

Bernward Tewes
unter Mitarbeit von
Hans-Joachim Mittag und
Hans-Georg Sonnenberg

Einführung in SPSS

mit Ausblicken auf die freie Statistiksoftware R

**kultur- und
sozialwissenschaften**

Das Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere das Recht der Vervielfältigung und Verbreitung sowie der Übersetzung und des Nachdrucks, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Kein Teil des Werkes darf in irgendeiner Form (Druck, Fotokopie, Mikrofilm oder ein anderes Verfahren) ohne schriftliche Genehmigung der FernUniversität reproduziert oder unter Verwendung elektronischer Systeme verarbeitet, vervielfältigt oder verbreitet werden.

Inhaltsverzeichnis

1	Einleitung	5
2	Datenerfassung und Datenaufbereitung mit SPSS	8
2.1	SPSS-Fenster Daten-Editor	8
2.1.1	Variablenansicht	9
2.1.2	Datenansicht	14
2.2	Weitere SPSS-Fenster	16
2.3	Datenaufbereitung	18
2.3.1	Variable berechnen	19
2.3.2	Umkodieren in andere Variablen	21
2.3.3	Fälle auswählen	22
3	Kurzeinführung in die Programmiersprache R	25
3.1	Nutzung der Programmumgebung von R	25
3.2	Editor mit Steuerungsfunktionen: Tinn-R	30
4	Univariate Häufigkeitsanalyse	33
4.1	Häufigkeitsanalyse mit SPSS	33
4.2	Häufigkeitsanalyse mit R	38
5	Kenngößen univariater Verteilungen	41
5.1	Kenngößen univariater Verteilungen mit SPSS	41
5.2	Kenngößen univariater Verteilungen mit R	43
6	Quantile und Boxplots	46
6.1	Quantile und Boxplots mit SPSS	47
6.2	Quantile und Boxplots mit R	51
7	Konzentrationsmessung	53
7.1	Konzentrationsmessung mit R	53
8	Bivariate Häufigkeitsverteilungen und Zusammenhangsmaße für nominalskalierte Daten	56
8.1	Bivariate Häufigkeitsverteilungen mit SPSS	56
8.2	Bivariate Häufigkeitsverteilungen mit R	60
9	Zusammenhangsmessung bei metrisch und ordinal skalierten Merkmalen	62
9.1	Korrelationen mit SPSS	63
9.2	Korrelationen mit R	66

10 Aufgaben Block 1	68
11 Schätzung des Mittelwertes bei Normalverteilung	72
11.1 Mittelwertschätzung mit SPSS	73
11.2 Mittelwertschätzung mit R	74
12 χ^2-Test bei zwei nominalskalierten Merkmalen	76
12.1 χ^2 -Unabhängigkeitstest mit SPSS	77
12.2 χ^2 -Unabhängigkeitstest mit R	78
13 Ein- und Zweistichproben-t-Tests	80
13.1 t-Tests mit SPSS	81
13.2 t-Test mit R	85
14 Lineare Regression	88
14.1 Lineare Regression mit SPSS	89
14.2 Lineare Regression mit R	98
15 Varianzanalyse	102
15.1 Varianzanalyse mit SPSS	104
15.2 Varianzanalyse mit R	109
16 Aufgaben Block 2	112
Literatur	114

1 Einleitung

Die Berechnung einer einzelnen statistischen Kennzahl ist bei überschaubaren Datenmengen wie man sie in den Beispielen des Kurses *Statistik* (33209) findet mit keinem großen Aufwand verbunden und kann problemlos im Kopf oder mit Hilfe kleiner Tabellen durchgeführt werden. Ein Taschenrechner ist dabei recht hilfreich.

Aber solche Beispiele sind aus didaktischen Gründen gewählt worden und spiegeln letztlich nicht die Datensätze wider, auf die Sie in der empirischen Praxis treffen werden. Sie werden i.d.R. mehr als ein Merkmal zu untersuchen haben und mehr als eine Handvoll Beobachtungen/Fälle vorliegen haben. Hier ist dann der Computer gefragt.

Wenn man fit ist im Umgang mit einem Tabellenkalkulationsprogramm wie Microsoft Excel, ist ein solches Programm zumindest für die einfacheren statistischen Berechnungen eine Alternative. Es werden einige Funktionen zur Berechnung statistischer Kennzahlen sowie zur grafischen Darstellung angeboten. Sollten allerdings die Kenntnisse im Umgang mit einem solchen Programm nicht so gut sein und/oder die statistischen Analysen zahlreicher und komplexer werden, so ist der Griff zu einem speziellen Statistiksoftwarepaket sinnvoll.

Programme zur statistischen Analyse von Daten sind etwa ab den 60er Jahren des letzten Jahrhunderts entwickelt worden. Von den ersten Vertretern, die in den USA entstanden und (zunächst) für Großrechner entwickelt wurden, besitzen heute nur noch *SAS* (ursprünglich *Statistical Analysis System*) und *SPSS* (ursprünglich *Statistical Package for the Social Sciences*) eine Marktbedeutung.

SAS und SPSS

Das Programm *SAS* ist damals zunächst zur Analyse von Daten aus der landwirtschaftlichen Forschung an der NC State University, USA, entwickelt worden, das *SAS Institute* wurde 1976 gegründet. Zwar bietet *SAS* auch heute noch eine komplette Umgebung, um statistische Analysen durchzuführen, aber *SAS Institute* sieht sein Produkt eher als umfassendes Instrument für Unternehmensentscheidungen.

Der Ursprung von *SPSS* geht auf Nie, Hull und Bent (Stanford University, USA) zurück, die es ab 1965 entwickelten. 1968 wurde die *SPSS Inc.* gegründet und damit das Produkt kommerziell vertrieben. Mit dem Aufkommen leistungsfähiger Arbeitsplatzrechner sind neben der Großrechnerimplementation eine MS-DOS- und später eine Windows-Version entstanden. Die Windows-Variante bildete dann über viele Jahre den Entwicklungsschwerpunkt. 2008 wurde *SPSS* auf der Basis von Java neu programmiert, was dazu führte, dass nun sowohl für Windows als auch für Mac OS, Linux und weitere Unix-Versionen aktuelle Versionen zur Verfügung stehen. Da man sich ähnlich wie *SAS* auch als Software für betriebswirtschaftliche Entscheidungen im Unternehmen platzieren woll-

te, erhielt das Produkt Anfang 2009 einen neuen Namen: *PASW Statistics* (**P**redictive **A**nalYTics **S**oft**W**are). Jedoch blieb dies nur eine kurze Episode, denn mit der Übernahme der *SPSS Inc.* durch *IBM* erhielt das Produkt den Namen *SPSS* zurück und heißt heute *IBM SPSS Statistics*.

- S bzw. S-Plus In den 70er Jahren wurde in den AT&T Bell Laboratories eine Sprache namens *S* explizit für statistische und grafische Anwendungen entwickelt. Ab 1988 erfolgte dann ein kommerzieller Vertrieb dieser Entwicklung mit grafischer Benutzeroberfläche (GUI, **G**raphical **U**ser **I**nterface) als *S-Plus* durch *Statistical Sciences Inc.*, heute wird dieser von *TIBCO Software Inc.* durchgeführt.
- Statistica und Stata In den 80er Jahren sind mit *STATISTICA* und *Stata* zwei weitere Software-Pakete zur Auswertung von Daten mit statistischen Methoden entstanden. *STATISTICA* wurde von der 1984 gegründeten *StatSoft* entwickelt. Es stellt heute eine übersichtliche Benutzeroberfläche (Graphical User Interface, kurz GUI) zur Verfügung und wird ausschließlich für die Windows-Plattform angeboten. *Stata* gehört der 1985 gegründeten *StataCorp* und versteht sich eher als eine Art „statistisches Betriebssystem“, d.h. es bietet eine einfache Programmiersprache. Zwar ist es heute auch mit einer als GUI zu bezeichnenden Menüstruktur zu bedienen, der Schwerpunkt insbesondere auch von Nutzerseite liegt jedoch offenbar immer noch auf der Kommandozeile.
- Open-Source-Produkt R Während alle bisher angeführten Programme zumindest heute kommerziell vertrieben werden, ist in den 90er Jahren in Neuseeland ein Programm entstanden, das unter der GNU General Public License zur Verfügung gestellt wird und somit kostenlos zur Verfügung steht. Unter dem Namen *R* wurde ab 1992 eine Sprache entwickelt, die man als Dialekt von *S* bezeichnen kann (hoher Kompatibilitätsgrad). Bei der Implementation der Sprache in einer Entwicklungsumgebung sind die Entwickler Robert Gentleman und Ross Ihaka (Universität Auckland) aber andere Wege gegangen. Namensgebend für *R* war der erste Buchstabe der Vornamen der ersten Entwickler. Das Projekt ist schnell gewachsen, seit 1997 gibt es ein internationales Kernentwicklerteam. Als Open-Source-Produkt erhält es zudem über Erweiterungspakete verschiedenster Entwickler zu fast jedem Bereich der Statistik eine passende Prozedur.

Der Abriss zur Entwicklung von Statistiksoftware ist natürlich bei Weitem nicht vollständig. Er beschränkt sich auf mehr oder weniger universell einsetzbare Produkte, die auch heute noch eine Marktbedeutung besitzen. Warum ist nun die Wahl in diesem Kurs auf *SPSS* bzw. *R* als Ergänzung gefallen?

Nun, *SPSS* ist ein umfassendes Statistikpaket mit einer relativ einfach zu bedienenden Oberfläche, ermöglicht aber auch die Nutzung per Syntaxsprache. Ferner bietet es zahlreiche Module zur Erweiterung der schon recht umfangreichen Möglichkeiten der Basiskomponente. Ein neues Modul, *IBM SPSS Statistics Developer*, ermöglicht es nun sogar, *R*-Algorithmen in der *SPSS*-Programmumgebung ausführen zu lassen. Insgesamt ist es dem Programm nicht zu Unrecht gelungen, eine Art Marktführerschaft zu erringen.

R kann man wie *S* als eine statistisch orientierte Programmiersprache mit guten Grafikfunktion ansehen. Der Erfolg von *R* manifestiert sich u.a. auch in der bereits angesprochenen Zugriffsmöglichkeit innerhalb von *SPSS* auf *R*-Routinen sowie in der Verwendung als Softwarebasis für das Statistiklabor der Freien Universität Berlin, das sich als didaktische Statistik-Software versteht (www.statistiklabor.de). Da man mit *R* i.d.R. auch dann noch eine Lösung finden kann, wenn z.B. *SPSS* keine (offensichtliche) Möglichkeit mehr bereitstellt, sei es über die Erweiterungspakete oder selbst programmiert, ist es häufig mehr als nur eine Alternative. Die kostenlose Bereitstellung sowie die Tatsache, dass *R* zunehmend an Boden gewinnt und *R*-Kenntnisse heute häufiger schon als Merkmal in Stellenprofilen gewünscht werden, sind weitere Gründe, hier darauf einzugehen.