

Statistik mit SAS

Von Prof. Dr. rer. nat. Julius Dufner
Priv.-Doz. rer. nat. Uwe Jensen
Dr. rer. nat. Erich Schumacher

Universität Hohenheim

Mit zahlreichen Abbildungen,
Beispielen und Übungsaufgaben



Springer Fachmedien Wiesbaden GmbH

Prof. Dr. rer. nat. Julius D. Dufner

Geboren 1941 in Freiburg . Br. Von 1960 bis 1967 Studium der Mathematik und Physik an der Universität Freiberg. Ab 1967 Assistententätigkeit am Mathematischen Institut der Universität Freiburg Promotion 1971. Von 1972 bis 1974 Assistententätigkeit an der Pädagogischen Hochschule Freiburg, zweites Staatsexamen. Dozent an der FH Darmstadt von 1974 bis 1976, ab 1976 an der Berufspädagogischen Hochschule Esslingen, 1979 Professor. Von 1988 an Professor an der Universität Hohenheim.

Privatdozent Dr. rer. nat. Uwe Jensen

Geboren 1950 in Bremen. Von 1971 bis 1976 Studium der Mathematik, Physik und Betriebswirtschaft an der Technischen Universität Braunschweig, Diplom 1976. 1979 Promotion und 1987 Habilitation im Fach Mathematik an der Universität Stuttgart-Hohenheim. 1976/77 Industrietätigkeit in Frankfurt. Von 1977 bis 1980 Wiss. Angestellter, seit 1980 Akademischer Rat/Oberrat am Institut für Angewandte Mathematik und Statistik der Universität Hohenheim.

Dr. rer. nat. Erich Schumacher

Geboren 1945 in Bonladen. Studium der Mathematik an der Universität Stuttgart, Diplom 1969. Von 1970 bis 1974 Wiss. Assistent, seit 1975 Wiss. Angestellter am Institut für Angewandte Mathematik und Statistik der Universität Hohenheim. 1979 Promotion in Hohenheim.

Die Deutsche Bibliothek – CIP-Einheitsaufnahme

Dufner, Julius:

Statistik mit SAS : mit Beispielen und Übungsaufgaben / von Julius Dufner ; Uwe Jensen ; Erich Schumacher. – Stuttgart : Teubner, 1992

(Teubner-Studienbücher : Mathematik)

ISBN 978-3-519-02088-2 ISBN 978-3-322-94766-6 (eBook)

DOI 10.1007/978-3-322-94766-6

NE: Jensen, Uwe.; Schumacher, Erich:

Das Werk einschließlich aller seiner Teile ist urheberrechtlich geschützt. Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlages unzulässig und strafbar. Das gilt besonders für Vervielfältigungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Verarbeitung in elektronischen Systemen.

© Springer Fachmedien Wiesbaden 1992

Ursprünglich erschienen bei B. G. Teubner Stuttgart 1992

Gesamtherstellung: Druckhaus Beltz, Hemsbach/Bergstraße

Umschlagentwurf: P.P.K,S-Konzepte T. Koch, Ostfildern/Stuttgart

Vorwort

Aufgrund einer in den letzten Jahren sprunghaft gewachsenen Verfügbarkeit über Rechnerkapazitäten, insbesondere im Bereich der Personal Computer (PC), lassen sich heute auch umfangreiche und aufwendige statistische Datenanalysen innerhalb kürzester Zeit ausführen. Die zunehmende Bedeutung der Statistik in nahezu allen Wissenschaftsdisziplinen geht nicht zuletzt zurück auf diese gewachsenen Möglichkeiten, eine statistische Datenanalyse praktisch durchzuführen. Dafür ist ein geeignetes Statistik-Softwarepaket erforderlich. SAS (Statistical Analysis System) zählt zu den am weitesten verbreiteten und leistungsfähigsten Software-Systemen dieser Art.

Das Buch richtet sich an den Anwender statistischer Verfahren. Damit ist einerseits der Nichtmathematiker gemeint, der durch Beobachtungen oder aufgrund von Experimenten Daten gesammelt hat und diese für eine geordnete Darstellung aufbereiten möchte und Schlussfolgerungen aus den gewonnenen Daten ziehen will. Dazu werden Verfahren der beschreibenden und der beurteilenden Statistik herangezogen. Diese Verfahren sollen dann mit Hilfe einer leistungsfähigen Statistik-Software auf einem Rechner umgesetzt werden.

Andererseits richtet sich dieses Buch auch an den Mathematik-Studenten (Dozenten) mit Interesse an der angewandten Stochastik, der die in den Statistikvorlesungen vermittelten Verfahren mit Hilfe eines Computers realisieren möchte. Auch der erfahrene Statistiker kann, so hoffen wir, an der einen oder anderen Stelle Nutzen aus diesem Buch ziehen.

Vorausgesetzt wird in jedem Fall ein Grundkurs in Statistik oder mathematischer Stochastik, wie er eigentlich in allen natur- und sozialwissenschaftlichen Disziplinen im Grundstudium angeboten wird. Die benötigten Begriffe und Resultate werden zwar alle noch einmal zusammengestellt und knapp erläutert, nicht jedoch in der Form, wie es für ein Lehrbuch der Statistik angebracht wäre.

In der Bereitstellung und Verfügbarkeit einer großen Vielfalt von statistischen Verfahren durch Statistik-Software-Systeme, die in immer kürzeren Abständen um neue Module mit immer komplexeren Methoden bereichert wird, liegt auch eine gewisse Gefahr. Das Wissen des Anwenders um die Hintergründe dieser Verfahren hält oft nicht

Schritt mit dieser rasanten Entwicklung. Ein solches Hintergrundwissen erscheint unserer Meinung nach, zumindest zu einem gewissen Grad, auch für den Anwender erforderlich zu sein, damit er das seinem Problem angemessene Modell auswählen kann, die Modellvoraussetzungen versteht, aus den Resultaten der Rechnung die richtigen Schlüsse ziehen kann und nicht zu Fehlinterpretationen geführt wird.

Im vorliegenden Buch werden deswegen zu allen statistischen Verfahren die Modelle erläutert und die Voraussetzungen zur Anwendung des jeweiligen Verfahrens genannt. Dann wird, zumeist anhand eines Beispiels, die Durchführung mit Hilfe von SAS beschrieben durch Angabe des Programm-Textes und dessen Erläuterung. Ein solches Programm führt zu graphischen Darstellungen und/oder zu Ausgabedateien, die im Text kurz Output genannt werden. Daher schließt sich an die Durchführung mit Hilfe von SAS eine ausführliche Erläuterung und Interpretation des Output an.

Der Leser soll dadurch in die Lage versetzt werden, sein statistisches Problem mit Hilfe von SAS zu lösen, weitgehend ohne auf die für den Anfänger abschreckend umfangreichen SAS-Handbücher zurückgreifen zu müssen. Das Buch geht insbesondere auf die Anwendung von SAS auf dem PC ein. Hierzu sind Grundkenntnisse des Betriebssystems DOS von Vorteil. Obwohl sich das Buch auf die PC-Version von SAS bezieht, ist es mit wenigen Einschränkungen auch für den Benutzer der Großrechner-Version geeignet.

Das vorliegende Buch ist weder ein Lehrbuch der Statistik noch eine systematische Einführung in SAS. Schwerpunkt der Darstellungen sind die Konzepte der Statistik, SAS dient als Werkzeug zur Realisierung dieser Konzepte. Daher ist das Buch auch gegliedert nach methodischen Gesichtspunkten der Statistik. SAS wird nur soweit vorgestellt, wie es zur Umsetzung der einzelnen statistischen Methoden notwendig ist. Deshalb kann dieses Buch auch kein Ersatz für die äußerst umfangreichen SAS-Handbücher sein, die immer dann herangezogen werden sollten, wenn man zusätzliche Möglichkeiten ausschöpfen möchte. Auf einige dieser zusätzlichen Möglichkeiten wird im Text durch Verweise auf die entsprechenden SAS-Dokumentationen hingewiesen.

Nach einer Einführung in SAS in den ersten beiden Kapiteln wird die beschreibende Statistik in Kapitel 3 an Hand einer Reihe von Beispielen behandelt. In Kapitel 4 werden die Grundlagen der Wahrscheinlichkeitstheorie und Statistik in knapper Form zusammengestellt und soweit

beschrieben, wie es im weiteren benötigt wird. In Kapitel 5 werden einige grundlegende Verfahren der Statistik vorgestellt. Dazu zählen neben den Ein- und Zweistichprobentests unter Normalverteilungsannahme auch Anpassungstests und die nichtparametrischen bzw. verteilungsfreien Verfahren. Die letzten beiden Kapitel 6 und 7 beinhalten eine Reihe von Verfahren der Varianzanalyse und Regressionsrechnung, die unter dem Begriff lineare Modelle zusammengefaßt werden können. Darin werden auch einige Probleme angeschnitten, die mehr den fortgeschrittenen Statistiker ansprechen, wie z.B. spezielle Randomisationsstrukturen, unbalancierte Daten bei Mehrfachklassifikation, Kovarianzanalyse. Das abschließende Literaturverzeichnis haben wir zur besseren Orientierung um einige Hinweise zu Lehrbüchern und weiterführender Literatur ergänzt.

Bei der erforderlichen Auswahl der Themen haben wir uns von dem Prinzip leiten lassen, einerseits möglichst einfache und grundlegende Verfahren der Statistik vorzustellen und andererseits einige komplexere Methoden zu behandeln, die unserer Erfahrung nach häufig in der Praxis verwandt werden. Gerade in diesem letzten Punkt stützt sich die subjektive Auswahl auf unsere mehrjährige Beratungspraxis und die Zusammenarbeit mit "Anwendern" in Hohenheim. Natürlich konnten dabei einige für die Anwendung interessante Gebiete, wie z.B. multivariate Methoden und Zeitreihenanalyse, nicht in dieses Buch aufgenommen werden.

Wir haben uns bemüht, Computer-Englisch und Abkürzungen weitgehend zu vermeiden. Allerdings erschien es uns sinnvoll, einige Wörter wie z. B. Output im Text wie ein deutsches Wort zu verwenden, da eine direkte Übersetzung, etwa Ausstoß, umständlich und sinnentstellend erscheint. Zu den übernommenen Anglizismen zählt auch, daß im gesamten Text ein Dezimalpunkt statt des im Deutschen üblichen Kommas verwendet wird. Von SAS reservierte Schlüsselwörter (DATA, PROC, UNIVARIATE,...) werden in Großbuchstaben wiedergegeben. Programmtexte und Ausgabedateien sind durch einen Rahmen hervorgehoben. Da die Ausgabedateien der einheitlichen Darstellung wegen ebenfalls in Proportionalchrift gesetzt wurden, können kleine Abweichungen in der Form gegenüber der Bildschirmausgabe auftreten. Disketten mit allen Beispiel-Programmtexten können von uns gegen eine Schutzgebühr bezogen werden.

Schließlich ist es uns eine angenehme Pflicht denen zu danken, die am Zustandekommen dieses Buches beteiligt waren. Dazu zählen eine Reihe von Studenten und uns verbundene Kollegen, die durch fortwährende Diskussionen und Anregungen direkt oder indirekt an der Gestaltung des Buches mitgewirkt haben. Unser Dank gilt Herrn Heinz Becker, der bei der Überprüfung der Programmtexte behilflich war. Ganz herzlich möchten wir uns auch bei unserer EXPertin Frau Regina Schulze bedanken, die uns bei der Erstellung des Textes im Textverarbeitungssystem EXP unterstützt hat. Gerne erwähnen wir auch dankend die angenehme Zusammenarbeit mit Herrn Dr. Spuhler vom Teubner Verlag.

Den Benutzern dieses Buches empfehlen wir, die Beispiele auch als Übungsaufgaben anzusehen und diese durch Variieren, Umstellen und Ergänzungen zu einer eigenen kleinen Programmsammlung auszubauen. Dabei wünschen wir viel Erfolg und möglichst wenige rote Fehlermeldungen.

Stuttgart-Hohenheim, im Sommer 1992

Julius Dufner, Uwe Jensen, Erich Schumacher

Inhaltsverzeichnis

Kapitel 1 SAS für Personal Computer

1.1	Das SAS-Softwaresystem	13
1.2	SAS auf dem PC	14
1.3	Der interaktive Display-Manager-Modus	16
1.3.1	Starten von SAS.....	16
1.3.2	Die primären Fenster	16
1.3.3	Display-Manager-Kommandos.....	17
1.3.4	Sekundäre Fenster	19
1.4	Der nichtinteraktive Modus	21

Kapitel 2 Das SAS-Programmsystem

2.1	Ein einführendes Beispiel	22
2.1.1	DATA step und PROC step.....	23
2.1.2	SAS-Programm	24
2.1.3	Realisierung im Display-Manager-Modus.....	25
2.2	Ergänzungen	31
2.2.1	SAS-Programm	32
2.2.2	Realisierung im Display-Manager-Modus.....	33
2.2.3	Regeln zur Programmgestaltung.....	36
2.3	Externe Daten	37
2.3.1	ASCII-Dateien	37
2.3.2	DOS-Dateien anderer Softwaresysteme.....	38
2.3.3	Transfer PC - Großrechner.....	39
2.4	Die Programmiersprache SAS	39
2.4.1	SAS-Anweisungen	40
2.4.2	SAS-Programme	40
2.4.3	Beschreibung der benutzten Anweisungen	41
2.4.3.1	DATA step	41
2.4.3.2	PROC step.....	45
2.4.3.3	Anweisungen an beliebiger Stelle eines SAS-Programms.....	46

Kapitel 3 Beschreibende Statistik

3.1	Eindimensionale Stichproben	49
3.1.1	Graphische Darstellungen	50
3.1.1.1	Histogramme.....	50
3.1.1.2	Ausgabe von SAS-Graphiken	54
3.1.1.3	Stabdiagramme	56
3.1.1.4	Kreisdiagramme	60
3.1.2	Statistische Maßzahlen.....	61
3.1.2.1	Lagemaße.....	62
3.1.2.2	Streuungsmaße	63
3.1.2.3	Formmaße.....	63
3.1.2.4	Statistische Maßzahlen mittels SAS	65
3.2	Zwei- und mehrdimensionale Stichproben	68
3.2.1	Punktendiagramme	68
3.2.2	Zusammenhangsmaße.....	70
3.2.3	Anpassung von Regressionsfunktionen	74
3.2.3.1	Prinzip der kleinsten Quadrate	74
3.2.3.2	Lineare Anpassung	77
3.2.3.3	Nichtlineare Anpassung	86
3.2.3.4	Ergänzungen zum DATA step	102

Kapitel 4 Grundlagen der Wahrscheinlichkeitstheorie und Statistik

4.1	Wahrscheinlichkeitstheorie	105
4.1.1	Ereignisse, Stichprobenraum.....	106
4.1.2	Wahrscheinlichkeiten.....	106
4.1.3	Zufallsvariable.....	107
4.1.4	Einige spezielle Wahrscheinlichkeitsverteilungen	112
4.1.4.1	Diskrete Verteilungen	112
4.1.4.2	Stetige Verteilungen	115
4.1.5	Grenzwertsätze.....	119
4.1.6	Testverteilungen.....	121
4.1.6.1	Die Chi-Quadrat (χ^2)-Verteilung	121
4.1.6.2	Die Student'sche t-Verteilung	122
4.1.6.3	Die F(isher)-Verteilung	123

4.2	Grundlagen der beurteilenden Statistik.....	124
4.2.1	Parameterschätzung	124
4.2.1.1	Punktschätzungen	124
4.2.1.2	Intervallschätzungen - Vertrauensintervalle	128
4.2.2	Tests	129

Kapitel 5 Beurteilende Statistik - Grundlegende Verfahren

5.1	Tests bei Normalverteilungsannahme.....	132
5.1.1	Einstichproben-Tests	132
5.1.1.1	Test des Erwartungswertes – Einstichproben t-Test	132
5.1.1.2	Test der Varianz	138
5.1.2	Zweistichproben-Tests	141
5.1.2.1	Vergleich verbundener (gepaarter) Stichproben	141
5.1.2.2	Vergleich unabhängiger Stichproben – Der t-Test.....	141
5.2	Anpassungstests	148
5.2.1	Übersicht über einige Anpassungstests	148
5.2.2	Der Shapiro-Wilk Test.....	155
5.3	Verteilungsfreie Verfahren - Nichtparametrische Methoden	159
5.3.1	Einstichproben-Tests	159
5.3.1.1	Der Binomialtest	159
5.3.1.2	Test auf Zufälligkeit	162
5.3.2	Zwei- und k-Stichprobentests.....	165
5.3.2.1	Vergleich zweier verbundener Stichproben	165
5.3.2.2	Vergleich zweier unverbundener Stichproben.....	169
5.3.2.3	Vergleich mehrerer unabhängiger Stichproben - Der Kruskal-Wallis Test.....	173
5.3.2.4	Vergleich mehrerer verbundener Stichproben - Der Friedman Test.....	176
5.3.3	Kontingenztafeln – Unabhängigkeits- und Homogenitätstests	179
5.3.3.1	Der Unabhängigkeitstest	180
5.3.3.2	Der exakte Test von Fisher	184
5.3.3.3	Der Homogenitätstest.....	188

Kapitel 6 Varianzanalyse

6.1	Einfaktorielle Varianzanalyse - fixe Effekte	191
6.1.1	Varianzanalysemodell und F-Test.....	192
6.1.2	Gütefunktion und Wahl des Stichprobenumfangs	196
6.1.3	Durchführung in SAS – Beispiel 6_1	198
6.1.4	Abweichungen von den Modellvoraussetzungen.....	201
6.1.5	Überprüfung von Modellvoraussetzungen	203
6.1.5.1	Test der Normalverteilungsannahme.....	203
6.1.5.2	Der modifizierte Levene-Test	205
6.1.6	Überparametrisierung des Modells	208
6.2	Multiple Mittelwertvergleiche	209
6.2.1	Schätzung der Modellparameter.....	210
6.2.2	Vertrauensintervall und Test für eine Paardifferenz	211
6.2.3	Multiple Tests und simultane Vertrauensintervalle	212
6.2.3.1	Bonferroni- und Sidak-Test.....	212
6.2.3.2	Scheffe-Test	213
6.2.3.3	Tukey-Test	214
6.2.3.4	Dunnnett-Test für Vergleiche mit einer Kontrolle	215
6.2.4	Sidak- , Scheffe-Tests und lineare Kontraste in SAS.....	216
6.2.4.1	Sidak- und Scheffe- Tests in SAS	216
6.2.4.2	Lineare Kontraste in SAS.....	218
6.2.5	Wachstumsversuch, Tukey- und Dunnnett-Tests in SAS.....	220
6.2.5.1	Vollständig zufällige Zuteilung mittels PROC PLAN.....	221
6.2.5.2	Auswertung in SAS	222
6.2.6	Vergleich simultaner Testprozeduren	227
6.2.6.1	Die Tests nach Bonferroni, Sidak, Scheffe, Tukey	227
6.2.6.2	Lineare Kontraste.....	228
6.2.6.3	Sequentielle Testprozeduren.....	229
6.2.6.4	Zusammenfassung	232
6.3	Einfaktorielle Varianzanalyse - zufällige Effekte	232
6.4	Zweifaktorielle Varianzanalyse - Kreuzklassifikation	235
6.4.1	Zweifaktorielle Varianzanalyse, fixe Effekte	236
6.4.1.1	Modell, F-Tests und paarweise Vergleiche.....	237
6.4.1.2	Durchführung in SAS – Beispiel 6_4	240
6.4.2	Zweifaktorielle Varianzanalyse, zufällige Effekte	244
6.4.2.1	Modell und F-Tests	244
6.4.2.2	Durchführung in SAS	246

6.4.3	Zweifaktorielles gemischtes Modell.....	248
6.4.3.1	Gemischtes Modell und F-Tests.....	249
6.4.3.2	Durchführung in SAS	250
6.4.4	Eine Beobachtung pro Zelle.....	251
6.4.4.1	Modell und F-Tests	252
6.4.4.2	Durchführung in SAS	254
6.4.5	Höherfaktorielle kreuzklassifizierte Versuche	255
6.4.5.1	3-faktorielle kreuzklassifizierte Varianzanalyse	255
6.4.5.2	Durchführung in SAS	256
6.4.5.3	r-faktorielle kreuzklassifizierte Varianzanalyse	256
6.5	Zweifaktorielle hierarchische Varianzanalyse	257
6.5.1	Modell und F-Tests	258
6.5.2	Durchführung in SAS – Beispiel 6_5.....	260
6.5.2.1	Tests.....	260
6.5.2.2	Schätzung der Varianzkomponenten	263
6.5.3	Höherfaktorielle Modelle.....	264
6.6	Versuchsplanung - spezielle Randomisationsstrukturen	265
6.6.1	Complete Randomized Designs.....	266
6.6.2	Randomisierte vollständige Blockanlagen	266
6.6.2.1	Modell , F-Tests und paarweise Vergleiche	268
6.6.2.2	Durchführung in SAS – Beispiel 6_6.....	269
6.6.2.3	Modell mit zufälligen Blockeffekten	272
6.6.3	2-faktorielle Anlage in Blöcken	272
6.6.4	Split-Plot Anlage in Blöcken	274
6.6.4.1	Modell und F-Tests	274
6.6.4.2	Multiple Vergleiche	277
6.6.4.3	Durchführung in SAS – Beispiel 6_7.....	280
6.7	Unbalancierte Daten	288
6.7.1	Zweifaktorielle Kreuzklassifikation, unbalancierte Daten, keine leeren Zellen.....	289
6.7.1.1	Modell.....	289
6.7.1.2	Einführendes Beispiel und R- Notation	291
6.7.1.3	Typ I- Quadratsummenzerlegung	295
6.7.1.4	Typ II- Quadratsummen	297
6.7.1.5	Typ III- Quadratsummenzerlegung.....	299
6.7.1.6	Durchführung in SAS – Beispiel 6_8.....	301
6.7.2	Paarweise Vergleiche adjustierter Erwartungswerte.....	303
6.7.2.1	Adjustierte Erwartungswerte – LSMeans	303

6.7.2.2	Durchführung in SAS – Beispiel 6_8	305
6.7.3	Modelle mit leeren Zellen und die Typ IV- Zerlegung	307
6.7.3.1	Schätzbare Funktionen und testbare Hypothesen	308
6.7.3.2	Typ IV- Quadratsummen.....	310
6.7.3.3	Typ IV-Zerlegung – Beispiel 6_9.....	310
6.7.3.4	Durchführung in SAS – Beispiel 6_9	314
6.7.4	Auswertung mehrfaktorieller Modelle in SAS.....	319

Kapitel 7 Lineare Regressionsanalyse

7.1	Einfache lineare Regression.....	322
7.1.1	Schätzung der Modellparameter.....	324
7.1.2	Univariate Vertrauensintervalle und Tests	327
7.1.3	Simultane Vertrauensbereiche und Tests.....	328
7.1.4	Durchführung in SAS – Beispiel 7_1	329
7.1.5	Überprüfung der Modellannahmen.....	335
7.1.6	Ergänzungen	336
7.1.6.1	Prognose- Intervall für eine Beobachtung.....	336
7.1.6.2	Regression ohne Absolutglied	337
7.2	Multiple lineare Regressionsanalyse	340
7.2.1	Schätzung der Modellparameter.....	341
7.2.2	Univariate Vertrauensintervalle und Tests	344
7.2.3	Simultane Vertrauensbereiche und Tests.....	345
7.2.4	Überprüfung der Modellannahmen.....	347
7.2.5	Durchführung in SAS – Beispiel 7_2.....	348
7.2.6	Techniken zur Modellauswahl.....	354
7.3	Kovarianzanalyse	357
7.3.1	Einfache Kovarianzanalyse	357
7.3.1.1	Schätzung der Modellparameter.....	359
7.3.1.2	Tests und paarweise Vergleiche	361
7.3.1.3	Durchführung in SAS – Beispiel 7_3.....	363
7.3.1.4	Überprüfung von Modellannahmen.....	368
7.3.2	Erweiterungen des Kovarianzanalysemodells	371
Anhang.....		372
Literaturverzeichnis.....		383
Sachverzeichnis.....		392