modell dar. Beispiele sind körperliche Nachbildungen (Holz-, Plastik- oder Gipsmodell eines Baukörpers oder Stadtteils), Landkarten bzw. Konstruktionszeichnungen. Innerhalb der Wirtschaftswissenschaften haben physische Modelle praktisch keine Bedeutung erlangt. Das spezifisch *Wirtschaftliche* ist rein geistiger Natur und schon deshalb nicht physisch abbildbar.

Die *symbolischen (sprachlichen) Modelle* sind für die Wirtschaftswissenschaft besonders wichtig. Mit Hilfe einer Sprache, mit ihrem System symbolischer Zeichen und dem zugehörigen System syntaktischer und semantischer Regeln wird die Struktur des zu untersuchenden Tatbestandes approximiert und in ihrer Problematik untersucht. Dient als Sprache die übliche Alltagssprache oder eine daraus entwickelte Fachsprache, so handelt es sich um ein *verbales Modell* oder um eine *Theorie*. Zunächst besteht ein verbales Modell also aus einer Ansammlung symbolischer Zeichen und Wörter. Aus diesen ergibt sich nicht sofort ein Sinn, wie beispielsweise an der Wortfolge „Weiß wohnt in Hamburg meine Oma Hund“ zu erkennen ist. Die Ergänzung einer fehlenden syntaktischen Gliederung in Subjekt, Prädikat und Objekt in „Meine Oma ist weiß und ihr Hund wohnt in Hamburg“ würde den Satz zwar verständlich aber nicht sinnvoll machen. Erst die Berücksichtigung

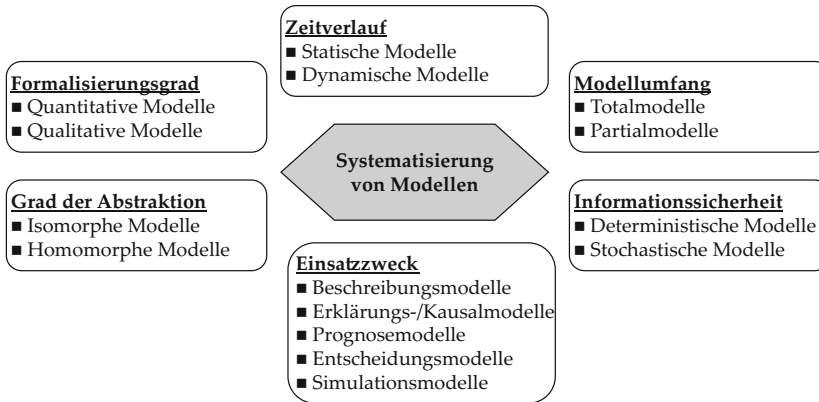


Abb. 1.4 Systematisierung von Modellen

der Semantik bzw. die Verknüpfung der Inhalte mit der entsprechenden Wortbedeutung verleiht dem verbalen Modell „Meine Oma wohnt in Hamburg und ihr Hund ist weiß“ einen Sinn.

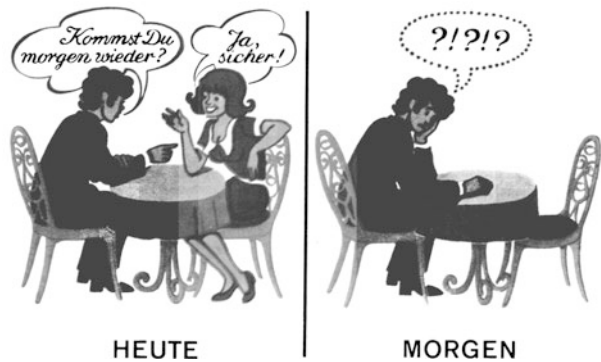
Gleiches gilt für künstliche Sprachen, wie logische und mathematische Systeme, die auch Kalküle oder *Symbolmodelle* genannt werden. Auch diese benötigen Zeichenfolgen (Variablen) sowie deren syntaktische und semantische Gliederung in Gleichungssystemen. Auf unser Zahnpaste Beispiel bezogen könnte ein mögliches verbales Modell bzw. eine Theorie lauten:

- Der Absatz der Zahnpaste hängt negativ von der Höhe des Preises und positiv von den eigenen Werbeausgaben der jeweiligen Periode (z. B. Kalenderwoche) ab.
- Das analoge formale Symbolmodell sähe so aus: $y_i = f(p_i, w_i) = \alpha_1 \cdot p_i + \alpha_2 \cdot w_i + \beta$.
 p : Preis zum Zeitpunkt i ; w_i : Werbeausgaben zum Zeitpunkt i ; α sind die entsprechenden Wirkungsgrade; β ist eine mögliche Konstante.

Bei beiden Modellarten handelt es sich um homomorphe *Partialmodelle*, da nur ein gewisser Teilausschnitt des unternehmerischen Handelns – hier Absatz eines Produktes – untersucht wird. Es war beispielsweise nicht von Interesse, auch die Beschäftigtenentwicklung und andere Größen mit zu berechnen. Dies würde man von *Totalmodellen* hingegen verlangen, was in den meisten Fällen aufgrund der damit verbundenen Komplexität aufwändig und auch sehr kostenintensiv ist. Vornehmlich in Wirtschaftsforschungsinstituten versucht man mit Totalmodellen zu arbeiten.

Bei den Verfahren der Statistik sind es die stochastischen, homomorphen Partialmodelle, die – zum Leidwesen vieler Studierender – Anwendung finden. Was bedeutet eigentlich der Begriff *stochastisch*? Nun, zunächst offenbart uns der Fremdwörterduden die Stochastik als das Teilgebiet der Statistik, das sich mit der Analyse *zufallsabhängiger Ereignisse* befasst und welches wir oben bereits als Induktive Statistik kennen gelernt haben. Mit

Abb. 1.5 Was heißt schon sicher? (Swoboda 1971, S. 31)



dem Begriff des *Zufalls* müssen wir uns immer dann beschäftigen, wenn wir über das Eintreten bestimmter Ereignisse keine vollständige Kenntnis haben, die Ereignisse also nicht *deterministisch* sind. Über die Zukunft lassen sich beispielsweise genauso wenig sichere Aussagen treffen wie über eine Grundgesamtheit, die wir durch eine Stichprobe nur zum Teil erfragen konnten. Als sicher kann bestenfalls – und das auch nicht immer – die Vergangenheit gelten. Am bemitleidenswerten Verehrer in Abb. 1.5 zeigen sich die alltagssprachlich bedingten Missverständnisse der Begriffe *Gewissheit* und *Sicherheit*.

Die Betriebs- und Volkswirtschaftslehre können sich nicht mit der Erkenntnis zufrieden geben, dass alles im Leben nun mal unsicher sei und man damit zu leben habe. Vielmehr wird im Rahmen der Induktiven Statistik bzw. der Stochastik der Versuch unternommen, den Grad der Sicherheit des Eintretens eines bestimmten Ereignisses zu schätzen. Zwar wäre obigem Verehrer wenig geholfen, wenn die Auserwählte ihr Kommen mit einer 95-prozentigen Wahrscheinlichkeit (also höchstwahrscheinlich) angegeben hätte. Es käme aber deutlich zum Ausdruck, dass das im Alltag verwendete *ja* und *nein*, *ganz sicher* oder *bestimmt nicht* immer mit einem gewissen Zweifel belegt ist. Diesen Zweifel oder diese Unsicherheit der Statistik anzulasten wäre insofern ungerechtfertigt, als die Statistik eben versucht, das Ausmaß von Sicherheit und Unsicherheit zu quantifizieren und nicht über die Zufälle, das Eintreten des Unwahrscheinlichen und die Überraschungen des Lebens hinweg zu sehen (Swoboda 1971, S. 30).

Ein anderer wichtiger Gliederungsgesichtspunkt ist der *Einsatzzweck* eines Modells. So kann unterschieden werden zwischen:

- Beschreibungsmodell,
- Erklärungsmodell/Prognosemodell,
- Entscheidungsmodell/Optimierungsmodell,
- Simulationsmodell.

Welchen Einsatzzweck ein Modell erfüllen muss, hängt dabei letztlich von der Fragestellung selbst bzw. deren Komplexität ab.

Ein *Beschreibungsmodell* versucht zunächst nichts anderes als die Realität durch ein Modell zu beschreiben. Allgemeingültige Hypothesen über Wirkungszusammenhänge im realen System enthält es hingegen nicht. So ist eine Bilanz oder eine Gewinn- und Verlustrechnung eines Unternehmens nichts anderes als der Versuch, die finanzielle Situation eines Unternehmens modellhaft darzustellen. Annahmen über Wirkungszusammenhänge zwischen einzelnen Bilanzpositionen werden dabei nicht aufgestellt oder überprüft.

In *Erklärungsmodellen* werden dagegen zunächst theoretische (hypothetische) Annahmen über Wirkungszusammenhänge aufgestellt und mit Hilfe empirischen Datenmaterials überprüft. So lassen sich auf quantitativer Basis Gesetzmäßigkeiten innerhalb des betrieblichen Geschehens aufdecken und zum Teil auf die Zukunft übertragen. Im letzteren Fall – also auf die Zukunft gerichteter Aussagen – spricht man von *Prognosemodellen*, die deshalb auch zur Gruppe der Erklärungsmodelle gezählt werden (Domschke und Drexl 2011, S. 1 ff.). Auf unser Zahnpaste Beispiel bezogen, stellt die Ermittlung der Erhöhung des Absatzes um beispielsweise 10.000 Tuben bei einer Preissenkung von 10 €-Cent ein Erklärungsmodell dar. Von einem Prognosemodell würde man sprechen, wenn durch eine in dieser Kalenderwoche (zum Zeitpunkt t) durchgeführte Erhöhung des Preises um 10 €-Cent eine Verringerung des Absatzes in der *nächsten* Kalenderwoche (also zum Zeitpunkt $t + 1$) um 8500 Einheiten prognostiziert werden könnte.

Unter *Entscheidungsmodellen* (Optimierungsmodellen) versteht Grochla (1969, S. 382) „auf die Ableitung von Handlungsmaßnahmen gerichtete Satzsysteme“. Charakteristisch für Entscheidungsmodelle ist die Generierung von optimalen Entscheidungen. Grundlage ist in der Regel die Existenz einer mathematischen Zielfunktion, die der Anwender des Modells unter Einhaltung bestimmter mathematischer Nebenbedingungen optimieren möchte. Derartige Modelle finden vornehmlich im Operations Research und weniger in der statistischen Datenanalyse Anwendung (vgl. z. B. Runzheimer et al. 2005).

In *Simulationsmodellen* werden Abläufe und Vorgänge – z. B. in einem Produktionssystem – nachgespielt. Der Computer mit seinem Zufallszahlengenerator eröffnet dabei die Möglichkeit, deren Abhängigkeit von stochastischen Einflussfaktoren (z. B. schwankende Ankunfts- oder Abfertigungsraten) offen zu legen. Aber auch Rollenspiele bei Führungsseminaren oder die Familienaufstellung der Psychologen können als Simulationen gelten.

1.4.3 Vom Modell zur „Business Intelligence“

Mit Hilfe statistischer Verfahren können selbst schwierigste Sachverhalte in ebenso komplexen statistischen Methoden verarbeitet werden. Diese Methoden gehen zum Teil weit über die in diesem Lehrbuch gezeigten Verfahren hinaus. Begnadet ist der Wissenschaftler und auch Praktiker, der diese Verfahren beherrscht. Allerdings kennt auch jeder die folgende oder eine ähnliche Situation: Ein engagierter, aber etwas vergeistigter Professor versucht einer Gruppe von Praktikern die Vorzüge des *Heckman Selection Model* mit Hilfe des dazugehörigen Artikels (siehe Heckman 1976) zu erklären. Die meisten Zuhö-