

# Statistisches Praktikum mit SPSS 13 für Windows

Reihe Benutzereinführung

Band 26

2007



**Bernhard Baltes-Götz**

# **Statistisches Praktikum mit SPSS 13 für Windows**

Benutzereinführung

Band 26

2007 (Rev. 070601)

Herausgeber:       Universitäts-Rechenzentrum Trier  
                          Universitätsring 15  
                          D-54286 Trier  
                          WWW: <http://www.uni-trier.de/urt/urthome.shtml>  
                          E-Mail: [urt@uni-trier.de](mailto:urt@uni-trier.de)  
                          Tel.: (0651) 201-3417, Fax.: (0651) 3921

Leiter:                Dr. Peter Leinen

Autor:                Bernhard Baltes-Götz (E-Mail: [baltes@uni-trier.de](mailto:baltes@uni-trier.de))

Druck:                Druckerei der Universität Trier

Copyright ©         2007; URT

---

## Vorwort

SPSS (frühere Bedeutung: **S**tatistical **P**ackage for the **S**ocial **S**ciences, jetzige Interpretation: **S**uperior **P**erforming **S**oftware **S**ystems) ist ein weitgehend komplettes und leicht zu bedienendes Statistik-Programmpaket, das in den Geo-, Wirtschafts- und Sozialwissenschaften sehr verbreitet ist und alle wichtigen Computertypen bzw. Betriebssysteme unterstützt (MS-Windows, MacOS, UNIX).

Im vorliegenden Manuskript wird ein Einblick in die statistische Datenanalyse mit der SPSS-Version 13 für Windows vermittelt, wobei großer Wert auf die methodologische Einordnung der beschriebenen EDV-Techniken gelegt wird.

Wesentliche Teile des Manuskripts sind wegen der weitgehend konsistenten Bedienungslogik auch für andere SPSS-Versionen unter MS-Windows oder alternativen Betriebssystemen verwendbar.

Das Manuskript wurde ursprünglich als Begleitlektüre zum Kurs *Statistisches Praktikum mit SPSS für Windows* am Universitäts-Rechenzentrum Trier (URT) erstellt, kann jedoch auch im Selbststudium verwendet werden. Dass dabei die meisten Themen in konkreter Arbeit am Rechner nachvollzogen werden sollten, folgt aus der Kurskonzeption:

## Zielgruppe/Voraussetzungen

- Der Kurs ist konzipiert für Personen, die in wesentlichem Umfang bei Forschungsarbeiten mit SPSS mitwirken wollen, also z.B. im Rahmen einer Diplomarbeit oder Dissertation die Durchführung einer eigenen Studie planen oder bereits begonnen haben. Wer lediglich einfache Teilaufgaben zu erledigen hat (z.B. wenige Auswertungen mit einer bereits vorhandenen und fehlerbereinigten SPSS-Datendatei), der sollte eventuell die 2-stündige SPSS-Kurzeinführung des Rechenzentrums besuchen oder das zugehörige Manuskript lesen. Es ist im Internet ausgehend von der Startseite der Universität Trier (<http://www.uni-trier.de/>) auf folgendem Weg zu finden:

[Weitere Serviceangebote](#) > [EDV-Dokumentationen](#) > [Elektronische Publikationen](#) > [Bedienungsanleitungen zu Statistikprogrammen](#) > [SPSS für Windows](#)

- Im Kurs wird eine methodische Grundausbildung (empirische Forschung, Statistik) vorausgesetzt, wie sie üblicherweise in den Studiengängen empirisch orientierter Fächer vermittelt wird.
- An EDV-Voraussetzungen werden elementare Fertigkeiten im Umgang mit PCs unter MS-Windows erwartet.

## Kursinhalte

- Wir konzentrieren uns darauf, in anderen Veranstaltungen (z.B. zur empirischen Forschung oder Statistik) erlernte Begriffe und Methoden mit dem EDV-Werkzeug SPSS in der Praxis anzuwenden. Zwar werden im Kursverlauf viele methodische Themen in knapper Form behandelt, doch kann damit eher vorhandenes Wissen aufgefrischt als neues erworben werden. Insbesondere kann die Anwendung der vielfältigen statistischen Auswertungsmethoden nur exemplarisch stattfinden. Eine explizite Behandlung ist nur bei wenigen, besonders häufig eingesetzten Verfahren möglich (z.B. Kreuztabellenanalyse).

Zu zahlreichen Auswertungsmethoden bietet das Rechenzentrum Spezialveranstaltungen an, in denen die wesentlichen methodologischen Grundlagen und natürlich die praktische Durchführung mit SPSS erläutert werden. Informationen über das URT-Kursprogramm finden Sie

z.B. auf dem WWW-Server der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend über:

[Rechenzentrum > Schulung/Kurse](#)

Zu den meisten Kursen sind ausführliche Manuskripte entstanden, die Sie auf dem genannten WWW-Server von der Startseite ausgehend folgendermaßen erreichen:

[Weitere Serviceangebote > EDV-Dokumentationen > Elektronische Publikationen](#)

- Im Sinne einer praxisnahen, projektorientierten Ausbildung beschreibt das Manuskript eine vollständige empirische Studie von der ersten Idee über die Kodierung, Erfassung, Kontrolle und Modifikation der Daten bis zur statistischen Auswertung und zur Verwertung der Ergebnisse.
- Zwar werden auch in EDV-handwerklicher Sicht die SPSS-Optionen nicht annähernd vollständig behandelt, doch sollten Sie nach dem Kurs mit den erworbenen Grundkenntnissen unter Verwendung der aufgezeigten Informationsmöglichkeiten selbständig und erfolgreich mit SPSS arbeiten können.

### **Zugriff auf die Dateien zum Kurs**

Leser(innen) im Selbststudium werden in der Regel keine eigene Datenerhebung realisieren, können jedoch alle Projekt-Arbeitsschritte ab der Datenprüfung anhand von Dateien, die auf Servern des Rechenzentrums zur Verfügung stehen, konkret durchführen. Im Internet finden Sie diese Dateien ausgehend von der Startseite der Universität Trier (<http://www.uni-trier.de/>) auf folgendem Weg:

[Weitere Serviceangebote > EDV-Dokumentationen > Elektronische Publikationen > Bedienungsanleitungen zu Statistikprogrammen > SPSS für Windows](#)

Im Campusnetz der Universität Trier sind die Dateien noch bequemer über eine Netzfrequenz zugänglich, nachdem Sie sich bei einem Windows-Rechner mit Einbindung in die Domäne URT angemeldet haben. Führen Sie dort nach

**Start > Ausführen**

den Befehl

**k baltes**

aus, um die Netzfrequenz als Laufwerk **K:** in Ihr Windows-System einzubinden. Anschließend finden Sie die Dateien im Verzeichnis

**K:\SPSS\Statistisches Praktikum mit SPSS für Windows\Manuskript**

---

# Inhaltsverzeichnis

<b>1</b>	<b>Von der Theorie zu den SPSS-Variablen</b>	<b>1</b>
1.1	Statistik und EDV als Hilfsmittel der Forschung	1
1.2	Planung und Durchführung einer empirischen Untersuchung im Überblick	2
1.2.1	Forschungsziele bzw. -hypothesen	2
1.2.2	Untersuchungsplanung	2
1.2.3	Durchführung der Studie (inklusive Datenerhebung)	4
1.2.4	Datenerfassung und -prüfung	4
1.2.5	Datentransformation	5
1.2.6	Statistische Datenanalyse	5
1.3	Beispiel für eine empirische Untersuchung	5
1.3.1	Die allgemeinspsychologische KFA-Hypothese	6
1.3.2	Untersuchungsplanung	6
1.3.3	Eine differentialpsychologische Hypothese	9
1.3.4	Zum Einfluss demographischer Merkmale	10
1.3.5	Zu Übungszwecken erhobene Merkmale	11
1.3.6	Der Fragebogen	11
1.4	Strukturierung und Kodierung der Daten	13
1.4.1	Fälle und Merkmale in SPSS	13
1.4.2	Strukturierung	14
1.4.2.1	Variablen zur Fallidentifikation	14
1.4.2.2	Abgeleitete Variablen gehören nicht in den Kodierplan	15
1.4.2.3	Mehrfachwahlfragen	15
1.4.2.3.1	Vollständige Sets aus dichotomen Variablen	15
1.4.2.3.2	Sparsame Sets aus kategorialen Variablen	16
1.4.2.4	Offene Fragen	17
1.4.3	Kodierung	18
1.4.3.1	Die wichtigsten Variablentypen in SPSS	18
1.4.3.2	Das Problem fehlender Werte	19
1.4.3.2.1	System-Missing (SYSMIS)	19
1.4.3.2.2	Fehlende Werte bei Mehrfachwahl-Fragen und offenen Fragen	20
1.4.3.2.3	Vereinfachung der Erfassung durch Datentransformationstechniken	20
1.4.3.3	Fehlerquellen bei der manuellen Datenerfassung minimieren	22
1.4.3.4	SPSS-Variablennamen	23
1.4.3.5	Kodierplan	24
1.5	Durchführung der Studie (inklusive Datenerhebung)	25
<b>2</b>	<b>Einstieg in SPSS für Windows</b>	<b>26</b>
2.1	SPSS für Windows an der Universität Trier	26
2.2	Programmstart und Benutzeroberfläche	27
2.2.1	SPSS starten	27
2.2.2	Die wichtigsten SPSS-Fenster	27
2.2.3	Was man mit SPSS so alles machen kann	28
2.3	Das Hilfesystem	29
2.3.1	Systematische Informationen	29
2.3.2	Gezielte Suche nach Begriffen	29
2.3.3	Kontextsensitive Hilfe zu den Dialogboxen	30
2.3.4	Lernprogramm	30
2.3.5	Fallstudien	31
2.3.6	Statistik-Assistent	31

<b>2.4</b>	<b>Weitere Informationsquellen</b>	<b>32</b>
2.4.1	Handbücher und Manuskripte	32
2.4.2	SPSS im Internet	33
2.4.3	Benutzerberatung	33
<b>2.5</b>	<b>SPSS für Windows beenden</b>	<b>33</b>
<b>3</b>	<b>Datenerfassung und SPSS-Dateneditor</b>	<b>34</b>
<b>3.1</b>	<b>Methoden zur Datenerfassung</b>	<b>34</b>
3.1.1	Automatisierte Verfahren	34
3.1.1.1	Online-Datenerhebung	34
3.1.1.2	Automatisches Einscannen von schriftlichen Untersuchungsdokumenten	36
3.1.2	Manuelle Verfahren	36
3.1.2.1	Erstellung einer Text-Datendatei mit einem beliebigen Texteditor	37
3.1.2.2	Einsatz eines speziellen Datenerfassungsprogramms	38
<b>3.2</b>	<b>Erfassung mit dem SPSS-Dateneditor</b>	<b>40</b>
3.2.1	Dateneditor und Arbeitsdatei	40
3.2.2	Variablen definieren	41
3.2.2.1	Das Datenfenster-Registerblatt Variablenansicht	41
3.2.2.2	Die SPSS-Variablenattribute	42
3.2.2.3	Variablendefinition durchführen	44
3.2.2.4	Übung	47
3.2.3	Variablen einfügen, löschen oder verschieben	47
3.2.3.1	Variablen einfügen	47
3.2.3.2	Variablen löschen	48
3.2.3.3	Variablen verschieben	48
3.2.4	Attribute auf andere Variablen übertragen	48
3.2.4.1	Variablendeklarationen vervielfältigen	48
3.2.4.2	Alle Attribute einer Variablen übertragen	50
3.2.4.3	Einzelne Attribute einer Variablen übertragen	50
3.2.4.4	Übung	50
3.2.5	Sichern der Arbeitsdatei als SPSS-Datendatei	50
3.2.6	Rohdatendatei, Transformationsprogramm und Fertigdatendatei	51
3.2.7	Dateneingabe	53
3.2.8	Daten korrigieren	54
3.2.8.1	Wert in einer Zelle ändern	54
3.2.8.2	Einen Fall einfügen	54
3.2.8.3	Einen Fall löschen	54
3.2.8.4	Einen Fall verschieben	55
3.2.9	Weitere Möglichkeiten des Dateneditors	55
3.2.10	Übung	55
<b>4</b>	<b>Univariate Verteilungs- und Fehleranalysen</b>	<b>56</b>
<b>4.1</b>	<b>Erfassungsfehler</b>	<b>56</b>
4.1.1	Überprüfung von Gültigkeitsregeln	56
4.1.2	Überprüfung von Einzelwerten	56
<b>4.2</b>	<b>Öffnen einer SPSS-Datendatei</b>	<b>57</b>
<b>4.3</b>	<b>Statistische Auswertungen durchführen: Häufigkeitsanalysen zur Prüfung der Variablen FNR</b>	<b>58</b>
<b>4.4</b>	<b>Arbeiten mit dem Ausgabefenster (Teil I)</b>	<b>60</b>
4.4.1	Arbeiten im Navigationsbereich	61
4.4.1.1	Fokus positionieren	61
4.4.1.2	Ausgabeblöcke bzw. Teilausgaben aus- oder einblenden	61
4.4.1.3	Ausgabeblöcke oder -teile markieren	61
4.4.2	Viewer-Dokumente drucken	61
4.4.3	Ausgaben sichern und öffnen	62
4.4.4	Objekte via Zwischenablage in andere Anwendungen übertragen	62
4.4.5	Übungen	63



---

<b>4.5</b>	<b>Graphische Darstellungen in Statistik-Dialogboxen anfordern: Häufigkeits- bzw. Fehleranalyse für die Variablen GESCHL und FB</b>	<b>63</b>
<b>4.6</b>	<b>Häufigkeits- bzw. Fehleranalysen für die restlichen Projektvariablen</b>	<b>65</b>
4.6.1	Übung	65
4.6.2	Diskussion ausgewählter Ergebnisse	68
<b>4.7</b>	<b>Suche nach Daten</b>	<b>70</b>
<b>4.8</b>	<b>Arbeiten mit dem Ausgabefenster (Teil II)</b>	<b>70</b>
4.8.1	Nachbearbeitung von Tabellen	70
4.8.1.1	Pivot-Editor starten	71
4.8.1.2	Modifikation von Zellinhalten	71
4.8.1.3	Tabellenvorlagen	72
4.8.2	Weitere Gestaltungsmöglichkeiten im Navigationsbereich	73
4.8.2.1	Blöcke bzw. Teilausgaben kopieren, verschieben oder löschen	73
4.8.2.2	Befördern und Degradieren	73
4.8.3	Ausgaben exportieren	73
4.8.4	Mehrere Ausgabefenster verwenden	75
<b>5</b>	<b>Speichern der SPSS-Kommandos zu wichtigen Anweisungsfolgen</b>	<b>76</b>
<b>5.1</b>	<b>Zur Motivation</b>	<b>76</b>
<b>5.2</b>	<b>Dialogunterstützte Erstellung von SPSS-Programmen</b>	<b>77</b>
<b>5.3</b>	<b>Arbeiten mit dem Syntax-Fenster</b>	<b>80</b>
<b>5.4</b>	<b>Elementare Regeln zur SPSS-Syntax</b>	<b>80</b>
<b>6</b>	<b>Datentransformation</b>	<b>82</b>
<b>6.1</b>	<b>Vorbemerkungen</b>	<b>82</b>
6.1.1	Rohdatendatei, Transformationsprogramm und Fertigdatendatei	82
6.1.2	Hinweise zum Thema Datensicherheit	83
6.1.3	Initialisierung neuer numerischer Variablen	84
<b>6.2</b>	<b>Alte Werte einer Variablen auf neue abbilden (Umkodieren)</b>	<b>85</b>
6.2.1	Das praktische Vorgehen am Beispiel einer künstlichen Gruppenbildung	85
6.2.2	Technische Details	87
6.2.3	Übungen	88
<b>6.3</b>	<b>Zur Rolle des EXECUTE-Kommandos</b>	<b>89</b>
<b>6.4</b>	<b>Berechnung von Variablen nach mathematischen Formeln</b>	<b>91</b>
6.4.1	Beispiel	91
6.4.2	Technische Details	93
6.4.2.1	Numerischer Ausdruck	93
6.4.2.1.1	Numerische Funktionen	93
6.4.2.1.2	Regeln für die Bildung numerischer Ausdrücke	95
6.4.2.2	Sonstige Hinweise	96
6.4.3	Übungen	97
<b>6.5</b>	<b>Bedingte Datentransformation</b>	<b>97</b>
6.5.1	Beispiel	98
6.5.2	Bedingungen formulieren	99
6.5.2.1	Vergleich	100
6.5.2.2	Logischer Ausdruck	100
6.5.2.3	Regeln für die Auswertung logischer Ausdrücke	101
6.5.3	Übung	102
<b>6.6</b>	<b>Häufigkeit bestimmter Werte bei einem Fall ermitteln</b>	<b>102</b>
<b>6.7</b>	<b>Erstellung der Fertigdatendatei mit dem Transformationsprogramm</b>	<b>104</b>
6.7.1	Transformationsprogramm vervollständigen	104
6.7.2	Transformationsprogramm ausführen	107

<b>7</b>	<b>Prüfung der zentralen Projekt-Hypothesen</b>	<b>108</b>
7.1	Entscheidungsregeln beim Hypothesentesten	108
7.2	Zu den Voraussetzungen unserer Hypothesentests	112
7.3	Verteilungsanalyse zu AERGF, AERGFAM und LOT	114
7.3.1	Diagnose von Ausreißern	114
7.3.2	Die SPSS-Prozedur zur explorativen Datenanalyse	115
7.3.3	Ergebnisse für AERGF	116
7.3.4	Ergebnisse für AERGFAM und LOT	119
7.4	Prüfung der differentialpsychologischen Hypothese	120
7.4.1	Regression von AERGFAM auf LOT	120
7.4.2	Methodologische Anmerkungen	121
7.4.2.1	Explorative Analysen im Anschluss an einen „gescheiterten“ Hypothesentest	121
7.4.2.2	Post hoc - Poweranalyse	122
7.4.2.3	Paarweiser oder fallweiser Ausschluss fehlender Werte	124
7.5	Prüfung der KFA-Hypothese	125
7.6	Übung	127
7.7	Arbeiten mit dem Ausgabefenster (Teil III)	127
7.7.1	Pivot-Editor starten	127
7.7.2	Dimensionen verschieben	128
7.7.3	Gruppierungen	128
7.7.4	Kategorien aus- und einblenden	130
<b>8</b>	<b>Grafische Datenanalyse</b>	<b>132</b>
8.1	Streudiagramm anfordern	132
8.2	Streudiagramm modifizieren	134
8.3	Grafiken verwenden	138
8.4	Übung	138
<b>9</b>	<b>Fälle auswählen</b>	<b>140</b>
9.1	Auswahl über eine Bedingung	140
9.2	Bericht anfordern	141
<b>10</b>	<b>Analyse von Kreuztabellen</b>	<b>143</b>
10.1	Beschreibung der bivariaten Häufigkeitsverteilung	144
10.2	Die Unabhängigkeits- bzw. Homogenitätshypothese	149
10.3	Testverfahren	150
10.3.1	Asymptotische $\chi^2$ - Tests	150
10.3.2	Exakte Tests	153
10.3.3	Besonderheiten bei $(2 \times 2)$ -Tabellen	155
10.3.3.1	Ein klarer Fall für Fischers Test	155
10.3.3.2	Einseitige Hypothesen	155
10.3.3.3	Kontinuitätskorrektur nach Yates	156
<b>11</b>	<b>Fälle gewichten</b>	<b>157</b>
11.1	Beispiel	157
11.2	Übung	158

---

<b>12</b>	<b>Auswertung von Mehrfachwahlfragen</b>	<b>158</b>
12.1	Häufigkeitstabellen	159
12.2	Kreuztabellen	161
12.3	Ein sparsames Set kategorialer Variablen expandieren	164
<b>13</b>	<b>Datendateien im Textformat einlesen</b>	<b>166</b>
13.1	Import von positionierten Textdaten (feste Breite)	166
13.2	Import von separierten Daten Textdaten	172
13.3	Überprüfung der revidierten differentialpsychologischen Hypothese	174
<b>14</b>	<b>Einstellungen modifizieren</b>	<b>176</b>
<b>15</b>	<b>Anhang</b>	<b>178</b>
15.1	<b>Weitere Hinweise zur SPSS-Kommandosprache</b>	<b>178</b>
15.1.1	Hilfsmittel für das Arbeiten mit der SPSS-Kommandosprache	178
15.1.2	Interpretation von Syntaxdiagrammen	178
15.1.3	Aufbau von SPSS-Programmen	179
15.1.4	Aufbau eines einzelnen SPSS-Kommandos	180
15.1.5	Regeln für Variablenlisten	182
15.1.5.1	Abkürzende Spezifikation einer Serie von Variablen	182
15.1.5.2	Der Platzhalter varlist	182



---

# 1 Von der Theorie zu den SPSS-Variablen

## 1.1 Statistik und EDV als Hilfsmittel der Forschung

Die Erfahrungswissenschaften bemühen sich um allgemeingültige Aussagen deskriptiver, explanatorischer oder prognostischer Art. In vielen Anwendungsbereichen sind dabei *deterministische* Gesetze (z.B. Ohmsches Gesetz der Elektrizität, Hebelgesetz der Mechanik) kaum zu finden, und man muss sich auf die Untersuchung *probabilistischer* Gesetze beschränken.

Beispiel: Welchen Effekt hat das Rauchen auf die Entstehung von Lungenkrebs?

Wie wir wissen, hat das (aktive oder passive) Rauchen auch bei vergleichbarer Dosierung der Schadstoffe keinesfalls für alle Personen dieselben Folgen.

In einer solchen Situation können statistische Methoden dabei helfen, rationale Entscheidungen zu treffen, denn:

"Statistics is a body of methods for making wise decisions in the face of uncertainty"  
(Wallis & Roberts, 1956, S. 1).

Die statistischen Methoden zur Entscheidungshilfe lassen sich in zwei Gruppen einteilen:

- **Deskriptive Statistik**  
Sie dient zur Darstellung und Zusammenfassung von *Stichprobendaten*. Hier kann man auch die *explorativen* Verfahren einordnen, deren Popularität in den letzten Jahren deutlich zugenommen hat.
- **Inferenzstatistik (schlussfolgernde Statistik)**  
Hier geht es darum, aus Stichprobendaten Informationen über die zugrunde liegende *Population* zu gewinnen, wobei folgende Aufgaben zu unterscheiden sind:
  - **Parameterschätzung**  
Beispiel: Wie hoch ist bei Rauchern das Risiko, an Lungenkrebs zu erkranken?  
Hier ist eine Wahrscheinlichkeit zu schätzen.  
Neben den Punktschätzungen sind die Intervallschätzungen von großer Bedeutung. Zu einer gewünschten Sicherheit (z.B. 95%) erhält man aus den Stichprobendaten ein Vertrauensintervall, das den fraglichen Populationsparameter mit der festgelegten Wahrscheinlichkeit enthält.
  - **Hypothesentests** (konfirmatorische Verfahren)  
Beispiel: Ist bei Rauchern das Lungenkrebsrisiko größer als bei Nichtrauchern?  
Hier ist eine Entscheidung zwischen zwei Hypothesen zu treffen:
    - *Nullhypothese*  
Im Beispiel: Das Lungenkrebsrisiko ist bei Rauchern *nicht* größer als bei Nichtrauchern.
    - *Alternativhypothese*  
Im Beispiel: Das Lungenkrebsrisiko ist bei Rauchern erhöht.

Die in den Beispielen zur Inferenzstatistik genannten Fragen sind anhand weniger, unrepräsentativer Einzelbeobachtungen (z.B. der steinalte Kettenraucher) nicht zu klären. Solche Anekdoten lassen keine sinnvollen Schlüsse und Entscheidungen zu, sondern demonstrieren lediglich die in obigem Zitat angesprochene Unsicherheit.

Eine grundlegende Strategie der statistisch arbeitenden Forschung, trotz Unsicherheit zu guten Entscheidungen zu kommen, besteht darin, zu einer Fragestellung hinreichend viele **unabhängige** Beobachtungen zu machen und diese mit statistischen Verfahren zu analysieren. Zur Klärung der Raucherproblematik wird man vielleicht bei ca. 500 Personen (= **Beobachtungseinheiten**,

**Merkmalsträgern, Fällen)** die **Merkmale** *Nikotinkonsum* und *Lungenkrebs-Erkrankungen* beobachten. Da außerdem eine Beteiligung weiterer Bedingungen an der Lungenkrebs-Entstehung anzunehmen ist, wird man in einer wohldurchdachten Studie noch viele zusätzliche Merkmale erheben (z.B. Alter, Geschlecht, Beruf, Schadstoffbelastung des Wohnortes).

Eine praktikable Auswertung solcher Datenmengen ist aber nur mit EDV-Hilfe möglich. Mit **SPSS für Windows** steht ein bequemes, leistungsfähiges und sehr bewährtes Analysesystem für die statistische Forschung zur Verfügung. Es bietet fast alle wichtigen statistischen Verfahren sowie gute graphische Darstellungsmöglichkeiten und unterstützt alle in der Windows-Welt gebräuchlichen Verfahren zur Kooperation mit anderen Programmen (z.B. Zwischenablage, ODBC).

## **1.2 Planung und Durchführung einer empirischen Untersuchung im Überblick**

Zunächst wollen wir uns einen Überblick über die verschiedenen Phasen eines empirischen Forschungsprojektes und damit auch über unser Kursprogramm verschaffen. Dabei werden zahlreiche Aufgaben, Methoden und Probleme angesprochen, über die Sie sich im Bedarfsfall in den Lehrveranstaltungen oder in der Literatur zur empirischen Forschung informieren können (siehe z.B. Bortz & Döring 1995; Pedhazur & Pedhazur Schmelkin 1991; Schnell, Hill & Esser 2005).

Die anschließende Darstellung soll als Übersicht dienen und ist daher relativ knapp gehalten. Ihr folgt unmittelbar die konkrete und ausführliche Anwendung auf unsere Beispielstudie.

Weil die dargestellten Aufgaben teilweise interdependent sind, bilden sie keine strenge, bei allen empirischen Studien gleichförmig ablaufende Sequenz.

### **1.2.1 Forschungsziele bzw. -hypothesen**

Einer empirischen Untersuchung wird in der Regel eine längere Phase der intensiven theoretischen Auseinandersetzung mit dem Thema vorangehen. Daraus ergeben sich Forschungsinteressen, die - u.a. in Abhängigkeit vom Forschungsstand - eher von explorativer (hypothesensuchender) oder eher von konfirmatorischer (hypothesenprüfender) Natur sind. Oft werden *beide* Forschungsstrategien vertreten sein. Die zu prüfenden Hypothesen sollten wegen ihrer Steuerungsfunktion für spätere Schritte möglichst exakt formuliert werden.

### **1.2.2 Untersuchungsplanung**

Wenn Sie eine Theorie bzw. eine Hypothesenfamilie empirisch prüfen oder einen Gegenstandsbereich empirisch explorieren möchten, haben Sie bei der Untersuchungsplanung zahlreiche Aufgaben zu lösen:

- **Festlegung der Beobachtungseinheit(en) und der zu untersuchenden Merkmale**

In der Regel ergibt sich aus der Fragestellung unmittelbar, welche Beobachtungseinheiten (Merkmalsträger) einer Studie zugrunde liegen sollten (z.B. Personen, Volkswirtschaften, Orte, Betriebe, Bodenproben, Jahre), und welche Merkmale bei jeder Beobachtungseinheit festgestellt werden sollten.

Gelegentlich bieten sich hierarchisch geschachtelte Untersuchungseinheiten auf mehreren Ebenen an (siehe z.B. Raudenbush & Bryk 2002). So hat man es etwa bei einer Studie zur Arbeitszufriedenheit und Produktivität von Arbeitnehmern aus verschiedenen Firmen in Abhängigkeit von Person- und Organisationsmerkmalen mit Beobachtungseinheiten auf zwei Ebenen zu tun:

- Arbeitnehmer
- Firmen

Bei der Untersuchung von Schülern aus verschiedenen Klassen, die wiederum zu Schulen gehören, kommt sogar eine Hierarchie mit drei Ebenen in Frage.

Bei der späteren inferenzstatistischen Auswertung ist zu beachten, dass die meisten Verfahren unabhängige *Residuen* voraussetzen. Die bei einer hierarchischen Datenstruktur auf der untersten Ebene naturgemäß anzutreffende Abhängigkeit der *Beobachtungen* muss in den Auswertungsverfahren geeignet modelliert werden.

Das Demonstrationsprojekt dieses Einstiegskurses kommt allerdings mit einer konventionellen, flachen Datenstruktur aus, und die Behandlung der speziellen Optionen und Probleme der Mehrebenenanalyse bleibt einem speziellen Kurs vorbehalten.

- **Entscheidung für ein Untersuchungsdesign**

Sie können z.B. einen (quasi-)experimentellen Untersuchungsplan entwerfen oder eine reine Beobachtungsstudie wählen, die quer- oder längsschnittlich angelegt sein kann.

- **Operationalisierung der zu untersuchenden Merkmale**

Sie werden bestrebt sein, zur Operationalisierung von theoretischen Begriffen (z.B. sozioökonomischer Status, Ärger, Optimismus) reliable und valide Messmethoden zu wählen bzw. zu entwerfen, die außerdem nicht zu aufwändig sind. Das Skalenniveau der gewählten Messmethoden muss die Voraussetzungen der geplanten Auswertungsverfahren erfüllen.

Wenn Sie das Glück haben, gut messbare quantitative Merkmale untersuchen zu können (z.B. Alter), dann sollten Sie die verfügbare Information *nicht* durch eine *künstliche* und *willkürliche* Klassenbildung reduzieren (z.B. durch Bildung der Altersklassen < 20, 21-40, 41-60, > 60). Häufig sind Modelle für metrische Daten einfacher und erfolgreicher als solche für vergrößerte Daten. Vor allem können Sie mit SPSS zu einer metrischen Variablen nach Belieben klassifizierte Varianten erzeugen, wenn Sie dies für spezielle Analysen wünschen. Eine Ausnahme von dieser Regel ist vielleicht bei der Befragung von Personen nach ihrem Einkommen zu machen. Um bei dieser sensiblen Frage Widerstände zu vermeiden, muss man sich eventuell auf die Erhebung von groben Einkommensklassen beschränken.

Bei den Überlegungen zur Operationalisierung spielen auch die verfügbaren technischen Hilfsmittel für die Datenerhebung und -erfassung eine Rolle. Mit Hilfe der Computertechnik ist eine interaktive, individualisierte und dabei auch noch ökonomische Datenerfassung zu realisieren. Bei besonderen Ansprüchen (z.B. zeitgenaue Steuerung experimenteller Abläufe) kommen spezielle Rechner im Forschungslabor zum Einsatz. Für eine kontinuierliche, alltagsbegleitende Datenerfassung können oft Rechner im Taschenformat (z.B. PDAs) genutzt werden. Einfache Befragungen werden mittlerweile routinemäßig via Internet realisiert, wenn die Zielgruppe auf diesem Weg erreichbar ist.

- **Empirisch prüfbare Hypothesen formulieren**

Aus einer in theoretischen Begriffen formulierten Hypothese ergibt sich im Verlauf der Untersuchungsplanung durch zahlreiche Konkretisierungen und Operationalisierungen eine in empirischen Begriffen formulierte und damit statistisch prüfbare Hypothese, die möglichst exakt notiert werden sollte. Dabei muss z.B. klar erkennbar sein, ob eine *gerichtete* oder eine *ungerichtete* Hypothese vorliegt.

- **Statistische Versuchsplanung**

Für jede Hypothese ist ein **statistisches Entscheidungsverfahren** zu wählen, dessen Voraussetzungen an Skalenniveau und Verteilungsverhalten der beteiligten Merkmale (voraussichtlich) erfüllt sind.

Zu jedem geplanten Test ist das **Fehlerrisiko erster Art** ( $\alpha$ -Fehler) festzulegen, wobei z.B. die übliche 5%-Konvention übernommen werden kann.

Es ist zu überlegen, wie eine repräsentative und zur Durchführung der geplanten Auswertungsverfahren hinreichend große **Stichprobe** rekrutiert werden kann. Bei ausgeprägt konfirmatorisch angelegten Studien ist bei der Stichprobenumfangsplanung insbesondere das **Fehlerrisiko zweiter Art** (der  $\beta$ -Fehler) zu berücksichtigen.<sup>1</sup>

- **Strukturierung und Kodierung der Daten**

Wer ganz sicher gehen will, dass die bei einer Studie erhobenen Informationen sicher und bequem in die EDV übernommen werden können, sollte die Daten schon in der Planungsphase gegenüber der zuständigen Software deklarieren. Beim Entwurf eines Online-Formulars oder bei der Vorbereitung einer Datenerfassung per Scanner geschieht die Datendeklaration gegenüber der verwendeten Software, also auf jeden Fall in der Planungsphase. Diese Software kann in der Regel die erfassten Merkmale als SPSS-Datendatei exportieren, so dass keine erneute Datendeklaration gegenüber SPSS erforderlich ist.

Häufig werden die Daten mit schriftlichen Untersuchungsdokumenten erhoben und anschließend manuell erfasst (z.B. mit dem SPSS-Dateneditor). Man sollte auch bei diesem Vorgehen die Daten schon vor der Erhebung deklarieren (z.B. gegenüber SPSS). Anfänger werden bei der Arbeit mit einem Computer-Programm, das die vorwiegend forschungslogisch und kaum durch EDV-Restriktionen diktierte Datenstruktur explizit einfordert, konzeptionelle Probleme eventuell eher entdecken als bei der schriftlichen Beschreibung ihres Forschungsvorhabens.

Bei einer flachen Datenstruktur (ohne geschachtelte Beobachtungseinheiten, siehe oben) sind oft nur Kodierungsregeln festzulegen. Hierunter fällt z.B. die Vereinbarung, dass beim Merkmal Geschlecht die Ausprägung *weiblich* durch eine Eins und die Ausprägung *männlich* durch eine Zwei erfasst werden soll. Bei einer hierarchischen Datenstruktur (z.B. mit Firmen und Mitarbeitern als geschachtelten Beobachtungseinheiten) werden meist die Beobachtungseinheiten der untersten Ebene zu den Fällen (bzw. Zeilen) der Datenmatrix.

Die Festlegungen zur Strukturierung und Kodierung der Projektdaten sollten in einem **Kodierplan** dokumentiert werden. Er ist bei einer manuellen Datenerfassung als genaue Arbeitsvorschrift unverzichtbar und eignet sich generell zur Dokumentation der Daten (eventuell für einen größeren Nutzerkreis).

Wir werden uns in Abschnitt 1.4 mit der Strukturierung und Kodierung von Daten ausführlich beschäftigen.

### 1.2.3 Durchführung der Studie (inklusive Datenerhebung)

Nach Abschluss der Planungs- und Vorbereitungsphase kann die Studie durchgeführt werden.

### 1.2.4 Datenerfassung und -prüfung

Wir verwenden bei unserem Demonstrationsprojekt zur Datenerhebung einen traditionellen Fragebogen. Damit fallen als nächstes folgende Arbeiten an:

- **Datenerfassung**

Das Eintragen der Rohdaten in eine Datei auf der Festplatte eines Computers kann mit dem Dateneditor von SPSS geschehen, mit einem speziellen Datenerfassungsprogramm oder (fehleranfällig!) mit einem normalen Texteditor. In jedem Fall ist bei der Erfassung der in der Planungsphase oder spätestens nach der Datenerhebung erstellte Kodierplan

---

<sup>1</sup> Bei der  $\beta$ -Fehler - basierten Kalkulation der Stichprobengröße kann z.B. das exzellente Programm Gpower eingesetzt werden. Eine Literaturangabe und eine kostenlose Bezugsquelle finden Sie in Abschnitt 7.



genau einzuhalten. Hier ist z.B. für jedes Merkmal festgelegt, wie seine Ausprägungen kodiert werden sollen.

Bei schriftlichen Befragungen großer Stichproben kann eine Anlage zum automatischen Einscannen und Interpretieren von Untersuchungsdokumenten rentabel eingesetzt werden. Voraussetzung ist dann u.a. die Beachtung einiger Regeln beim Entwurf der Untersuchungsmaterialien.

- **Überprüfung auf Erfassungsfehler**

Je fehleranfälliger die gewählte Erfassungsmethode war, desto mehr Aufwand ist bei der Datenprüfung angebracht.

Bei einer Online-Datenerhebung entfällt die Datenerfassung und -prüfung. Im Abschnitt 3.1.1 folgende weitere Informationen zu den Techniken der automatischen Datenerhebung- bzw. -erfassung.

### 1.2.5 Datentransformation

Nach der Erfassung und Prüfung liegen bei vielen Studien die Daten immer noch nicht in auswertbarer Form vor. Vielfach müssen Variablen überarbeitet (z.B. rekodiert) oder aus Vorläufern neu berechnet werden (z.B. durch Mittelwertsbildung). Solche Transformationen nehmen bei vielen Projekten einen erheblichen Umfang an, wobei sowohl akribische Fleißarbeit als auch kreative Begriffsbildung gefragt sind.

### 1.2.6 Statistische Datenanalyse

Nach langer Mühe können mit Hilfe von SPSS z.B. die gesuchten Schätzwerte ermittelt und die geplanten Hypothesentests durchgeführt werden. Bei einer eher explorativen Untersuchungsanlage ist eine längere, kreative Auseinandersetzung mit den Daten erforderlich, wobei zahlreiche Datentransformationen und statistische Analysen ausgeführt werden.

## 1.3 *Beispiel für eine empirische Untersuchung*<sup>1</sup>

Um die im Rahmen einer empirischen Untersuchung mit SPSS zu erledigenden Arbeiten unter realistischen Bedingungen üben zu können, wird im Verlauf des Kurses eine kleine psychologische Fragebogenstudie durchgeführt. Dabei werden Sie alle Phasen der empirischen Forschung von der ersten Idee bis zur statistischen Hypothesenprüfung mit Computerhilfe kennen lernen und die erforderlichen Arbeiten zum großen Teil selbständig durchführen. Als Beispiel wurde u.a. deshalb eine psychologische Fragebogenstudie gewählt, weil die Kursteilnehmer dabei in wenigen Minuten interessante empirische Daten selbst erzeugen können. Damit ist auch die Phase der *Datenerhebung* in den Übungsablauf einbezogen, die ansonsten aus Zeitgründen ausgespart werden müsste.

Bezogen auf das in Abschnitt 1.2 vorgestellte Schema beschäftigen wir uns nun mit dem theoretischen Hintergrund unserer Studie und mit Fragen der Untersuchungsplanung.

---

<sup>1</sup> Hierbei werden in stark vereinfachter Form Ideen aus der Forschungsabteilung von Herrn Prof. Dr. J. Brandstädter (Universität Trier) aufgegriffen, dem ich an dieser Stelle herzlich für die Erlaubnis und für die Überlassung von Untersuchungsmaterial danken möchte.

### 1.3.1 Die allgemeinspsychologische KFA-Hypothese

Nach einer Theorie von Kahneman<sup>1</sup> & Miller (1986) hängt die Stärke unserer emotionalen Reaktion auf ein positives oder negatives Ereignis u.a. davon ab, welche alternativen (aber nicht eingetretenen) Ereignisse wir uns vorstellen können, mit anderen Worten: welche **kontrafaktischen Alternativen** mental verfügbar sind. Wir wollen uns auf den Fall ungünstiger Ereignisse beschränken. Hierfür stellen Kahneman & Miller die folgende Hypothese auf:

**Bei einem negativen Ereignis erhöht die mentale Verfügbarkeit (Vorstellbarkeit) kontrafaktischer (also positiver) Ereignisalternativen den erlebten Ärger.**

Im weiteren Verlauf wollen wir unser Projekt kurz als *KFA-Studie* bezeichnen.

Weil diese Hypothese für beliebig aus der Population herausgegriffene Personen Gültigkeit beansprucht, kann sie als *allgemeinspsychologisch* bezeichnet und von *differentialpsychologischen* Hypothesen unterschieden werden, die sich mit Unterschieden zwischen Personen beschäftigen (siehe Abschnitt 1.3.3).

### 1.3.2 Untersuchungsplanung

Hinsichtlich des Untersuchungsdesigns haben wir uns aufgrund praktischer Erwägungen bereits auf eine **querschnittlich angelegte Fragebogenstudie** mit den Kursteilnehmern als **Beobachtungseinheiten** festgelegt.

Nun geht es um die **Operationalisierung** der theoretischen Begriffe bzw. um den Entwurf des Fragebogens. Wir wollen die Untersuchungsteilnehmer bitten, sich in eine Geschichte einzufühlen, bei der zwei Personen objektiv denselben Schaden erleiden, jedoch in unterschiedlichem Grad eine kontrafaktische (also günstige) Alternative vor Augen haben. Dann sollen die Probanden für jeden Geschädigten angeben, wie stark sie sich in dessen Lage ärgern würden. Die genaue Instruktion ist dem unten wiedergegebenen Fragebogen (Teil 2) zu entnehmen.

Indem wir jede Person den *beiden* imaginierten Behandlungen aussetzen, gewinnen wir jeweils *zwei* Beobachtungswerte, die eine statistische Analyse der allgemeinspsychologischen Hypothese mit hoher Teststärke (kleinem  $\beta$ -Fehler) ermöglichen sollen. Gegen diese Befragungstechnik lässt sich einwenden, dass durch die Präsentation der *beiden* Varianten ein Kontrast künstlich induziert, zumindest jedoch verstärkt wird (Artefakt!). Wer diese Artefaktgefahr für real hält, sollte an Stelle des Messwiederholungsfaktors KFA einen Gruppierungsfaktor setzen und jede Person nur zu *einer* Schädigungsvariante befragen.

Die beiden Ärgermessungen werden durch Ratingskalen realisiert, wobei das Antwortformat der Anschaulichkeit halber an ein Thermometer mit den Ankerpunkten 0° und 100° erinnert. Wir gehen davon aus, dass die Ärgermessungen annähernd Intervallniveau besitzen.

In Abschnitt 1.3.1 wurde die KFA-Hypothese noch ohne Bezug auf unsere Untersuchungsplanung formuliert. Jetzt nehmen wir eine Konkretisierung vor durch ...

- Verwendung von direkt beobachtbaren Begriffen
  - Bezug auf Verteilungsparameter (Erwartungs- bzw. Mittelwert)
- Eingangs wurde betont, dass unserer Hypothesen in der Regel probabilistischer Natur sind. Auch bei einer allgemeinspsychologischen Hypothese wird man kaum auf einer Gültigkeit für *alle* Personen einer Population bestehen (womöglich sogar mit derselben Effektstärke).

---

<sup>1</sup> Kahneman erhielt 2002 den Nobelpreis für Wirtschaft, womit vor allem seine erfolgreiche Anwendung psychologischer Erkenntnisse (u.a. zu Urteilen und Entscheidungen unter Unsicherheit) in wirtschaftswissenschaftlichen Theorien gewürdigt wurde.

Die konkretisierte Hypothese soll über die im statistischen Entscheidungsverfahren tatsächlich beteiligten Verteilungsparameter reden.

Außerdem soll hier der Klarheit halber (in einer für Forschungsberichte kaum zu empfehlenden Ausführlichkeit) dargelegt werden, dass bei einem inferenzstatistischen Entscheidungsverfahren *zwei* konkurrierende Hypothesen beteiligt sind:

**Nullhypothese:** Die Untersuchungsteilnehmer erleben in der Rolle des Geschädigten mit hochgradig verfügbarer kontrafaktischer Alternative im Mittel *nicht* mehr Ärger als in der Rolle des Geschädigten mit "weit entfernter" kontrafaktischer Alternative.

**Alternativhypothese<sup>1</sup>:** Die Untersuchungsteilnehmer erleben in der Rolle des Geschädigten mit hochgradig verfügbarer kontrafaktischer Alternative im Mittel *mehr* Ärger.

Wir wollen unser Entscheidungsproblem mit einem **t-Test für abhängige bzw. gepaarte Stichproben** lösen, falls die Verteilungsvoraussetzungen dieses Verfahrens erfüllt sind. Da gerichtete Hypothesen vorliegen, ist **einseitig** zu testen. Dabei wird eine Irrtumswahrscheinlichkeit erster Art in Höhe von  $\alpha = 5\%$  akzeptiert.

Unsere Studie soll aus praktischen Gründen mit der **studentischen Stichprobe** der Kursteilnehmer durchgeführt werden. Damit können unter induktivistischer Perspektive die Ergebnisse günstigstenfalls auf die Population der Studierenden generalisiert werden.

Da aus statistischer Sicht eine Stichprobe nie zu groß sein kann, sollen nach Möglichkeit *alle* Kursteilnehmer als Probanden gewonnen werden. Es ist aus praktischen Gründen nicht möglich, weitere Untersuchungsteilnehmer zu rekrutieren. Der Übung halber soll aber trotzdem eine  $\beta$ -Fehler - basierte Kalkulation des **Stichprobenumfangs** vorgenommen werden. SPSS unterstützt solche Berechnungen im Zusatzprogramm **SamplePower**, das uns leider nicht zur Verfügung steht. Stattdessen verwenden wir das exzellente Power-Analyse-Programm **GPower 3** (Faul et al., im Druck), das für MS-Windows und MacOS kostenlos über folgende Webseite zu beziehen ist:

<http://www.psych.uni-duesseldorf.de/abteilungen/aap/gpower3/>

Auf den Pool-PCs der Universität Trier unter dem Betriebssystem MS-Windows lässt sich GPower 3 über folgende Programmgruppe starten

**Start > Programme > Wissenschaftliche Programme > GPower**

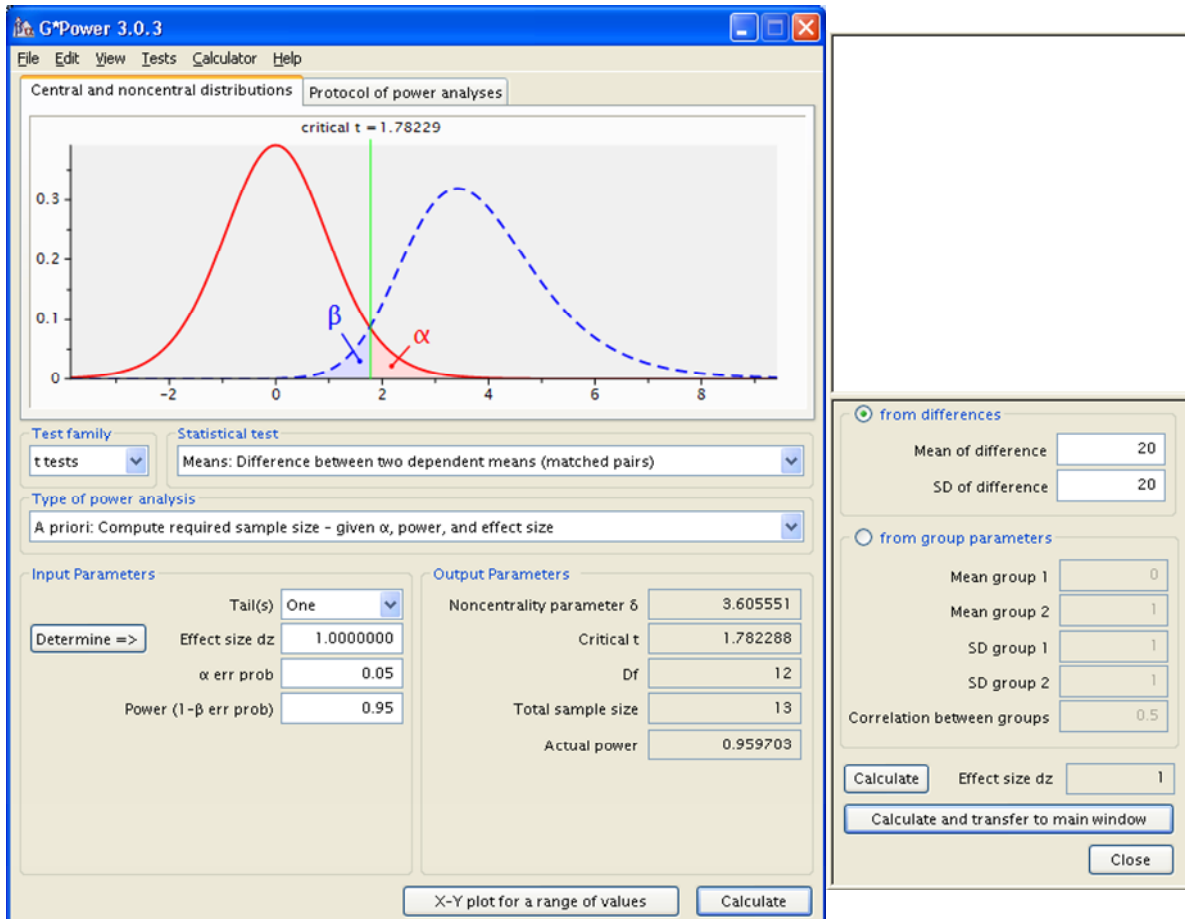
Wir wählen

- **t-Tests** als **Test family**
- **Means: Difference between two dependent means** als **Statistical test**
- **A priori** als **Type of power analysis**

und öffnen über den Schalter **Determine** ein Zusatzfenster, um die Effektstärke in der Population aufgrund theoretischer Annahmen und/oder bisheriger empirischer Erfahrungen festlegen zu können:

---

<sup>1</sup> Hier handelt es sich um einen statistischen Terminus, der nur zufällig mit unserer allgemeinpsychologischen Hypothese den Wortbestandteil *alternativ* gemeinsam hat.



Unsere KFA-Hypothese handelt vom *Ärgerzuwachs* aufgrund der kontrafaktischen Alternative und kann über die Differenz der beiden Ärgermessungen mit dem Einstichproben - t-Test beurteilt werden. Wir verwenden in GPower 3 diese Sichtweise, um die Effektstärke bequem festlegen zu können. Für den Einstichproben - t-Test ist die Teststärke  $d$  folgendermaßen definiert (vgl. z.B. Wentura 2004, S. 4):

$$d := \frac{\mu}{\sigma}$$

Darin sind:

- $\mu$  Mittelwert des betrachteten Merkmals (hier: der Differenz) in der Population
- $\sigma$  Standardabweichung in der Population

Als mittleren Ärgerzuwachs (Hauptparameter der KFA-Hypothese) erwarten wir ca. 20°. Als Standardabweichung der Differenz (Nebenparameter der KFA-Hypothese) erwarten wir aufgrund bisheriger Studien ebenfalls einen Wert von ca. 20. Mit dem Schalter **Calculate and transfer to main window** befördern wir die resultierende Effektstärke von 1 in das Hauptfenster. Nach einem Mausklick auf den Schalter **Calculate** erhalten wir das beruhigende Ergebnis, dass ...

- bei einem einseitigen Test (**Tail(s): One**)
- zum Niveau  $\alpha = 0,05$  (**alpha err prob**)
- für eine gewünschte Effektstärke (**Power**) von 0,95, also einen  $\beta$ -Fehler von 0,05

lediglich eine Stichprobe von 13 Fällen erforderlich ist. Sofern ein Effekt mit der angenommenen Stärke vorhanden ist, werden wir ihn also mit großer Wahrscheinlichkeit entdecken.

Wie die mit unserem Fragebogen erfassten Merkmale in der EDV-Welt repräsentiert werden können, wird in Abschnitt 1.4 behandelt.

Zuvor sollen noch einige zusätzliche Fragestellungen aufgegriffen und in den Untersuchungsplan aufgenommen werden.

### 1.3.3 Eine differentialpsychologische Hypothese

Neben der zentralen KFA-Hypothese soll in unserer Studie die folgende, auf Überlegungen von Scheier & Carver (1985) zurückgehende Hypothese überprüft werden:

**Der durch ein negatives Ereignis ausgelöste Ärger wird durch dispositionellen Optimismus gedämpft.**

Begründung: Dispositioneller Optimismus (im Sinne generalisierter positiver Ergebniserwartungen) führt zur Verwendung günstiger Bewältigungsstrategien (z.B. positive Reinterpretation). Während die allgemeinspsychologische KFA-Hypothese für eine beliebig aus der Allgemeinbevölkerung herausgegriffene Person einen bestimmten Effekt vorhersagt, geht es hier um Differentialpsychologie, also um Verhaltensunterschiede in Folge von relativ stabilen Personmerkmalen.

Als Quasiereignis soll der schon zur Prüfung der allgemeinspsychologischen Hypothese verwendete imaginierte Schadensfall dienen (Fragebogenteil 2, siehe unten).

Das arithmetische Mittel der für beide Situationsvarianten angegebenen Ärgerausprägungen soll uns als Ärgermaß dienen. Zur Erfassung von dispositionellem Optimismus wird der von Scheier & Carver (1985) entwickelte *Life Orientation Test* (LOT) eingesetzt (siehe Fragebogenteil 3). Wie aus den Antworten auf die 12 Fragen dieses Tests ein Optimismus-Messwert zu ermitteln ist, wird später erläutert. Wir gehen jedenfalls davon aus, dass diese Messmethode annähernd Intervallniveau besitzt.

Nach dieser **Operationalisierung** der theoretischen Begriffe kann die folgende **empirisch prüfbare Alternativhypothese** formuliert werden:

**Je höher der LOT-Wert einer Versuchsperson, desto weniger Ärger berichtet sie im Mittel für den imaginierten Schadensfall.**

Weil sich die Nullhypothese durch Negation der Alternativhypothese ergibt, muss sie nicht explizit notiert werden.

Weil die Messungen zum Ärger und zum Optimismus (hoffentlich) Intervallskalenniveau besitzen, kann die differentialpsychologische Hypothese mit einer **einfachen linearen Regressionsanalyse** geprüft werden, sofern deren Modell- und Verteilungsvoraussetzungen erfüllt sind.

Die Hypothese ist wiederum *einseitig* formuliert und soll mit einem  $\alpha$ -Fehler – Risiko von 5% geprüft werden.

Zur Berechnung der erforderlichen Stichprobengröße wählen wir in GPower 3:

- **t-Tests** als **Test family**
- **Correlation** als **Statistical test**
- **A priori** als **Type of power analysis**

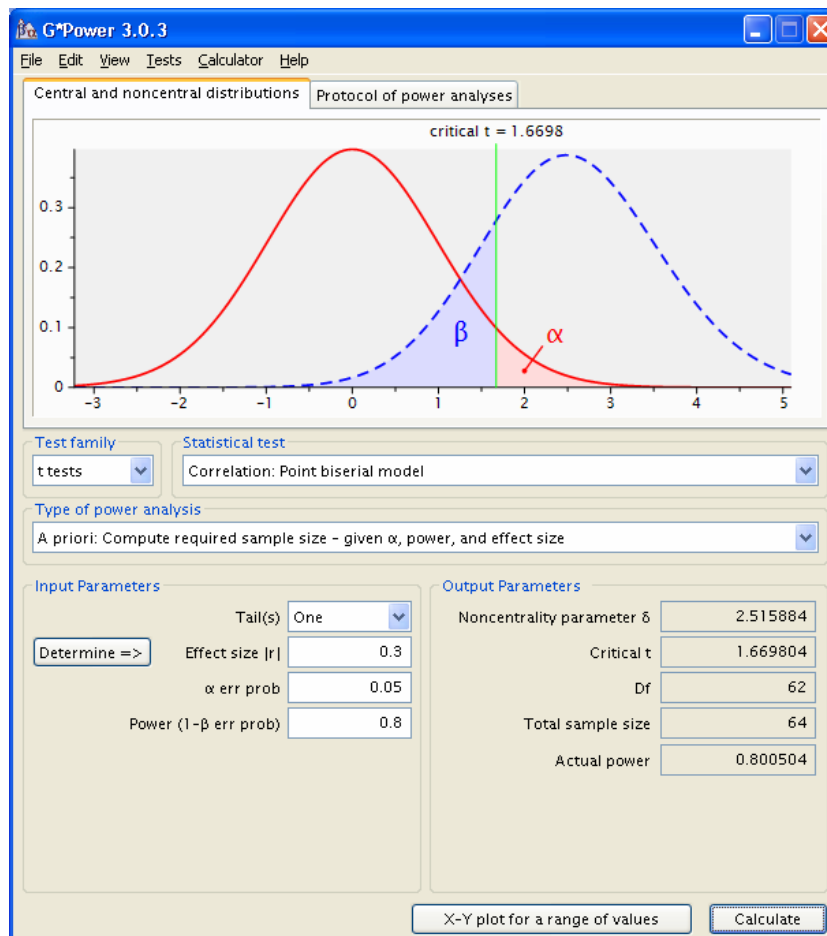
Begründung für die Deklaration eines Korrelationstests an Stelle einer Regressionsanalyse:

- Die beiden Tests verwenden dieselbe Prüfgröße und dieselbe Prüfverteilung (vgl. Abschnitt 7).
- Beim Korrelationstest unterstützt GPower 3 auch einseitige Fragestellungen (gerichtete Hypothesen). Zur gerichteten Hypothesen über Regressionskoeffizienten siehe Baltès-Götz (2006).

Das Effektstärkemaß ist bei einer Korrelation direkt über ihren Betrag definiert, und wir interessieren uns für einen Effekt der Stärke 0,3. Bei der Power geben wir uns mit einer Entdeckungswahrscheinlichkeit von 0,8 zufrieden. Trotzdem wird

- bei einem einseitigen Test
- zum Niveau  $\alpha = 0,05$

ein erforderlicher Stichprobenumfang von 13 Fällen berechnet:



Weil die Kursstichprobe in der Regel kleiner ist, stehen unserer Chancen einen Effekt von der vermuteten Stärke zu entdecken also eher schlecht. Bei einer gewünschten Power von 0,95 ( $\beta$ -Fehler 0,05) werden sogar 111 Fälle benötigt. In einem realen Forschungsprojekt zur Klärung der differentialpsychologischen Hypothese müsste der Stichprobenumfang folglich erhöht werden.

### 1.3.4 Zum Einfluss demographischer Merkmale

Auf die Erfassung demographischer Merkmale (siehe Fragebogenteil 1) kann man in keiner Studie verzichten, auch wenn sich keine expliziten Hypothesen darauf beziehen. Man benötigt sie auf jeden Fall zur Beschreibung der Stichprobe, damit sich später die Leser(innen) von Berichten ein Urteil über die Interpretier- bzw. Generalisierbarkeit der Ergebnisse bilden können. Wir wer-

den darüber hinaus einige demographische Merkmale auf Zusammenhänge mit unseren zentralen Projektvariablen untersuchen. Insofern finden sich auch in unserem überwiegend konfirmatorisch (hypothesenprüfend) angelegten Projekt einige explorative Elemente.

**1.3.5 Zu Übungszwecken erhobene Merkmale**

Ohne inhaltlichen Bezug zu den Fragestellungen des Projektes, sondern nur zu Übungszwecken sollen zusätzlich folgende Informationen erhoben werden:

- Größe und Gewicht (siehe Fragebogenteil 1)  
Mit diesen Merkmalen lassen sich manche statistische Verfahren gut demonstrieren. Außerdem sorgen sie für das Auftreten gebrochener Zahlen in unseren Daten.
- Motive zur Kursteilnahme (siehe Fragebogenteil 4)  
Hier wollen wir die Behandlung von Mehrfachwahlfragen sowie von offenen Fragen üben.

**1.3.6 Der Fragebogen**

**1) Angaben zur Person**

Geschlecht	Frau <input type="checkbox"/>	Mann <input type="checkbox"/>
Geburtsjahr		
Fachbereich		
Körpergröße	___, ___ ___	m
Körpergewicht	___ ___	kg

**2) Fragen zur Reaktion in ärgerlichen Situationen**

Versetzen Sie sich bitte möglichst gut in folgende Situation:

*Herr Meier und Herr Schulze waren mit demselben Taxi auf dem Weg zum Flughafen. Sie sollten zur selben Zeit, aber mit verschiedenen Maschinen abfliegen. Durch einen Stau kommen sie erst eine halbe Stunde nach der planmäßigen Abflugzeit am Flughafen an.*

*Herr Meier erfährt, dass seine Maschine pünktlich vor einer halben Stunde gestartet ist.*

*Herr Schulze erfährt, dass seine Maschine Verspätung hatte und erst vor zwei Minuten gestartet ist.*

Wie sehr würden Sie sich **ärgern**, wenn Sie in der Situation von ...

<b>Herrn Meier</b> wären?	0	10	20	30	40	50	60	70	80	90	100
<b>Herrn Schulze</b> wären?	0	10	20	30	40	50	60	70	80	90	100

Betrachten Sie bitte die Antwortskala als "Ärgerthermometer".

### 3) Aussagen zur Selbsteinschätzung

Teilen Sie bitte für die folgenden Selbstbeschreibungen durch Ankreuzen einer Antwortkategorie mit, inwiefern die Aussagen auf Sie persönlich zutreffen.

	völlig falsch	falsch	unentschieden	stimmt	stimmt genau
1. Auch in unsicheren Zeiten rechne ich im Allgemeinen damit, dass sich alles zum Besten wendet.	--	-	0	+	++
2. Ich kann mich leicht entspannen.	--	-	0	+	++
3. Wenn etwas schief gehen kann, dann passiert es mir auch.	--	-	0	+	++
4. Bei allem sehe ich stets die negative Seite.	--	-	0	+	++
5. Ich blicke kaum einmal mit Zuversicht in die Zukunft.	--	-	0	+	++
6. Ich bin gern mit Freunden zusammen.	--	-	0	+	++
7. Ich muss mich immer mit etwas beschäftigen.	--	-	0	+	++
8. Ich habe stets die Hoffnung, dass die Dinge in meinem Sinne gehen.	--	-	0	+	++
9. Die Dinge laufen immer so, wie ich es mir wünsche.	--	-	0	+	++
10. Ich bin nicht leicht aus der Ruhe zu bringen.	--	-	0	+	++
11. Ich glaube an den sprichwörtlichen "Silberstreifen am Horizont".	--	-	0	+	++
12. Dass mir einmal etwas Gutes widerfährt, damit rechne ich kaum.	--	-	0	+	++

### 4) Ihre Motive für die Teilnahme am SPSS-Kurs

a) Kreuzen Sie bitte in der folgenden Liste möglicher Motive für die Teilnahme am SPSS-Kurs alle für Sie zutreffenden Aussagen an und/oder nennen Sie Ihre sonstigen Motive.

Ich möchte SPSS kennen lernen, ...

- um eine eigene empirische Studie damit auszuwerten.
- weil in vielen Stellenanzeigen SPSS-Kenntnisse verlangt werden.
- weil ich mich um eine Stelle als EDV-Hilfskraft in der Forschung bewerben will (HIWI-Job).
- weil ich mich für EDV interessiere und ein modernes Programm kennen lernen möchte.
- weil ich mich für Statistik interessiere und mit Auswertungsverfahren experimentieren möchte.
- Andere Motive: \_\_\_\_\_

b) Möchten Sie im Kurs bestimmte statistische Methoden besonders gerne üben? Ja  Nein

Wenn „Ja“, welche? \_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_



## 1.4 Strukturierung und Kodierung der Daten

Wir werden die mit unserem Fragebogen erhobenen Informationen später manuell mit dem SPSS-Dateneditor erfassen und erstellen daher einen Kodierplan mit genauen Handlungsanweisungen für die Erfassung. Dabei müssen wir uns mit den Voraussetzungen beschäftigen, die SPSS für die Aufnahme unserer Daten bereitstellt. Diese sind in erster Linie durch die Logik der empirischen Forschung und nur in geringem Ausmaß durch EDV-Restriktionen festgelegt.

Bei der automatischen Erhebung bzw. Erfassung (Online-Formulare, Daten-Scanner) wird kein Kodierplan als Arbeitsvorschrift für Datenerfasser benötigt, jedoch kann auch hier eine Dokumentation der Daten nützlich sein (z.B. für die Kooperation in einer Arbeitsgruppe). Die in Abschnitt 1.4 behandelten Fragen werden bei den automatischen Methoden teilweise bei der Datendeklaration gegenüber der Umfrage- bzw. Scanner-Software geregelt, teilweise vom Automaten entschieden. Bei manchen Aufgaben sind Urteilsvermögen und Handarbeit eines Menschen durch keinen Automaten zu ersetzen, z.B. bei der Behandlung der Antworten auf offene Fragen (siehe Abschnitt 1.4.2.4). Insgesamt kann der Abschnitt 1.4 auch solchen Lesern zur Lektüre empfohlen werden, die zu einer Online- oder Scannerlösung tendieren.

### 1.4.1 Fälle und Merkmale in SPSS

Wir haben oben bereits daran erinnert, dass in einer empirischen Studie bei den einbezogenen **Fällen** bzw. **Beobachtungseinheiten** die Ausprägungen etlicher **Merkmale** festgestellt werden. Nun wollen wir uns ansehen, wie die Merkmalsausprägungen der Fälle im SPSS-System gespeichert werden. Die ganz konkrete Demonstration von KFA-Beispieldaten im **SPSS-Dateneditorfenster** wird das Verständnis der anschließenden, wieder eher allgemein-methodologisch geprägten, Ausführungen sicher unterstützen. U.a. werden dabei auch einige zentrale Begriffe des SPSS-Systems erläutert:

#### a) Variable

Der Begriff *Variable* wird in der Literatur zur statistischen Datenanalyse häufig synonym zu *Merkmal* gebraucht. Wir wollen ihn SPSS-konform in einer etwas technischeren Bedeutung verwenden: Schreibt man für ein Merkmal die Ausprägungen aller Fälle in der Stichprobe untereinander, so entsteht ein Spaltenvektor. Genau einen solchen Spaltenvektor wollen wir als *Variable* bezeichnen.

Zwar resultieren Variablen meist wie gerade beschrieben aus jeweils einem Merkmal, doch kann z.B. das Bemühen um eine rationelle Datenerfassung zu Ausnahmen führen. In Kürze wird eine Technik vorgeschlagen, die zur Erfassung von 100 Merkmalen mit Hilfe von fünf Variablen führt.

#### b) Datenmatrix und Dateneditor

Schreibt man alle Variablen nebeneinander, so entsteht die (Fälle  $\times$  Variablen) - Datenmatrix (Datentabelle). Sie kann bei der Datenerfassung im Fenster des SPSS-Dateneditors aufgebaut und dort auch während der laufenden Auswertungsarbeit ständig eingesehen oder bearbeitet werden. Die folgende Abbildung zeigt das Dateneditorfenster mit KFA-Beispieldaten aus einem früheren SPSS-Kurs:

	fnr	geschl	gebj	fb	gressse	gewicht	aergo	aergm	lot1	lot2	lot3	lot4	lot5	lot6	lot7	lot8	lot9	lot10	lot11	lot12	motiv1
1	1	1	69	1	163	51	5	8	4	2	4	5	4	5	3	4	4	0	4	4	1
2	2	1	70	1	158	56	5	8	4	3	5	4	4	4	3	4	2	3	4	4	1
3	3	1	69	1	174	58	4	8	4	2	3	4	4	4	5	4	3	1	3	4	0
4	4	2	67	1	182	77	6	2	4	4	4	5	4	5	3	3	2	4	4	4	1
5	5	1	67	1	180	69	8	8	3	1	4	4	4	5	5	4	3	4	4	5	1
6	6	1	66	1	175	72	8	10	2	2	4	5	4	5	1	4	4	3	3	5	0
7	7	1	75	1	167	50	6	8	3	3	3	2	3	4	3	3	2	2	3	4	1
8	8	1	74	1	163	54	5	6	4	3	3	3	5	5	3	4	4	2	4	5	1
9	9	2	67	1	185	70	4	4	3	3	4	4	5	5	2	3	2	3	4	5	0
10	10	1	64	3	164	57	6	10	4	2	4	5	5	5	4	3	4	2	5	5	1
11	11	1	70	6	176	54	2	6	4	2	3	2	4	5	4	4	3	2	3	5	1
12	12	2	72	6	190	96	10	10	4	3	2	4	4	5	3	3	2	4	3	3	1
13	13	1	70	1	162	58	8	10	3	2	2	2		5	3	4	2	2	4	3	
14	14	2	70	4	178	70	3	5	4	2	4	1	5	4	3	3	4	4	4	5	1

Jede Variable, d.h. jede Spalte der Datenmatrix, besitzt einen eindeutigen **Variablennamen**, über den sie bei der Anforderung statistischer oder graphischer Analysen angesprochen werden kann.

Nachdem Sie einen exemplarischen Eindruck vom *Ziel* der Strukturierungs- und Kodierungsmaßnahmen gewonnen haben, werden wir nun einige Details behandeln und einen Kodierplan für unser Projekt erstellen. Dabei soll u.a. angestrebt werden, den Aufwand und die Fehlergefahr beim Erfassen der Daten möglichst gering zu halten.

## 1.4.2 Strukturierung

Welche SPSS-Variablen im oben besprochenen Sinn sollen zur Aufnahme der mit unserem Fragebogen erfassten Informationen definiert werden? Obwohl die Antwort auf diese Frage trivial zu sein scheint, sind doch zu einigen Themen kurze Erläuterungen angebracht.

### 1.4.2.1 Variablen zur Fallidentifikation

Über die empirischen Variablen hinaus sollten in die Datenmatrix stets organisatorische Variablen aufgenommen werden, die eine Relation zwischen den schriftlichen oder sonstigen Untersuchungsdokumenten eines Falles und seinen Daten im Rechner herstellen. Eine solche Korrespondenz ist für eventuelle spätere Kontrollen oder Korrekturen der Daten unbedingt erforderlich. Meist verwendet man für diesen Zweck eine *einzelne* Variable, die z.B. FNR (für *Fallnummer*) genannt werden kann. Natürlich muss die Fallidentifikation auch auf den schriftlichen oder sonstigen Untersuchungsdokumenten eingetragen werden.

Bei personbezogenen Daten wählt man aus Datenschutzgründen zur Fallidentifikation z.B. eine zufällig vergebene Nummer ohne jeden Bezug zu den Personalien.

Möglicherweise erscheint Ihnen das Eintippen einer Identifikations-Variablen sinnlos, weil im SPSS-Dateneditor (siehe Abbildung in Abschnitt 1.4.1) die Zeilen bzw. Fälle ohnehin fortlaufend nummeriert sind. Die Nummern der Datenfensterzeilen stellen jedoch die gewünschte Korrespondenz zwischen den Datensätzen im Rechner und den nummerierten schriftlichen Untersuchungsunterlagen *nicht zuverlässig* her. Die Nummerierung der Datenfensterzeilen kann sich nämlich leicht ändern, z.B. wenn ein Sortieren der Fälle nötig wird, oder wenn Fälle gelöscht oder eingefügt werden.

### 1.4.2.2 Abgeleitete Variablen gehören nicht in den Kodierplan

Häufig sind in einem Forschungsprojekt nicht nur die direkt erfassten *Rohvariablen* von Interesse, sondern auch darauf aufbauende Variablen. Im KFA-Projekt soll etwa der Optimismus der Untersuchungsteilnehmer durch ihre mittlere Antwort auf die LOT-Fragen geschätzt werden. SPSS verfügt über leistungsfähige Befehle zur Berechnung neuer Variablen aus bereits vorhandenen, so dass derartige Routinearbeiten keinesfalls bei der Datenerfassung (z.B. per Taschenrechner) erledigt werden sollten.

Freilich müssen nach diesem Vorschlag *alle* Ausgangsvariablen aufgenommen werden, was aber vielfach ohnehin erforderlich ist (z.B. zur Überprüfung messtechnischer Eigenschaften).

Erfassen Sie also ausschließlich Rohvariablen, und führen Sie alle erforderlichen Transformationen später mit SPSS-Methoden durch. Wir werden uns im weiteren Kursverlauf mit den SPSS-Transformationsmethoden ausführlich beschäftigen. Im Kodierplan mit den Handlungsanweisungen für die Datenerfassung haben abgeleitete Variablen jedenfalls nichts zu suchen.

### 1.4.2.3 Mehrfachwahlfragen

Im Teil 4a unseres Fragebogens teilen die Untersuchungsteilnehmer für fünf konkrete Motive und eine Restkategorie mit, ob sie bei ihrer Entscheidung für die Kursteilnahme relevant waren. Damit erfahren wir von jeder Person sechs eigenständige Merkmalsausprägungen und benötigen (ohne Komprimierungsverfahren, siehe unten) folglich in der SPSS-Datentabelle sechs Variablen, um die Antworten aufzunehmen, die wir z.B. durch die Zahlen Eins (für *trifft zu*) und Null (für *trifft nicht zu*) kodieren können.

Beim Umgang mit einer solchen Mehrfachwahl-Frage müssen Sie sich vor allem vor dem aussichtslosen Versuch hüten, alle Auskünfte zum Fragebogenteil 4a in *eine* Variable zu verpacken. Dies käme dem unsinnigen Versuch gleich, *mehrere* Werte (z.B. Zahlen) in *eine* Zelle der SPSS-Datenmatrix einzutragen.

#### 1.4.2.3.1 Vollständige Sets aus dichotomen Variablen

In unserem Beispiel führt also eine Mehrfachwahl-Frage zu sechs dichotomen SPSS-Variablen, die jeweils die Information darüber enthalten, ob ein bestimmtes Motiv (bzw. ein sonstiges Motiv) vorlag oder nicht.

Das folgende Datenfenster zeigt die sechs Variablen, hier mit den Namen MOTIV1 bis MOTIV5 und ANDERE, bei einem Fall mit dem Antwortmuster (1,0,0,0,1,0):

The screenshot shows the SPSS Data Editor window for a file named 'kfa.sav'. The window title is 'kfa.sav - SPSS Daten-Editor'. The menu bar includes 'Datei', 'Bearbeiten', 'Ansicht', 'Daten', 'Transformieren', 'Analysieren', 'Grafiken', 'Extras', 'Fenster', and 'Hilfe'. Below the menu bar is a toolbar with various icons. The main area displays a data table with the following structure:

	motiv1	motiv2	motiv3	motiv4	motiv5	andere
1	1	0	0	0	1	0

The status bar at the bottom indicates 'SPSS Prozessor ist bereit'.

Wir werden in Abschnitt 12 ein so genanntes **Mehrfachantworten-Set** bestehend aus diesen sechs Variablen definieren und mit seiner Hilfe eine gemeinsame Auswertung der Variablen vornehmen. An dieser Stelle müssen Sie jedoch unbedingt akzeptieren, dass wir es mit *sechs* Merkmalen bzw. Variablen zu tun haben, die eine gewisse Verwandtschaft und ein gemeinsames dichotomes Format besitzen.

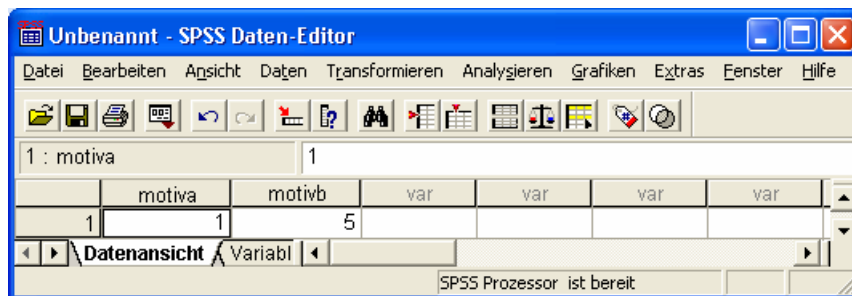
#### 1.4.2.3.2 Sparsame Sets aus kategorialen Variablen

Das im letzten Abschnitt beschriebene Standardverfahren zur Übersetzung einer Mehrfachwahlfrage in SPSS-Variablen ist angemessen, sofern nicht zu viele Antwortmöglichkeiten im Spiel sind. Wenn Sie etwa eine Liste mit 100 möglichen Freizeitaktivitäten präsentieren, dann führt das Schema zur Definition von 100 SPSS-Variablen. Unter der Annahme, dass jeder einzelne Untersuchungsteilnehmer maximal sieben verschiedene Optionen wählen wird, ist das Schema sicherlich unpraktisch bei der Datenerfassung. Für solche Situationen bietet sich ein alternatives Vorgehen an, das im eben konstruierten Freizeitbeispiel lediglich sieben Variablen bzw. Spalten in der SPSS-Datentabelle benötigt.

Auch dieses Komprimierungsverfahren soll an unserem Motivbeispiel demonstriert werden, obwohl es in diesem Fall (bei nur sechs Antwortmöglichkeiten) nicht geeignet ist. Unter der Annahme, dass pro Person maximal *zwei* verschiedene Motive zutreffen werden, definieren wir die beiden SPSS-Variablen MOTIVA und MOTIVB, die jeweils folgende Werte annehmen können:

- 1 für das Motiv *Eigene empirische Studie*
- 2 für das Motiv *Orientierung am Arbeitsmarkt*
- 3 für das Motiv *Bewerbung als EDV-Hilfskraft*
- 4 für das Motiv *Interesse an der EDV*
- 5 für das Motiv *Interesse an Statistik*
- 6 für andere Motive

Mit den Variablen MOTIVA und MOTIVB stehen für jede Person *zwei* Möglichkeiten zur Verfügung, um die „Hausnummern“ von angekreuzten Motiven zu erfassen. Das Antwortmuster (1,0,0,0,1,0) wird folgendermaßen übertragen:



Im Prinzip kann man im Beispiel die beiden Werte Eins und Fünf auch in umgekehrter Reihenfolge eintragen (MOTIVA = 5, MOTIVB = 1). Wesentlich ist nur, dass die Nummer jedes angekreuzten Motivs bei einer Variablen als Wert auftritt. Von einer Person, die zwei Motive angekreuzt hat, wissen wir *nicht*, welchem Motiv sie die größte Bedeutung beimisst. Daher können auch die resultierenden Variablen eine solche subjektive Ranginformation nicht enthalten. Allerdings wird man beim Erfassen der Systematik halber wohl so vorgehen, dass in MOTIVA die Nummer des ersten angekreuzten Motivs landet usw. (bei Anordnung von oben nach unten).

Wir sparen vier Variablen ein, wobei kein Informationsverlust eintritt, wenn tatsächlich pro Person maximal zwei Motive angekreuzt werden. Erweist sich ein sparsames Set während der Erfassung als unterdimensioniert, kann es bei Verwendung des SPSS-Dateneditors problemlos erweitert werden (z.B. um die Variable MOTIVC).

Auch bei der sparsamen Informationsanordnung kann man mit SPSS z.B. für jedes Motiv ermitteln, wie viel Prozent der Kursteilnehmer es angekreuzt haben. Vor einer solchen Auswertung ist wiederum ein Mehrfachantworten-Set zu definieren, diesmal bestehend aus den beiden Variablen MOTIVA und MOTIVB, wobei in der zugehörigen SPSS-Dialogbox eine *kategoriale* Kodierung der Variablen anzugeben ist.

Bei manchen Auswertungen ist es aber doch erforderlich, über Transformationsanweisungen das vollständige dichotome Set (mit einer Variablen pro Merkmal) herzustellen (siehe Abschnitt 12).

#### 1.4.2.4 Offene Fragen

Offene Fragen lösen vielfältige und oft schwer strukturierbare Antworten aus, und es bleibt dann offen, ob und wie die Antworten in SPSS-Variablen übersetzt werden sollen.

Ein Weg zur Systematisierung und Erfassung der Antworten besteht darin, eine Kategorienliste zu entwickeln und die vorhandenen bzw. fehlenden Nennungen der Listenelemente analog zu den Antworten auf eine Mehrfachwahl-Frage zu erfassen. Im Fall unseres Fragebogenteils 4b ist also durch Inspektion der ausgefüllten Fragebögen eine Liste mit speziell gewünschten statistischen Auswertungsverfahren erstellen, z.B. mit dem Ergebnis:

- Regressionsanalyse
- Kreuztabellenanalyse
- Faktorenanalyse
- Diskriminanzanalyse

Bei der Umsetzung in SPSS-Variablen wird man bei einem relativ kleinen Kategorienschema ein vollständiges Set mit dichotomen Variablen verwenden, ansonsten ein sparsames Set aus kategorialen Variablen (siehe oben). Aus der obigen vierelementigen Liste mit speziellen methodischen Interessen entsteht also ein vollständiges Set mit dichotomen Variablen, z.B.:

- REG für die Regressionsanalyse
- KT für die Kreuztabellenanalyse
- FAKT für die Faktorenanalyse
- DA für die Diskriminanzanalyse

Bei der Variablen REG ist eine Eins einzutragen, wenn ein Fall auf die offene Frage hin die Regressionsanalyse angegeben und damit sein Interesse an dieser Methode signalisiert hat. Anderenfalls wird eine Null notiert, die aber *nicht* als explizit bekundetes Desinteresse an der Regressionsanalyse zu interpretieren ist.

Das beschriebene Vorgehen erfordert zum Erstellen der Kategorienliste eine bei großen Stichproben recht aufwändige Vorauswertung der Fragebögen, die sich mit folgendem Trick vermeiden lässt: Man verwendet eine *dynamisch wachsende Liste* in Verbindung mit einem sparsamen Set kategorialer Variablen. In unserem Beispiel kann man z.B. über ein sparsames Set aus drei Variablen mit den Namen METH1 bis METH3 für jeden Fall maximal drei spezielle Auswertungsinteressen festhalten. Die Kategorienliste wird erst während der Datenerfassung entwickelt, indem man bei jedem Fall entscheidet, in welche bereits definierten oder neu aufzunehmenden Kategorien seine Antworten einzuordnen sind. Die Liste kann dynamisch um beliebig viele Kategorien erweitert werden, weil die drei Variablen beliebig viele verschiedene Werte als Kategoriennummern aufnehmen können. Selbstverständlich müssen die neu aufgenommenen Kategorien mit den vergebenen Nummern sorgfältig dokumentiert werden. Falls mehrere Personen an der Erfassung beteiligt sind, muss die eindeutige Zuordnung durch entsprechende Verabredungen sichergestellt werden.

Offene Fragen sind sicher vielfach sinnvoll, weil sie Informationen zutage fördern können, an die bei der Untersuchungsplanung niemand gedacht hat. Gelegentlich sind die Antworten jedoch so spärlich oder so schlecht strukturierbar, dass eine *statistische* Analyse nicht lohnend erscheint. So werden erfahrungsgemäß im Teil 4a des Beispielfragebogens kaum individuelle Motive zur Kursteilnahme angegeben, und wir ignorieren diese offene Frage im weiteren Projektverlauf.

### 1.4.3 Kodierung

Für jedes erhobene Merkmal muss festgelegt werden, wie die einzelnen Merkmalsausprägungen kodiert werden sollen. Dabei ist eine Kodierung durch einfach aufgebaute Werte anzustreben (z.B. durch positive, ganze Zahlen). Bei konkreten Überlegungen zur Kodierung müssen wir berücksichtigen, welche Variablentypen von SPSS unterstützt werden:

#### 1.4.3.1 Die wichtigsten Variablentypen in SPSS

An dieser Stelle beschränken wir uns auf die wichtigsten Variablentypen, mit denen die meisten Projekte auskommen:

- **Numerische Variablen**

Werte: reelle Zahlen

Z.B. geeignet für die Merkmale: - Alter  
- Größe  
- Gewicht

- **Zeichenkettenvariablen (synonym: alphanumerische Variablen, String-Variablen)**

Werte: Folgen von Zeichen (Buchstaben, Ziffern, Sonderzeichen), bis zur SPSS-Version 12 beschränkt auf die maximale Länge 255

Z.B. geeignet für die Merkmale: - Familienname  
- Man könnte das Merkmal Geschlecht alphanumerisch kodieren mit den Werten **weiblich** und **männlich**.

- **Datumsvariablen**

Werte: Datumsangaben

Z.B. geeignet für das Merkmal: Geburtsdatum

Anwendungsfälle für Datumsvariablen dürften in der Regel klar erkennbar sein. Ansonsten müssen Sie sich nur zwischen der numerischen und der alphanumerischen Kodierung entscheiden.

Bei Merkmalen mit **mindestens ordinalem Skalenniveau** ist offensichtlich nur die numerische Kodierung sinnvoll.

Bei Merkmalen mit **Nominalskalenniveau** hat man hingegen die Wahl zwischen numerischer und alphanumerischer Kodierung der Merkmalsausprägungen.

Beispiel Geschlecht: - numerische Kodierung: **1** für Frauen, **2** für Männer  
- alphanumerische Kodierung: **f** für Frauen, **m** für Männer

Beim Arbeiten mit SPSS empfiehlt es sich, auch nominalskalierte Merkmale numerisch zu kodieren, weil manche Auswertungsverfahren für diese Merkmale nur numerische Variablen akzeptieren (z.B. die Diskriminanzanalyse).<sup>1</sup>

<sup>1</sup> Offenbar überarbeitet SPSS sukzessive alle Prozeduren dahingehend, dass auch *kurze* String-Variablen (mit maximal achtstelligen Werten) akzeptiert werden, wenn in statistischer Hinsicht nur Nominalskalenniveau erforderlich ist. Diese Anpassung ist jedoch noch nicht für alle Prozeduren erfolgt.

### 1.4.3.2 Das Problem fehlender Werte

Trotz aller Sorgfalt sind in fast jedem Forschungsprojekt bei manchen Fällen einige Variablenausprägungen nicht bekannt, z.B. wegen technischer Probleme oder wegen nachlässig ausgefüllter Fragebögen. Bei der Kodierungsplanung muss daher für alle betroffenen Variablen festgelegt werden, was an Stelle fehlender oder ungültiger Werte in die zugehörigen Zellen der Datenmatrix eingetragen werden soll. Diese Ersatzwerte bezeichnet man häufig als *MD-Indikatoren*, wobei *MD* für *missing data* steht. Gelegentlich sind bei einer Variablen sogar mehrere MD-Indikatoren nötig, wobei z.B. ein erster Indikator signalisiert *Frage trifft nicht zu* und ein zweiter bedeutet *Keine auswertbare Antwort geliefert*.

Beispiel: Angenommen, wir hätten uns im demographischen Teil unseres Fragebogens danach erkundigt, ob ein Teilnehmer Wehr- bzw. Zivildienst abgeleistet hat. Dann könnten wir zu dieser Frage die SPSS-Variable DIENST definieren und dabei folgende Kodierungsregeln vereinbaren:

- *Nein* wird durch 0 kodiert.
- *Ja* wird durch 1 kodiert.
- *Ausmusterung* wird durch 2 kodiert.
- Frauen erhalten bei DIENST den Wert 8 (*Frage trifft nicht zu*).
- Verweigert ein Mann die Antwort, erhält er den Wert 9.

Beachten Sie bei der Verwendung von benutzerdefinierten MD-Indikatoren folgende Regeln:

- Es ist klar, dass alle MD-Indikatoren einer Variablen außerhalb des validen Wertebereichs liegen müssen. So wäre z.B. die 99 kein geeigneter MD-Indikator für unsere Variable Körpergewicht (gemessen in kg).
- Wählen Sie möglichst prägnante oder extreme Werte (also z.B. bei einer Variablen mit den validen Werten 1 und 2 den MD-Indikator 9). Dies bewirkt warnend auffällige Ergebnisse, falls Fälle mit fehlenden Werten nicht ordnungsgemäß von einer Analyse ausgeschlossen wurden.
- Der Einfachheit halber sollte für alle Variablen mit ähnlichem Wertebereich derselbe MD-Indikator verwendet werden.

**Wichtig:** Für jede betroffene Variable müssen dem SPSS-System alle benutzerdefinierten MD-Indikatoren bekannt gemacht werden (siehe Abschnitt 3.2.2).

#### 1.4.3.2.1 System-Missing (SYSMIS)

Neben den vom Benutzer variablenspezifisch vereinbarten MD-Indikatoren verwendet SPSS für alle numerischen Variablen automatisch einen weiteren MD-Indikator, der mit *System-Missing*, *systemdefiniert fehlend* oder *SYSMIS* bezeichnet wird. Er kommt immer dann zum Einsatz, wenn SPSS auf eines der folgenden Probleme trifft:

- Im Dateneditor bzw. beim Lesen einer bereits vorhandenen Datendatei (z.B. im Textformat) findet SPSS im Feld einer als numerisch definierten Variablen unzulässige Zeichen oder überhaupt keinen Eintrag.
- Beim Neuberechnen einer Variablen per Transformationsanweisung (siehe unten) fehlt ein Argument, oder der Funktionswert ist nicht definiert (z.B. bei Division durch 0).

Wir haben gerade erfahren, dass man beim Erfassen eines neuen Falles per SPSS-Dateneditor für eine Variable den Ersatzwert SYSMIS ganz einfach dadurch vereinbaren kann, dass man in die betroffene Zelle *nichts* einträgt.

**Tipp:** Bei der Datenerfassung mit dem SPSS-Dateneditor können Sie für numerische Variablen routinemäßig SYSMIS als MD-Indikator verwenden, bei Bedarf ergänzt durch zusätzliche benutzerdefinierte MD-Indikatoren. Man kann SYSMIS bequem dadurch vereinbaren, dass man die betroffene Zelle unverändert lässt. Weil SPSS den Ersatzwert SYSMIS automatisch richtig versteht, ist eine Deklaration nicht nötig und kann daher auch nicht vergessen werden.

Im Datenfenster und in der Ergebnisausgabe wird SYSMIS durch einen Punkt dargestellt (siehe Abbildung in Abschnitt 1.4.1, Variable LOT5 bei Fall 13).

#### 1.4.3.2.2 Fehlende Werte bei Mehrfachwahl-Fragen und offenen Fragen

Nachdem der Sinn und die Verwendung von MD-Indikatoren geklärt sind, geht es in diesem Abschnitt um eine spezielle Interpretationsunsicherheit im Zusammenhang mit fehlenden Werten, die bei Mehrfachwahl-Fragen aus der Verwendung eines probanden-freundlichen Antwortformates resultieren kann.

Im Fragebogenteil 4a zu den Motiven für die Kursteilnahme sorgt die sechste Ankreuzalternative (*Andere Motive*) durch Komplettieren der Antwortmöglichkeiten dafür, dass eine redliche Auskunftsperson mindestens eines der sechs Kästchen ankreuzen muss. Ohne diese Restkategorie könnten wir bei einem Fragebogen mit fünf leeren Motivkästchen folgende Möglichkeiten nicht unterscheiden:

- Bei der Person trifft tatsächlich keines der fünf vorgegebenen Motive zu.
- Die Person hat den Fragebogenteil 4a nicht bearbeitet (fehlende Daten).

Ursache für die Interpretationsunsicherheit ist offenbar das vereinfachte Antwortformat, das pro Motiv nur *ein* Kästchen vorsieht, statt jeweils ein Ja- *und* ein Nein-Kästchen vorzugeben. Damit ersparen wir den Untersuchungsteilnehmern zahlreiche Nein-Markierungen. Dies ist sinnvoll, damit deren Kooperationsbereitschaft nicht überstrapaziert wird, und die Fehlerquote gering bleibt.

Bei der offenen Frage in Teil 4b wird durch die vorgeschaltete Frage, ob überhaupt spezielle Methoden gewünscht sind, dafür gesorgt, dass bei Fragebögen ohne eingetragene Methodeninteressen folgende Möglichkeiten unterschieden werden können:

- Die Person hat kein Interesse an speziellen Auswertungsmethoden.
- Die Person hat den Fragebogenteil 4b nicht bearbeitet (fehlende Daten).

Durch das Bemühen um die Unterscheidbarkeit von verneinenden und fehlenden Antworten sollte das Fragebogendesign allerdings nicht zu umständlich bzw. pedantisch geraten.

#### 1.4.3.2.3 Vereinfachung der Erfassung durch Datentransformationstechniken

Im Zusammenhang mit dem MD-Problem bei den Variablen zu unserem Fragebogenteil 4 wage ich nun einige Vorschläge, die zwar dem Datenerfasser das Leben erleichtern, aber zugegebenermaßen die Kursteilnehmer(innen) beim ersten Entwurf eines Kodierplans durch einige zusätzliche Überlegungen belasten:

Bei der Mehrfachwahl-Frage nach den Kursmotiven haben wir geschickt durch die sechste Ankreuzalternative *Andere Motive* dafür gesorgt, dass Personen mit fehlenden Werten sicher zu identifizieren sind. Wir könnten den Erfasser im Kodierplan beauftragen:



- Schreibe bei den Variablen MOTIV1 bis MOTIV5 und ANDERE den Wert 1, wenn das zugehörige Kästchen markiert ist, sonst eine 0.
- Ist aber keines der sechs Kästchen markiert, dann versorge die Variablen MOTIV1 bis MOTIV5 und ANDERE mit dem vereinbarten MD-Indikator.

Die im zweiten Satz enthaltene Regel lässt sich mit (später anzuwendenden) SPSS-Transformationskommandos bequem automatisieren, so dass wir den Erfasser damit nicht belasten wollen. Damit wird die Lösung des MD-Problems zugunsten einer möglichst einfachen Erfassung in die spätere Projektphase der Datentransformation verschoben. Schlussendlich soll für die Variablen MOTIV1 bis MOTIV5 und ANDERE folgende Kodierung sichergestellt sein:

0	=	nein
1	=	ja
System-Missing	=	Wert unbekannt

Zur Erfassung der Informationen im Fragebogenteil 4b wollen wir eine dynamische Kategorienliste mit einem zugehörigem sparsamen Set kategorialer Variablen METH1 bis METH3 (vgl. Abschnitt 1.4.2.4) entwickeln. Der damit schon reichlich belastete Erfasser soll folgendermaßen vorgehen (bei Verwendung des SPSS-Dateneditors):

- Die Antwort auf die Frage, ob spezielle Methodenwünsche bestehen, wird konventionell in der Variablen SMG mit folgender Kodierungsvorschrift erfasst:

0	=	nein
1	=	ja
System-Missing	=	keine Antwort

- In die Dateneditorzellen zu den Variablen METH1 bis METH3 sollen die Kategoriennummern der gewünschten Methoden eingetragen werden. Bei weniger als drei Nennungen soll in den nicht benötigten Zellen nichts eingetragen werden, was zum MD-Indikator SYSMIS führt.

Diese Regel erleichtert die Erfassung und hat noch einen weiteren Vorteil: Sollte sich herausstellen, dass zusätzliche Variablen METH4 etc. benötigt werden, können wir diese ergänzen, ohne bei bereits erfassten Fällen irgendwelche Ersatzwerte (z.B. Nullen) nachtragen zu müssen.

Bei den Variablen METH1 bis METH3 soll später mit SPSS-Transformationsanweisungen dafür gesorgt, dass ihre Ausprägungen zuverlässig folgendermaßen interpretiert werden können:

0	=	Von der $i$ -ten ( $i = 1, \dots, 3$ ) Option zur Nennung einer interessierenden Methode wurde kein Gebrauch gemacht.
natürliche Zahl $\geq 1$	=	Die Methode mit dieser Kategoriennummer wurde angegeben.
System-Missing	=	Wert unbekannt

Dazu müssen unter den verschiedenen Wertekonstellationen der Variablen SMG und METH1 bis METH3 folgende Anpassungen vorgenommen werden:

		Mindestens eine speziell interessierende Methode angegeben?	
		Ja	Nein
SMG	1	METH1 ... METH3: SYSMIS → 0 Bem.: Korrektes Antwortverhalten. Variablen zu nicht benutzten Optionen (gem. Kodierplan bisher auf SYSMIS) werden auf 0 gesetzt.	SMG: 1 → SYSMIS Bem.: Irreguläres Antwortverhalten. METH1 bis METH3 behalten SYMIS. SMG wird ebenfalls auf SYMIS gesetzt.
	0	SMG: 0 → 1 METH1 ... METH3: SYSMIS → 0 Bem.: Leicht irreguläres Antwortverhalten. Wir sind großzügig und setzen SMG auf 1.	METH1 ... METH3: SYSMIS → 0 Bem.: Korrektes Antwortverhalten. Die Variablen zu allen Optionen (gem. Kodierplan bisher auf SYSMIS) werden auf 0 gesetzt.
	SYSMIS	SMG: SYSMIS → 1 METH1 ... METH3: SYSMIS → 0 Bem.: Leicht irreguläres Antwortverhalten. Wir sind großzügig und setzen SMG auf 1 sowie die Variablen zu nicht benutzten Optionen auf 0.	Bem.: Irreguläres Antwortverhalten. Alle Variablen behalten den Wert SYSMIS.

Vermutlich kam beim Lesen der letzten Ausführungen wenig Freude auf. Das MD-Problem verursacht oft erheblichen Aufwand, wobei auch Ermessenentscheidungen gefragt sind.

Jedenfalls ist die vorgeschlagene Methode zur Erfassung der Informationen aus dem Fragebogen teil 4b recht simpel und praktikabel.

#### 1.4.3.3 Fehlerquellen bei der manuellen Datenerfassung minimieren

Wenn die Daten manuell erfasst werden, ist bei den Kodierungsvereinbarungen darauf zu achten, dass dem Erfasser keine zeitaufwändigen und fehleranfälligen Arbeiten zugemutet werden, z.B.:

- Treten gebrochene Zahlen als Werte auf (z.B. bei unserer Frage nach der Körpergröße), so kann man durch Wechsel der Maßeinheit das lästige Dezimaltrennzeichen eliminieren.  
Beispiel: 1,65 m → 165 cm
- Bei bipolaren Skalen mit positiven und negativen Werten (z.B. bei unseren LOT-Fragen) empfiehlt sich eine Transformation zu ausschließlich positiven Werten z.B.:

--	→	1
-	→	2
0	→	3
+	→	4
++	→	5

Vorteil: Im Vergleich zu einer bipolaren Kodierung von -2 bis +2 spart man Tipparbeit und macht keine Fehler durch vergessene Vorzeichen bei den negativen Zahlen.

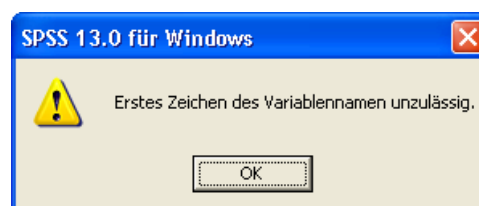
- Wurden einige Fragen aus messtechnischen Gründen umgepolt (negativ formuliert), was im KFA-Projekt bei einigen LOT-Fragen geschehen ist, so sollte diese Umpolung keinesfalls während der Erfassung rückgängig gemacht werden. Dies gelingt sehr viel bequemer und ohne Fehlerisiko mit den Transformationsmöglichkeiten von SPSS (siehe unten).

#### 1.4.3.4 SPSS-Variablenamen

Es empfiehlt sich, an dieser Stelle auch schon SPSS-Namen für die Variablen festzulegen und ebenfalls in den Kodierplan (siehe Abschnitt 1.4.3.5) aufzunehmen. Dabei sind die SPSS-Regeln für Variablenamen zu beachten:

- Maximal 64 Zeichen  
Die jahrzehntelange Beschränkung von SPSS-Variablenamen auf acht Zeichen ist seit der Version 12 überwunden, doch sollte man sich weiterhin möglichst kurz fassen. Lange Namen belegen viel Platz (z.B. in der Kopfzeile des Dateneditors) und sind beim Einsatz von SPSS-Syntax (siehe unten) recht umständlich.
- Das erste Zeichen muss ein Buchstabe sein.
- An den restlichen Positionen sind folgende Zeichen zugelassen: Buchstaben, Ziffern sowie die Symbole @, #, \_ und \$. Von der zweiten bis zur vorletzten Position ist außerdem der Punkt erlaubt.
- Aus den eben genannten Regeln ergibt sich insbesondere, dass Leerzeichen in Variablenamen verboten sind.
- Die von älteren SPSS-Versionen verschmähten Umlaute in Variablenamen werden mittlerweile akzeptiert. Allerdings sind Probleme zu erwarten, wenn eine SPSS-Datendatei zu einem Rechner mit einem anderen Betriebssystem transferiert wird. Der unter MS-Windows vereinbarte Variablenname „Größe“ kommt z.B. auf dem Macintosh als „Gr÷fle“ an, was durchaus zu Missverständnissen führen kann. Wir werden daher im Kurs Umlaute und „ß“ in Variablenamen vermeiden.
- Die folgenden Schlüsselwörter der SPSS-Kommandosprache (siehe unten) dürfen nicht als Variablenamen verwendet werden: ALL, AND, BY, EQ, GE, GT, LE, LT, NE, NOT, OR, TO, WITH.
- Die Groß-/Kleinschreibung ist irrelevant hinsichtlich der Identifikation von Variablen, jedoch verwendet SPSS bei Ausgaben die Schreibweise aus der Variablendeklaration. Wir schreiben in SPSS die Variablenamen aus Bequemlichkeitsgründen in Kleinbuchstaben. In Manuskript erscheinen sie zur Hervorhebung in Großbuchstaben.

Beim Versuch, einen irregulären Variablenamen zu vereinbaren, erhalten Sie im SPSS-Dateneditor eine meist informative Fehlermeldung, z.B.:



Tipps zur Benennung:

- Bilden Sie möglichst *informative* Namen, also z.B. FNR, GESCHL und GEBJ für *Fallnummer*, *Geschlecht* und *Geburtsjahr* an Stelle unpraktischer Bezeichnungen wie VAR1, VAR2, VAR3.
- Die eben genannte Regel muss in einem speziellen Fall relativiert werden: Bei Serien verwandter Variablen (z.B. die 12 LOT-Fragen im Teil 3 unseres Fragebogens) ist es in der Regel schwer, entsprechend viele individuelle Variablenamen zu bilden. Hier ist meist eine Indexschreibweise günstiger, bei der an einen informativen Namensstamm eine fortlaufende Nummer angehängt wird, z.B. LOT1, LOT2, ...

### 1.4.3.5 Kodierplan

Die Festlegungen zur Strukturierung und Kodierung der Projektdaten sollten in einem **Kodierplan** dokumentiert werden. Er hat zwei Funktionen:

- Während der Erfassung regelt er, wie die Daten eines Falles ins Dateneditorfenster einzutragen bzw. mit einem anderen Programm zu erfassen sind.
- Später dient der Kodierplan als kompakte Beschreibung der entstandenen Datendatei.

Bei unserer KFA-Studie kann für die geplante Erfassung mit dem SPSS-Dateneditor z.B. der folgende Kodierplan verwendet werden:

Merkmalsname	SPSS-Var.-name	Kodierung	Bemerkungen
Fallnummer	FNR	MD-Indikator: entfällt	
Geschlecht	GESCHL	1 = Frau 2 = Mann MD-Indikator: SYSMIS	
Geburtsjahr	GEBJ	<b>zweistellige</b> Eingabe! MD-Indikator: SYSMIS	
Fachbereich	FB	1 = I (Pädag., Philos., Psychol.) 2 = II (Sprachen) 3 = III (Hist. und polit. Wiss.) 4 = IV (BWL, Ethnol., Inform., Mathe, Soziol., VWL, Wirtsch.-Inf.) 5 = V (Jura) 6 = VI (Geowissenschaften) 7 = VII (Theologie) MD-Indikator: SYSMIS	
Körpergröße	GROESSE	Eingabe in <b>cm</b> ! MD-Indikator: SYSMIS	
Körpergewicht	GEWICHT	Eingabe in kg MD-Indikator: SYSMIS	
Ärger als Herr Meier (ohne KFA)	AERGO	0 = 0 1 = 10 . . . 10 = 100 MD-Indikator: SYSMIS	
Ärger als Herr Schulze (mit KFA)	AERGM	0 = 0 1 = 10 . . . 10 = 100 MD-Indikator: SYSMIS	
LOT-Fragen	LOT1 bis LOT12	1 = -- 2 = - 3 = o 4 = + 5 = ++ MD-Indikator: SYSMIS	
Kursmotive	MOTIV1 bis MOTIV5, ANDERE	0 = nicht angekreuzt 1 = angekreuzt	SYSMIS wird <b>nicht</b> vergeben! Die MD-Behandlung erfolgt später.
Spezielle Methoden gewünscht?	SMG	0 = nein 1 = ja MD-Indikator: SYSMIS	
Gewünschte statistische Methoden	METH1 bis METH3	1 = Meth.-Kat. 1 gew. 2 = Meth.-Kat. 2 gew. . . . Bei weniger als drei Nennungen: SYSMIS-Initialisierung belassen	Die Kategorienliste wird wäh- rend der Erfassung nach Bedarf entwickelt und dokumentiert. Die MD-Behandlung erfolgt später!

Dieser Kodierplan ist bei der Datenerfassung erfreulich einfach zu handhaben und leistet damit einen wichtigen Beitrag zur Integrität der auszuwertenden Daten.

Bei der Erfassung mit dem SPSS-Dateneditor (siehe Abschnitt 3.2) werden viele Regeln des Kodierplans in die Variablendeklaration einfließen. Dann wird eventuell die Frage auftauchen, ob man nicht auf einen Kodierplan verzichten und sein Regelwerk direkt im Deklarationsteil einer SPSS-Datendatei unterbringen kann. Allerdings enthält unser Beispiel viele Vorschriften (z.B. zweistellige Erfassung des Geburtsjahrs, Verlagerung der MD-Behandlung bei den Motiv-Fragen), die per Variablendeklaration nicht hinreichend klar dokumentiert werden können, um das Risiko von Erfassungsfehlern zu minimieren.

### 1.5 Durchführung der Studie (inklusive Datenerhebung)

Bei den obigen Überlegungen zur Strukturierung und Kodierung der Daten hat sich ergeben, dass der in Abschnitt 1.3 wiedergegebene Fragebogen ohne Korrekturen eingesetzt werden kann. Damit steht der Durchführung unserer Befragung nichts mehr im Wege.

Im realen Kursverlauf haben die Teilnehmer noch im Zustand der „naiven Unbefangenheit“ (vor Besprechung der KFA-Theorie) die Rolle der Probanden übernommen und so ihre eigenen, von zufälligen Stichprobeneffekten gefärbten Daten produziert. Die Leser(innen) im Selbststudium werden wohl aus praktischen Gründen in der Regel auf die Durchführung einer eigenen KFA-Erhebung verzichten. Im weiteren Verlauf des Manuskriptes werden die in einem früheren Kurs erhobenen Daten analysiert. Die zugehörigen Dateien können über das Internet bezogen werden (siehe Vorwort).

Hier ist der ausgefüllte Fragebogen derjenigen Untersuchungsteilnehmerin zu sehen, die bei der zufälligen Vergabe einer Fallidentifikation (vgl. Abschnitt 1.4.2.1) die Nummer 1 erhielt:

UNIVERSITÄTS-RECHENZENTRUM TRIER (URT)

**Statistisches Praktikum mit SPSS für Windows**

Beispielfragebogen

1

B. Baltes-Götz

- 2 -

**1) Angaben zur Person**

Geschlecht	Frau <input checked="" type="checkbox"/> Mann <input type="checkbox"/>
Geburtsjahr	1969
Fachbereich	I
Körpergröße	1.63 m
Körpergewicht	51 kg

**2) Fragen zur Reaktion in ärgerlichen Situationen**

Versetzen Sie sich bitte möglichst gut in folgende Situation:

*Herr Meier und Herr Schulze waren mit demselben Taxi auf dem Weg zum Flughafen. Sie sollten zur selben Zeit, aber mit verschiedenen Maschinen abfliegen. Durch einen Stau kommen sie erst eine halbe Stunde nach der planmäßigen Abflugzeit am Flughafen an.*

*Herr Meier erfährt, dass seine Maschine pünktlich vor einer halben Stunde gestartet ist.*

*Herr Schulze erfährt, dass seine Maschine Verspätung hatte und erst vor zwei Minuten gestartet ist.*

Wie sehr würden Sie sich ärgern, wenn Sie in der Situation von ...

Herrn Meier wären?	0	10	20	30	40	<input checked="" type="checkbox"/>	60	70	80	90	100
Herrn Schulze wären?	0	10	20	30	40	50	60	70	<input checked="" type="checkbox"/>	90	100

Betrachten Sie bitte die Antwortskala als "Ärgerthermometer".

**3) Aussagen zur Selbsteinschätzung**

Teilen Sie bitte für die folgenden Selbstbeschreibungen durch Ankreuzen einer Antwortkategorie mit, inwiefern die Aussagen auf Sie persönlich zutreffen.

	völlig falsch	falsch	unterschieden	stimmt	stimmt genau
1. Auch in unsicheren Zeiten rechne ich im Allgemeinen damit, dass sich alles zum Besten wendet.	--	-	o	<input checked="" type="checkbox"/>	++
2. Ich kann mich leicht entspannen.	--	<input checked="" type="checkbox"/>	o	+	++
3. Wenn etwas schief gehen kann, dann passiert es mir auch.	--	<input checked="" type="checkbox"/>	o	+	++
4. Bei allem sehe ich stets die negative Seite.	<input checked="" type="checkbox"/>	-	o	+	++
5. Ich blicke kaum einmal mit Zuversicht in die Zukunft.	--	<input checked="" type="checkbox"/>	o	+	++
6. Ich bin gern mit Freunden zusammen.	--	-	o	+	<input checked="" type="checkbox"/>
7. Ich muss mich immer mit etwas beschäftigen.	--	-	<input checked="" type="checkbox"/>	+	++
8. Ich habe stets die Hoffnung, dass die Dinge in meinem Sinne gehen.	--	-	o	<input checked="" type="checkbox"/>	++
9. Die Dinge laufen immer so, wie ich es mir wünsche.	--	-	o	<input checked="" type="checkbox"/>	++
10. Ich bin nicht leicht aus der Ruhe zu bringen.	--	-	<input checked="" type="checkbox"/>	+	++
11. Ich glaube an den sprichwörtlichen "Silberstreifen am Horizont".	--	-	o	<input checked="" type="checkbox"/>	++
12. Dass mir einmal etwas Gutes widerfährt, damit rechne ich kaum.	--	<input checked="" type="checkbox"/>	o	+	++

**4) Ihre Motive für die Teilnahme am SPSS-Kurs**

a) Kreuzen Sie bitte in der folgenden Liste möglicher Motive für die Teilnahme am SPSS-Kurs alle für Sie zutreffenden Aussagen an und/oder nennen Sie Ihre sonstigen Motive.

Ich möchte SPSS kennen lernen. ...

- um eine eigene empirische Studie damit auszuwerten.
- weil in vielen Stellenanzeigen SPSS-Kenntnisse verlangt werden.
- weil ich mich um eine Stelle als EDV-Hilfskraft in der Forschung bewerben will (HIWI-Job).
- weil ich mich für EDV interessiere und ein modernes Programm kennen lernen möchte.
- weil ich mich für Statistik interessiere und mit Auswertungsverfahren experimentieren möchte.
- Andere Motive: \_\_\_\_\_

b) Möchten Sie im Kurs bestimmte statistische Methoden besonders gerne üben? Ja  Nein

Wenn "Ja", welche? Faktorenanalyse  
Regressionsanalyse  
Korrelationsanalyse

Diese Nummer wurde nachträglich von der Untersuchungsleitung auf den Fragebogen geschrieben.

---

## 2 Einstieg in SPSS für Windows

In den bisher dargestellten Projektphasen von der theoretischen Ausarbeitung bis zur Erstellung des Kodierplans spielte SPSS noch keine wesentliche Rolle. Die im KFA-Projekt nun anstehende Datenerfassung wollen wir jedoch mit diesem Programm bewerkstelligen, so dass an dieser Stelle einige einführende Bemerkungen zu SPSS angemessen sind. In Abschnitt 2.1 geht es um die Verfügbarkeit von SPSS an der Universität Trier, und in den Abschnitten 2.2 bis 2.5 werden elementare Merkmale des Programms dargestellt.

### 2.1 SPSS für Windows an der Universität Trier

An der Universität Trier steht das Basissystem von SPSS für Windows mit den folgenden Erweiterungs-Modulen bzw. Zusatzprodukten zur Verfügung:

Erweiterungsmodule (in das Hauptprogramm integriert)	Zusatzprodukte (separat aufrufbar)
Regression Models Advanced Models Tables Trends Categories Conjoint Exact Tests Missing Values Analysis	Amos Answer Tree

Die aufgeführten SPSS-Produkte können auf folgende Weise genutzt werden:

#### a) Pool-PCs

Auf den Pool-PCs unter dem Betriebssystem MS-Windows finden Sie über

#### **Start > Programme**

die Programmgruppe **SPSS vom NT-Server des URT** mit Unterverzeichnissen zu allen installierten SPSS-Produkten.

#### b) Windows-Arbeitsplatzrechner im Campusnetz

Auch auf einem vernetzten Windows-Arbeitsplatzrechner kann wie auf einem Pool-PC die SPSS-Software über URT-Lizenzserver genutzt werden. Zur Installation der Programme stehen automatische Routinen zur Verfügung, die (im Rahmen einer normalen Anmeldung bei der Domäne URT) über

#### **Start > Systemsteuerung > Software > Neue Programme hinzufügen**

erreichbar sind.

#### c) Erwerb einer Einzelplatz-Mietlizenz

Beschäftigte und Studierende der Universität Trier können über das URT eine befristete SPSS-Mietlizenz zur Verwendung im Rahmen ihrer dienstlichen Tätigkeit bzw. ihrer Ausbildung erwerben, wobei auch eine Installation im Privathaushalt erlaubt ist. Nähere Informationen erhalten Sie in der URT-Benutzerberatung.

## 2.2 Programmstart und Benutzeroberfläche

### 2.2.1 SPSS starten

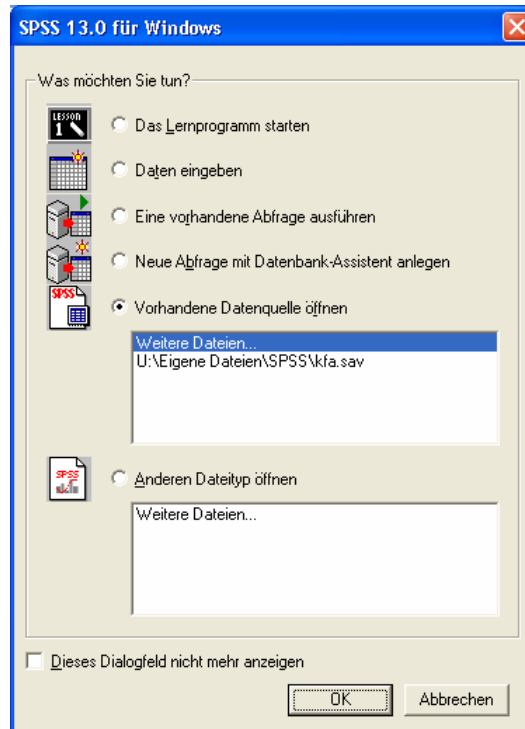
Nach erfolgreicher Anmeldung bei einem Pool-PC unter MS-Windows erreichen Sie SPSS 13 über das zugehörige Desktop-Symbol oder über das Startmenü:

**Start > Alle Programme > SPSS vom NT-Server des URT > SPSS 13.0 für Windows**

Auf einem PC mit lokaler SPSS-Installation können Sie das Programm in der Regel so starten:

**Start > Alle Programme > SPSS für Windows > SPSS 13.0 für Windows**

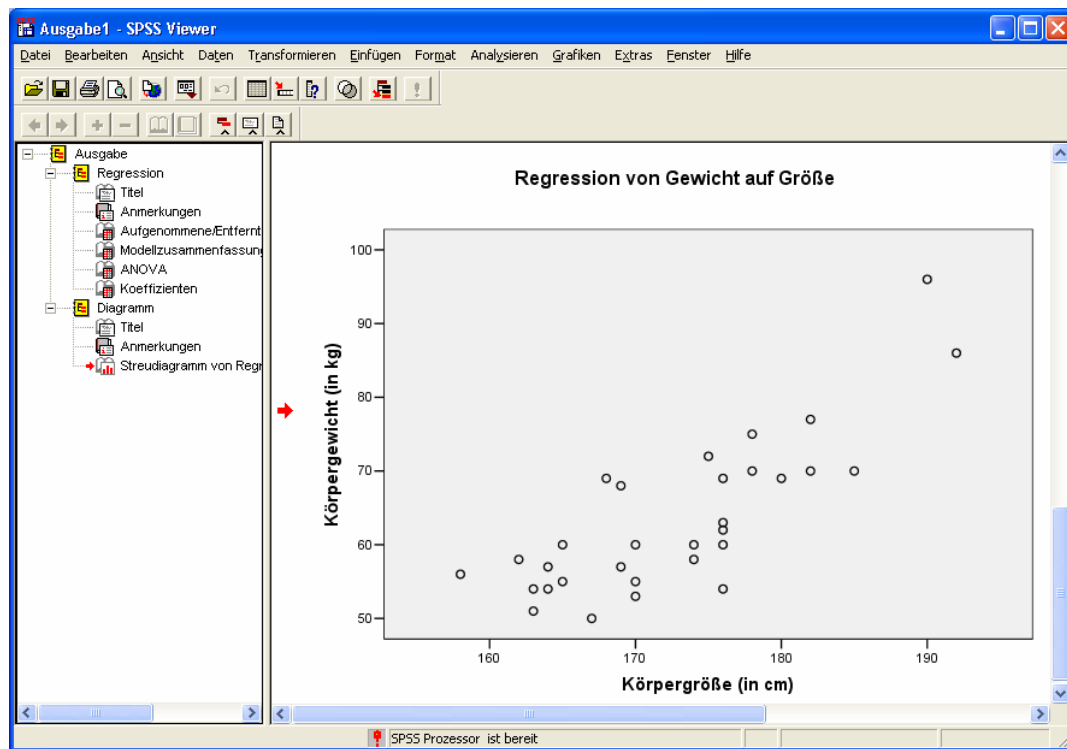
Nach dem Start erscheint der folgende Assistent:



Er ermöglicht z.B. ein bequemes Öffnen der in früheren Sitzungen benutzten Dateien.

### 2.2.2 Die wichtigsten SPSS-Fenster

Das Dateneditorfenster mit der (Fälle  $\times$  Variablen) - Datenmatrix haben Sie schon in Abschnitt 1.4.1 kennen gelernt. Nach der Datenerfassung können Sie mit Hilfe seiner Menüzeile statistische und grafische Datenanalysen anfordern, die dann im **Ausgabefenster**, auch *SPSS Viewer* genannt, erscheinen, z.B.:



Die SPSS-Fenster enthalten in der Kopfzone eine Menüzeile und verschiebbare Symbolleisten, im Fußbereich eine Statuszeile mit Informationen über wichtige Programmzustände.

### 2.2.3 Was man mit SPSS so alles machen kann

Wir sind im Moment dabei, einen ersten Eindruck vom Arbeitsplatz *SPSS für Windows* zu gewinnen. Einen guten Überblick vermitteln die Optionen in der Menüzeile des Dateneditorfensters:

- **Datei**  
Hier finden Sie u.a. Befehle zum Öffnen bzw. Sichern von Datendateien sowie zum Beenden von SPSS.
- **Bearbeiten**  
Über das **Bearbeiten**-Menü erreichen Sie Editorbefehle zum Ausschneiden, Kopieren, Einfügen, Löschen und Suchen von Daten sowie die **Optionen**-Dialogbox zur Anpassung von diversen SPSS-Einstellungen. Außerdem können Sie hier Modifikationen des Datenfensters rückgängig machen.
- **Ansicht**  
Hier können Sie u.a. die Statuszeile sowie die Symbolleisten aus- bzw. einschalten sowie die Schriftart der angezeigten Daten festlegen.
- **Daten**  
Über das **Daten**-Menü sind Dialoge zur Auswahl einer Teilstichprobe, zur Aggregation von SPSS-Dateien (z.B. mit Daten aus verschiedenen Stichproben) sowie zum Sortieren und Gewichten der Fälle erreichbar.
- **Transformieren**  
Hier finden Sie z.B. die Befehle zum Rekodieren von Variablen oder zum Berechnen neuer Variablen aus bereits vorhandenen.
- **Analysieren**  
Dieser Menüpunkt erschließt die statistischen Auswertungsmethoden, mit denen wir letztlich unsere Forschungsfragen klären wollen.



- **Grafik**  
An dieser Stelle bietet SPSS vielfältige Möglichkeiten zur grafischen Präsentation von Datenstrukturen an.
- **Extras**  
Hier finden sich diverse Funktionen (z. B. zur Anzeige von Informationen über die Variablen im Datenfenster) sowie ein Editor zur Modifikation der SPSS-Menüs.
- **Fenster**  
Über dieses Menü sind die offenen SPSS-Fenster erreichbar.
- **Hilfe**  
Hiermit starten Sie die Online-Hilfe, die Informationen über das gesamte SPSS-System bereithält und außerdem ein Lernprogramm sowie einen Statistik-Assistenten bietet.

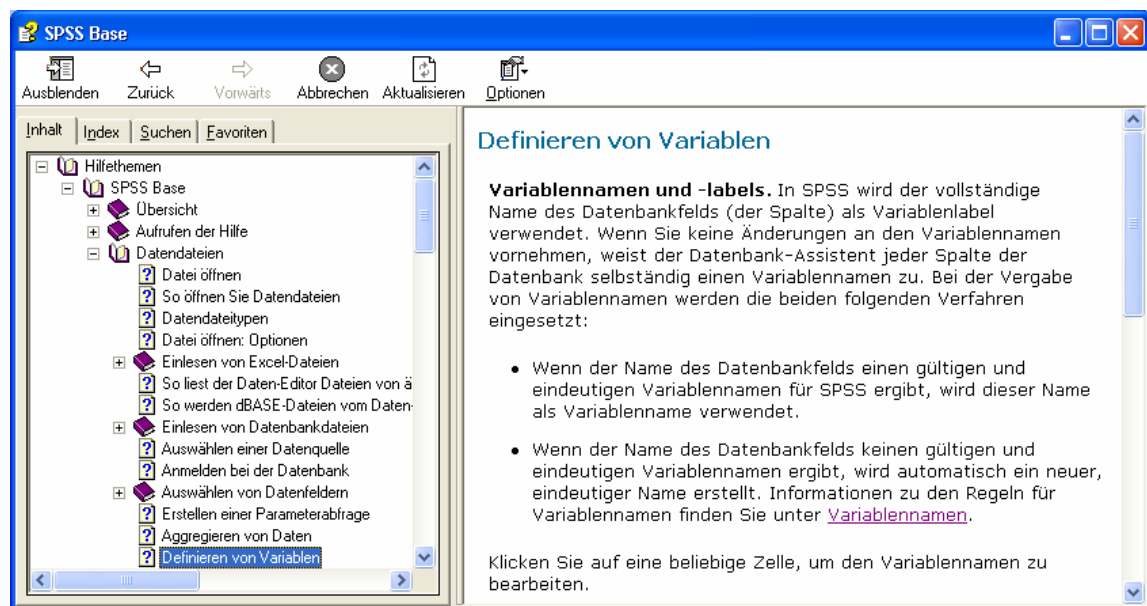
Bei leerem Datenfenster sind die meisten Menüoptionen nicht verfügbar.  
Die anderen SPSS-Fenster bieten angepasste Menüzeilen.

## 2.3 Das Hilfesystem

Bei der Arbeit mit SPSS für Windows können Sie stets auf ein mächtiges Hilfesystem zurückgreifen, dessen wichtigste Möglichkeiten nun vorgestellt werden.

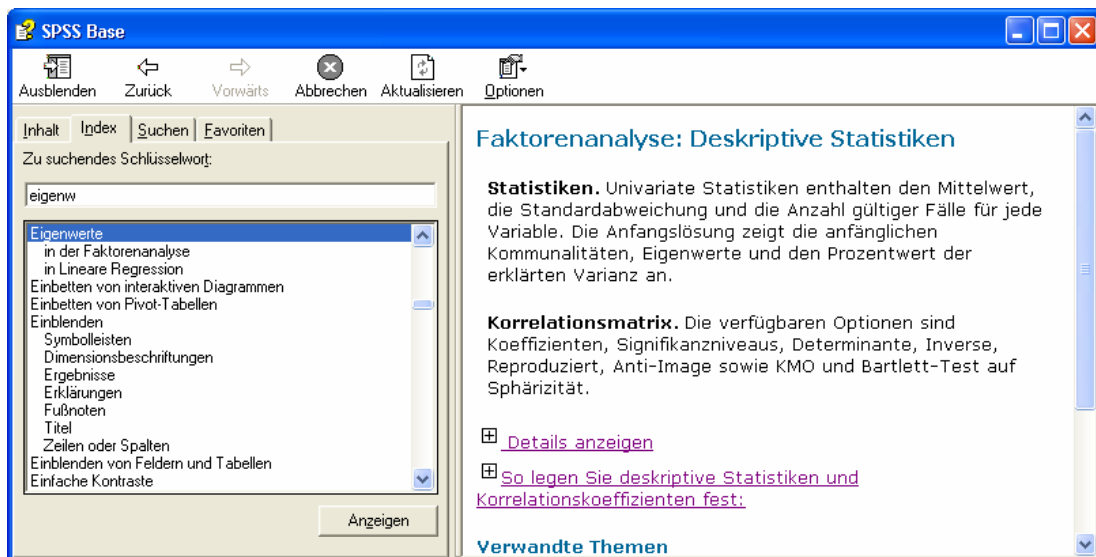
### 2.3.1 Systematische Informationen

Nach dem Menübefehl **Hilfe > Themen** finden Sie auf der **Inhalt**-Registerkarte des folgenden Fensters Informationen über die installierten SPSS-Module in systematischer Form:



### 2.3.2 Gezielte Suche nach Begriffen

Die Registerblätter **Index** und **Suchen** im Hilfefenster eignen sich für die Suche nach Informationen zu bestimmten Begriffen, z.B.:



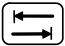
### 2.3.3 Kontextsensitive Hilfe zu den Dialogboxen

In fast jeder Dialogbox können Sie mit der Standardschaltfläche **Hilfe** Informationen zu all ihren Optionen anfordern.

### 2.3.4 Lernprogramm

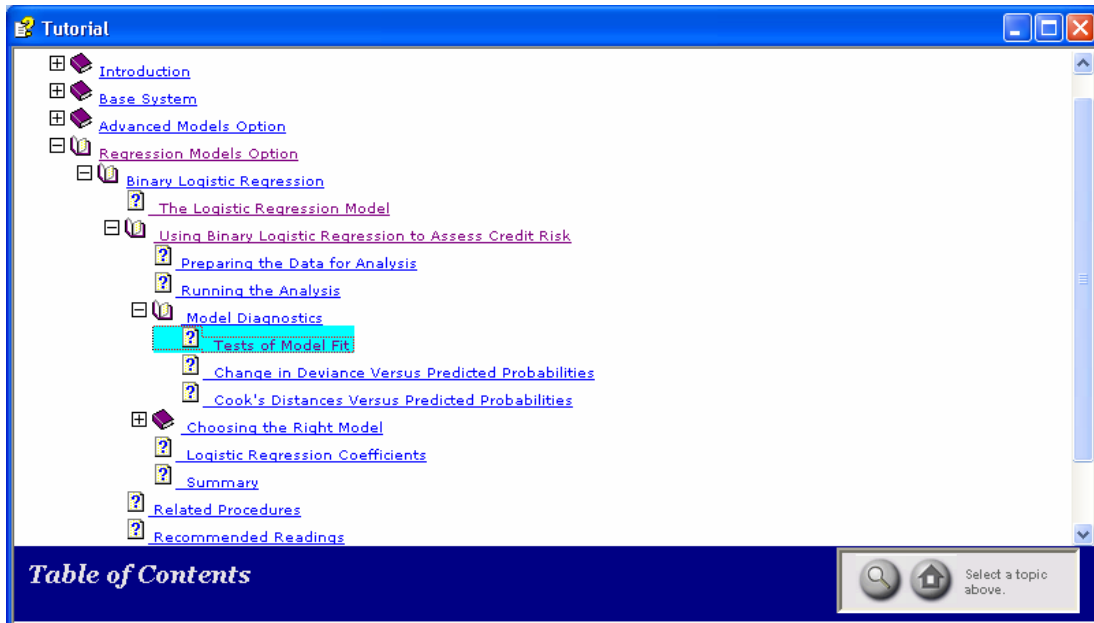
Neben dem eher zum Nachschlagen geeigneten Hilfenfenster mit seinen systematischen Beschreibungen und seinem vollständigem Index gibt es ein weiteres Informationsangebot, das eher didaktisch orientiert und daher auf das Wichtigste beschränkt ist: das interaktive SPSS-Lernprogramm. Es wird mit **Hilfe > Lernprogramm** gestartet und sollte mehr oder weniger linear durchgearbeitet werden. In den einzelnen Kapiteln werden konkrete Arbeitsabläufe geübt, z.B.:



Sie können das Lernprogramm als eigenständige Windows-Anwendung parallel zu SPSS ausführen und damit die Lektionen sofort nachvollziehen, indem Sie zwischen SPSS und dem Lernprogramm hin und her wechseln, z.B. mit der Tastenkombination **ALT** .

### 2.3.5 Fallstudien

Nach **Hilfe > Fallstudien** startet ein Tutorial, das mit der interaktiven Technik des Lernprogramms arbeitet, aber den Schwerpunkt auf statistische Analysen legt.

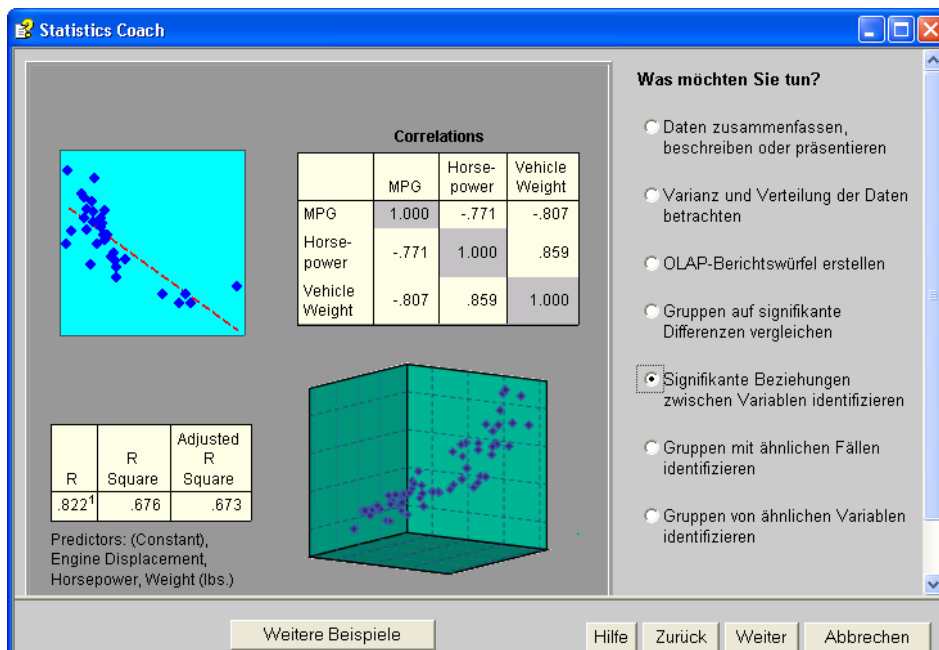


Viele Auswertungsprozeduren werden über ein komplettes Anwendungsbeispiel und Informationen zu folgenden Themen erschlossen:

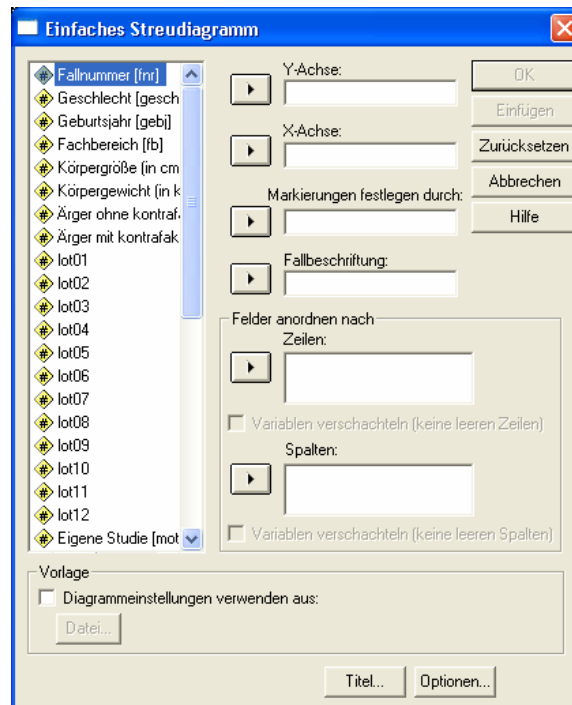
- Einsatzmöglichkeiten
- Anforderung der Analyse
- Interpretation der Ergebnisse
- Verwandte Verfahren
- Literaturangaben

### 2.3.6 Statistik-Assistent

Der über **Hilfe > Statistics Coach** verfügbare Assistent versucht, den Anwender durch eine Sequenz von Fragen zur richtigen Statistik- bzw. Grafikdialogbox zu führen, z.B.:



In einem Test mit dem abgebildeten Einstieg und der anschließenden Vereinbarung, dass zwei metrische Variablen grafisch analysiert werden sollen, hat der Assistent am Ende seiner Befragung tatsächlich ein geeignetes Verfahren vorgeschlagen und die zugehörige Dialogbox geöffnet:



## 2.4 Weitere Informationsquellen

### 2.4.1 Handbücher und Manuskripte

Es stehen u.a. zur Auswahl:

- SPSS-Originalhandbücher  
Mit SPSS wird eine umfangreiche Sammlung von PDF-Handbüchern zu den einzelnen Modulen und zu den statistischen Algorithmen ausgeliefert. Allein die Dokumentation der Kommandosprache, über die man die meisten Leistungen des SPSS-Systems abrufen kann (siehe unten), umfasst ca. 2000 Seiten. Dieses PDF-Dokument ist auch im Hilfesystem verfügbar (**Hilfe > Command Syntax Reference**).
- Sekundärliteratur  
Im Buchhandel und in wissenschaftlich orientierten Bibliotheken finden sich zahlreiche Sekundär-Handbücher zu SPSS.
- Auf die URT-Manuskripte zur Verwendung spezieller Analysemethoden in SPSS wurde schon im Vorwort hingewiesen.

Nach dem Absolvieren des vorliegenden Kurses sind für die meisten SPSS-Anwender(innen) solche Handbücher besonders nützlich, welche die jeweils benötigten statistischen Methoden auf einem angemessenen Niveau behandeln und die konkrete Realisation mit SPSS gut unterstützen (z.B. durch eine Erläuterung der Ergebnistabellen). Leider habe ich mir aus Zeitgründen von den zahlreichen Statistik-Lehrbüchern mit SPSS-Unterstützung nur wenige Titel näher ansehen können, so dass die folgende Liste sicher unvollständig ist:

- Backhaus et al. (2006). *Multivariate Analysemethoden*  
Cohen, et al. (2003). *Applied Multiple Regression/Correlation Analysis ...*  
Norušis (2005). *SPSS 14.0. Statistical Procedures Companion*  
Norušis (2005). *SPSS 14.0. Advanced Statistical Procedures Companion*  
Tabachnik & Fidell (2007). *Using multivariate statistics*

Die vollständigen bibliographischen Angaben finden sich im Literaturverzeichnis.

#### 2.4.2 SPSS im Internet

SPSS ist im Internet vielfach präsent, besonders zu erwähnen sind:

- **Die WWW-Homepage der SPSS Inc.:** <http://www.spss.com/>
- **Die Usenet-Diskussionsgruppe comp.soft-sys.stat.spss**  
Hier werden technische und statistische Themen lebhaft diskutiert, wobei SPSS-Mitarbeiter zu wichtigen Fragen kompetent Stellung nehmen.

#### 2.4.3 Benutzerberatung

Bei Problemen mit der Anwendung von SPSS können Sie sich an die URT-Benutzerberatung wenden.

Ort: im Gebäude E (Räume 002 - 014).  
Zeiten: Montag bis Freitag: 10.00-11.30 Uhr, Montag bis Donnerstag: 14-16 Uhr

#### 2.5 SPSS für Windows beenden

Die Beendigung einer SPSS-Sitzung wird mit

##### **Datei > Beenden**

eingeleitet. Falls Sie während der Sitzung Dokumente erstellt bzw. verändert und noch nicht gesichert haben (z.B. im Daten- oder im Ausgabefenster), werden Sie von SPSS an das Speichern erinnert.

---

## 3 Datenerfassung und SPSS-Dateneditor

Wie bei unserer KFA-Studie liegen auch in vielen anderen Projekten nach Abschluss der Datenerhebung schriftliche Untersuchungsdokumente vor, die nun erfasst, d.h. in eine Computerdatei übertragen werden müssen. Bevor in Abschnitt 3.2 die konkrete Erfassung der KFA-Daten mit dem SPSS-Dateneditor beschrieben wird, sollen in Abschnitt 3.1 einige alternative Erfassungsmethoden vorgestellt werden.

### 3.1 Methoden zur Datenerfassung

#### 3.1.1 Automatisierte Verfahren

Zunächst geht es um zwei Optionen zur Rationalisierung der Datenerfassung, die sich zunehmender Beliebtheit erfreuen.

##### 3.1.1.1 Online-Datenerhebung

Wenn die nötigen technischen und organisatorischen Voraussetzungen gegeben sind, sollte eine **Online-Datenerhebung** eingesetzt werden. Hiermit sind Verfahren gemeint, bei denen die Untersuchungsteilnehmer(innen) ihre Daten (aktiv oder passiv) direkt in eine EDV-Anlage einspeisen (z.B. Internet-Umfrage, automatische Aufzeichnung physiologischer Daten). Nach Abschluss der Datenerhebung kann sofort die Auswertung beginnen, weil die Daten automatisch in einer Datei landen, die meist direkt in SPSS genutzt werden kann. Auf eine gelegentliche Kontrolle (z.B. wegen möglicher Defekte in der Aufzeichnungsapparatur) sollte man aber trotzdem nicht verzichten. Die *Datenerfassung* als eigenständige Arbeitsphase entfällt bei den Online-Verfahren.

Mit der zunehmenden Verbreitung des Internets verbessern sich Chancen für den Einsatz dieser Kommunikations-Infrastruktur bei einer Vielzahl von Untersuchungen. Allerdings sind u.a. die folgenden Einschränkungen zu beachten:

- Man erreicht (noch) nicht jede Population.
- Für umfangreiche Befragungen ist die Technik weniger geeignet, weil die Unterbrechung und spätere Fortsetzung der Teilnahme umständlich ist, bei manchen Systemen sogar unmöglich.
- Wenn sich die Online-Umfrageteilnehmer in einer relativ öffentlichen Situation befinden (z.B. PC-Pool einer Hochschule), ist die Auskunftsbereitschaft bei persönlichen Fragen eventuell beschränkt.

Das URT betreibt Online-Umfragesysteme auf HTML- und PDF-Basis (**GlobalPark Umfragecenter 5.1**, **Teleform 9.1**), wobei sich z.B. der KFA-Fragebogen mit beiden Systemen gut realisieren lässt. Bei der HTML-basierten GlobalPark-Lösung wird auf Seiten der Umfrageteilnehmer lediglich ein Web-Browser vorausgesetzt:

**Umfrage - Mozilla Firefox**

http://www.unipark.de/uc/stat\_prakt\_spss/ospe.php?SES=(

## 2) Fragen zur Reaktion in ärgerlichen Situationen

Versetzen Sie sich bitte möglichst gut in folgende Situation:

*Herr Meier und Herr Schulze waren mit demselben Taxi auf dem Weg zum Flughafen. Sie sollten zur selben Zeit, aber mit verschiedenen Maschinen abfliegen. Durch einen Stau kommen sie erst eine halbe Stunde nach der planmäßigen Abflugzeit am Flughafen an.*

**Herr Meier** erfährt, dass seine Maschine pünktlich vor einer halben Stunde gestartet ist.  
**Herr Schulze** erfährt, dass seine Maschine Verspätung hatte und erst vor zwei Minuten gestartet ist.

Wie sehr würden Sie sich ärgern, wenn Sie in der Situation von ...

	10°	20°	30°	40°	50°	60°	70°	80°	90°	100°
<b>Herrn Meier</b> wären?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
<b>Herrn Schulze</b> wären?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Betrachten Sie bitte die Antwortskala als "Ärgerthermometer".

Zurück Weiter

Fertig

Bei der auf Teleform (siehe unten) basierenden PDF-Lösung kann das Design des Fragebogens über einen graphischen Editor gestaltet werden. Die Untersuchungsteilnehmer benötigen über den Web-Browser hinaus noch den kostenlos verfügbaren und sehr weit verbreiteten Acrobat-Reader der Firma Adobe:

http://urt-ds2.uni-trier.de/urt/PdfForms/data/KFA/Form/KFA.pdf - Microsoft Internet Explorer URT

Adresse http://urt-ds2.uni-trier.de/urt/PdfForms/data/KFA/Form/KFA.pdf

10. Ich bin nicht leicht aus der Ruhe zu bringen.

11. Ich glaube an den sprichwörtlichen "Silberstreifen am Horizont"

12. Dass mir einmal etwas Gutes widerfährt, damit rechne ich kaum.

### 4) Ihre Motive für die Teilnahme am SPSS-Kurs

a) Kreuzen Sie bitte in der folgenden Liste möglicher **Motive** für die Teilnahme am SPSS-Kurs alle für Sie zutreffenden Aussagen an und/oder nennen Sie Ihre sonstigen Motive.

Ich möchte SPSS kennen lernen, ...

um eine eigene empirische Studie damit auszuwerten.  
 weil in vielen Stellenanzeigen SPSS-Kenntnisse verlangt werden.  
 weil ich mich um eine Stelle als EDV-Hilfskraft in der Forschung bewerben will (HIWI-Job).  
 weil ich mich für EDV interessiere und ein modernes Programm kennen lernen möchte.  
 weil ich mich für Statistik interessiere und mit Auswertungsverfahren experimentieren möchte.  
 Andere Motive \_\_\_\_\_

b) Möchten Sie im Kurs bestimmte statistische **Methoden** besonders gerne üben?  Ja  Nein

Wenn "Ja", welche? Kreuztabellenanalyse

9740305277

Formular drucken Daten übertragen

2 von 2

Fertig

### 3.1.1.2 Automatisches Einscannen von schriftlichen Untersuchungsdokumenten

Auch nach einer schriftlichen Befragung im konventionellen Stil lässt sich das manuelle Erfassen der Daten vermeiden. Diese lästige und fehleranfällige Arbeit kann man einer EDV-Anlage zum automatischen Einscannen und Interpretieren der schriftlichen Untersuchungsdokumente übertragen. Allerdings muss die EDV-Anlage erst mit einigem Aufwand in ihre Arbeit eingewiesen werden, so dass bei kleineren Projekten kaum ein Rationalisierungsgewinn zu erzielen ist. Eine weitere Voraussetzung dieses Verfahrens ist die Beachtung einiger Regeln beim Entwurf der Untersuchungsmaterialien. Insgesamt gesehen ist das Einscannen von Fragebögen sicher für viele Forschungsprojekte eine attraktive und rentable Erfassungsmethode.

An der Universität Trier steht für diesen Zweck im Grafikraum des Rechenzentrums (E-020) das Programm **Teleform 9.1** mit der erforderlichen Hardware (Scanner mit automatischem Einzelblatteinzug) zur Verfügung. Das Programm kann neben Markierungen in den Kästchen zu Einfach- oder Mehrfachwahlfragen (OMR) und gedruckten Zeichen (OCR) auch Handschrift lesen (ICR). Es enthält einen Formulargenerator, so dass Fragebogendesign und -deklaration in einem Arbeitsschritt erfolgen.

Eine besondere Attraktion besteht in der Möglichkeit, zu einem Teleform-Projekt ein interaktives PDF-Formular mit identischem Design zu erstellen und für eine Online-Umfrage zu verwenden. Damit können Sie entscheiden, ob Sie Ihre Daten

- mit einem gedruckten Fragebogen erheben und per Scanner erfassen,
- per Online-Umfrage erfassen (siehe oben)
- oder parallel über beide Kanäle erfassen wollen.

Das Teleform-System führt die Daten aus beiden Quellen zusammen und exportiert sie z.B. in eine SPSS-Datendatei, wobei Feldbezeichnungen und sonstige Informationen übernommen werden.

### 3.1.2 Manuelle Verfahren

Bei kleineren Studien (z.B. im Rahmen einer Diplomarbeit) dominieren noch immer die manuellen Erfassungsmethoden, wobei die Daten gemäß Kodierplan via Tastatur in einen Rechner gelangen.

Zunächst einige Empfehlungen, die für alle manuellen Erfassungsmethoden gelten:

- Schon beim Entwurf des Kodierplans ist darauf zu achten, dass dem Erfasser keine unnötigen und fehleranfälligen Arbeiten zugemutet werden (siehe oben).
- Übertragen Sie Daten von Fragebögen oder ähnlichen Untersuchungsmaterialien *direkt* in den Rechner. Das gelegentlich empfohlene Verfahren, die Daten zunächst von den Untersuchungsdokumenten auf so genannte Kodierbögen zu übertragen, um sie dann von dort endgültig zu erfassen, erhöht den Aufwand und die Fehlerwahrscheinlichkeit.

Von den möglichen manuellen Erfassungsmethoden sollen drei in diesem Manuskript vorgestellt werden:

- Erstellung einer Text-Datendatei mit einem beliebigen Texteditor  
Die Erfassung in eine Text-Datendatei hat nur einen einzigen Vorteil: Man kann sie mit fast jedem beliebigen Texteditor durchführen, z.B. auch mit dem vertrauten Textverarbeitungsprogramm. Ihr wesentlicher Nachteil ist die hohe Fehleranfälligkeit. Diese veraltete Erfassungsmethode wird hier nur beschrieben, um Sie davon abzuhalten. Allerdings gibt es noch einen zweiten Grund, das Innenleben von Text-Datendateien zu beschreiben: Es sind noch sehr viele Exemplare im Umlauf, die Sie eventuell auswerten müssen. Daher kommen wir nicht umhin, später das Einlesen von Text-Datendateien zu behandeln.



- Erfassung mit dem SPSS-Dateneditor  
Der SPSS-Dateneditor ist ein integraler Bestandteil des SPSS-Systems, so dass wir uns mit seiner Bedienung auf jeden Fall vertraut machen müssen. Er ist nicht perfekt optimiert für die Erfassung größerer Datenmengen, kann aber in kleinen bis mittleren Projekten verwendet werden.  
Relativ ähnliche Arbeitsbedingungen für die Datenerfassung bieten Tabellenkalkulationsprogramme wie z.B. MS-Excel.
- Einsatz eines speziellen Datenerfassungsprogramms  
Ein spezielles Datenerfassungsprogramm (z.B. SPSS Data Entry, MS-Access) bietet Vorteile gegenüber dem SPSS-Dateneditor, erfordert aber auch zusätzlichen Einarbeitungsaufwand.

Aufgrund des relativ geringen Datenaufkommens in unserem KFA-Projekt ist der SPSS-Dateneditor das optimale Erfassungswerkzeug. Weil in Abschnitt 3.2 die Erfassung der KFA-Daten mit dem SPSS-Dateneditor ausführlich beschrieben wird, müssen im aktuellen Abschnitt nur die beiden anderen manuellen Erfassungsmethoden vorgestellt werden.

Auch wenn das verwendete Erfassungsprogramm keine SPSS-Datendatei erzeugt, stellt die Übernahme der Daten selten ein Problem dar:

- SPSS unterstützt beim Datenimport zahlreiche Formate.
- Auf den Pool-PCs der Universität Trier steht mit dem Programm **StatTransfer** ein Konvertierungsspezialist zur Verfügung, der Dateien gängiger Datenbanken oder Statistikprogramme in das SPSS-Format übersetzen kann.

### 3.1.2.1 Erstellung einer Text-Datendatei mit einem beliebigen Texteditor

Bei dieser veralteten, zeitaufwendigen und vor allem sehr fehleranfälligen Methode muss festgelegt werden, wie die Beobachtungswerte eines Falles in der Textdatei angeordnet werden sollen. Im Wesentlichen stehen zwei Alternativen zur Auswahl: positionierte Daten und separierte Daten.

#### Positionierte Daten

In einer Datei mit fest positionierten bzw. formatierten Daten beginnt der Datensatz jedes Falles in einer neuen Datenzeile. Ferner hat jede Variable einen festen Standort im Datensatz eines Falles (z.B. in Zeile 1, Spalten 12-13). Damit sind die Datensätze aller Fälle identisch aufgebaut. So sehen die festformatig per Texteditor erfassten KFA-Daten<sup>1</sup> aus, die wir im Manuskript analysieren werden:

<sup>1</sup> Da unser Kodierplan für die Erfassung per SPSS-Dateneditor konzipiert ist, enthält er keine Zeilen-Spalten-Positionen für die Variablen. Diese wurden eigens für die Erstellung der Daten-Textdatei festgelegt. Dies geschah im Rahmen des folgenden SPSS-Programms, welches die Textdatei über das WRITE-Kommando aus der vorhandenen SPSS-Datendatei erstellt hat:

```
write outfile='kfar.txt'
/fnr 1-2 '1' geschl 5 gebj 6-7 fb 8 groesse 9-11 gewicht 12-13
/fnr 1-2 '2' aergo aergm 5-8 lot01 to lot12 10-21 motiv1 to keine 23-28
smg 30 meth1 to meth3 31-40.
exe.
```

Unser Kodierplan sieht außerdem die systematische Verwendung des MD-Indikators SYSMIS vor. Dies ist jedoch bei Text-Datendateien nicht sinnvoll. Hier sollten benutzerdefinierte MD-Indikatoren verwendet werden.

```

11 169116351
12 5 8 422125344342 100000 1 1 2 3
21 170115856
22 5 8 431224342342 100000 1 1 2 0

.
.
.

301 167117060
302 910 551115443131 100000 0 0 0 0
311 167116968
312 7 9 412544231132 100010 1 1 3 0

```

### Separierte Daten

In einer Datei mit separierten Daten müssen die Variablenausprägungen jedes Falles in derselben Reihenfolge vorliegen, und je zwei Werte müssen durch ein Separatorzeichen voneinander getrennt werden. Beim Trennzeichen hat man die freie Auswahl, entscheidet sich aber meist zwischen folgenden Kandidaten:

- Tabulatorzeichen
- Komma
- Semikolon
- Leerzeichen

Beim Einlesen separierter Daten durch SPSS wird eine Serie aufeinander folgender Leerzeichen behandelt wie ein einzelnes Leerzeichen. Ansonsten schließen zwei aufeinander folgende Trennzeichen einen fehlenden Wert ein, den SPSS beim Einlesen durch SYSMIS kodiert.

Obwohl nicht zwingend vorgeschrieben, sollte man alle Daten eines Falles in eine einzige Zeile schreiben und für jeden Fall eine neue Zeile beginnen.

In der ersten Zeile einer Textdatei mit separierten Daten können die Variablenamen an SPSS übergeben werden, was im folgenden Beispiel mit Tabulator-separierten KFA-Daten demonstriert wird:

FNR	GESCHL	GEBJ	FB	GROESSE	GEWICHT	AERGO	AERGM	.	.	.
1	1	69	1	163	51	5	8	.	.	.
2	1	70	1	158	56	5	8	.	.	.
3	1	69	1	174	58	4	8	.	.	.
4	2	67	1	182	77	6	2	.	.	.
5	1	67	1	180	69	8	8	.	.	.
.	.	.	.	.	.	.	.	.	.	.
.	.	.	.	.	.	.	.	.	.	.
29	1	68	1	176	63	7	9	.	.	.
30	1	67	1	170	60	9	10	.	.	.
31	1	67	1	169	68	7	9	.	.	.

#### 3.1.2.2 Einsatz eines speziellen Datenerfassungsprogramms

Wenn bei *größeren* Projekten eine manuelle Datenerfassung unumgänglich ist (vgl. Abschnitt 3.1.1), dann sollte in der Regel ein spezielles Datenbankprogramm verwendet werden. Man arbeitet hier bequem mit einer *Erfassungsmaske*, die einen *einzelnen* Fall in übersichtlicher Form auf dem Bildschirm präsentiert. Zudem werden die eingegebenen Daten in der Regel *sofort auf Plausibilität überprüft*: Falsche Eingaben werden mit entsprechendem Protest abgewiesen.

Ein Nachteil dieser Methode besteht im hohen Aufwand:

- Es muss ein spezielles Programm erlernt werden.
- Für jedes Projekt sind einige Konfigurationsarbeiten erforderlich (z.B. Gestaltung der Erfassungsmaske, Definition der Regeln zur Plausibilitätskontrolle)

Anschließend werden zwei Spezialprogramme zur Datenerfassung exemplarisch beschrieben. Sofern ein Arbeitsplatz mit permanenter Internet-Verbindung zur Verfügung steht, kann übrigens auch ein Online-Umfragesystem für die manuelle Dateneingabe mit Erfassungsmaske und Plausibilitätskontrolle eingesetzt werden.

### a) SPSS Data Entry

Mit Data Entry kann man eine analog zum Fragebogen aufgebaute Eingabemaske entwerfen, um dem Erfasser die Orientierung zu erleichtern, z.B.:

Allerdings sollte der Erfasser keine Zeit mit der Mausbedienung verlieren, sondern die erforderlichen Tastaturbefehle zur Bearbeitung des Formulars beherrschen.

Mit den folgenden Vorzügen macht Data Entry die Erfassung rationeller und sicherer:

- **Filterfragen (Skip & Fill)**  
In Abhängigkeit vom erfassten Wert einer Filtervariablen verzweigt Data Entry zu unterschiedlichen Folgevariablen und versorgt dabei übersprungene Variablen mit einem festgelegten MD-Indikator.
- **Plausibilitätsprüfungen**  
Man kann z.B. dafür sorgen, dass bei der Variablen GESCHL nur die Werte 1, 2 und SYSMIS akzeptiert werden.

Neben der *Datenerfassung* will Data Entry auch das *Fragebogendesign* unterstützen. Man kann entweder *ein* Formular zur Verwendung bei der Datenerhebung (z.B. durch schriftliche Befragung) *und* bei der EDV-Erfassung entwerfen, oder für beide Anwendungsfälle angepasste Formulare verwenden. Dazu bietet Data Entry Beispielfragebögen bzw. Musterbibliotheken (z.B. mit demographischen Fragen) an.

Weitere Funktionen von Data Entry sind:

- Existierende SPSS-Datendateien auf Fehler prüfen
- Einen Fragebogen zu einer existierenden SPSS-Datendatei erstellen

Eine ausführliche Beschreibung zu Data Entry finden Sie auf dem WWW-Server der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend über:

[Weitere Service-Angebote > EDV-Dokumentationen > Elektronische Publikationen > Datenerfassung](#)

## b) INPUT II

Ein anderer Weg zum maßgeschneiderten Datenbankprogramm mit maskengesteuerter Dateneingabe, Plausibilitätskontrolle und Filterführung ist die Verwendung des Programmgenerators **INPUT II**, der an der Universität Trier für die speziellen Bedürfnisse wirtschafts- und sozialwissenschaftlicher Forschungsarbeit mit SPSS entwickelt wurde. Der (nur im Campusnetz der Universität Trier verfügbare) Generator erzeugt aus einer Datensatzbeschreibung ein spezialisiertes Erfassungsprogramm, das dann ohne Lizenzgebühren auf jedem beliebigen PC (z.B. auch zu Hause) eingesetzt werden kann. Das DOS-basierte Erfassungsprogramm ist zwar nicht mehr ganz auf der Höhe moderner Softwaretechnik, begnügt sich dafür aber auch mit einer minimalen Hardwareausstattung.

Eine INPUT II - Beschreibung finden Sie auf dem WWW-Server der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend über:

[Weitere Service-Angebote > EDV-Dokumentationen > Elektronische Publikationen > Datenerfassung](#)

## 3.2 Erfassung mit dem SPSS-Dateneditor

Für die nächsten Schritte im KFA-Projekt benötigen Sie eine SPSS-Sitzung mit einem leeren Datenfenster. Diese Situation liegt z.B. vor, nachdem Sie SPSS gestartet und den Startassistenten mit dem Ziel **Daten eingeben** verlassen haben. Nötigenfalls können Sie ein leeres Datenfenster mit dem folgenden Menübefehl anfordern:

**Datei > Neu > Daten**

Im realen SPSS-Kurs werden wir nun mit dem SPSS-Dateneditor unsere Variablen deklarieren und anschließend die Daten erfassen.

Wenn Sie dieses Manuskript im Selbststudium lesen, können und sollten Sie trotzdem die folgenden Arbeitsschritte zur Variablendeklaration konkret nachvollziehen und die Daten des im Manuskript abgedruckten ersten Falles eintragen (siehe Seite 25). Alle Projektphasen nach der Datenerfassung können Sie durch Verwendung der SPSS-Datendatei **kfar.sav** mitmachen, deren Inhalt im weiteren Verlauf erklärt wird. Wie Sie diese Datei von einem Server des Rechenzentrums beziehen können, ist im Vorwort zu erfahren.

### 3.2.1 Dateneditor und Arbeitsdatei

Wir haben schon in Abschnitt 1.4.1 festgestellt, dass über das Dateneditorfenster<sup>1</sup> die rechteckige (Fälle × Variablen) - Datenmatrix zugänglich ist. SPSS speichert die Daten während der Sitzung in einer temporären Datei, bezeichnet als **Arbeitsdatei** oder **Arbeitsdatendatei**, die nach Möglichkeit im Hauptspeicher des PCs gehalten wird. Die im Dateneditorfenster sicht- und modifizierbare Arbeitsdatei enthält:

---

<sup>1</sup> Wie Sie sicher schon bemerkt haben, wird im Manuskript gelegentlich für *Dateneditorfenster* die kürzere Bezeichnung *Datenfenster* verwendet.

- Die **rechteckige (Fälle × Variablen)-Datenmatrix**  
Wir wollen statistische und graphische Analysen für die Variablen anfordern, d.h. für die Spalten der (Fälle × Variablen)-Datenmatrix in der Arbeitsdatei. Dazu ist jede Variable über ihren eindeutigen Variablennamen ansprechbar.
- Einen so genannten **Deklarationsteil**  
Dort merkt sich SPSS verarbeitungsrelevante Merkmale der Variablen (z.B. MD-Indikatoren). Über die **Variablenansicht** des Datenfensters (siehe unten) können Sie die Merkmale der Variablen jederzeit einsehen und ändern.

Mit Hilfe des Dateneditors oder durch Transformationskommandos (siehe unten) können während einer Sitzung u.a. folgende Modifikationen der Arbeitsdatei vorgenommen werden:

- Erweiterung um neue Variablen
- Änderung von Variablenattributen (z.B. Namen, MD-Indikatoren)
- Löschen von Variablen
- Erweiterung um neue Fälle
- Änderung von Variablenausprägungen eines Falles
- Löschen von Fällen

Weil die Begriffe **Dateneditor** und **Arbeitsdatei** für den Umgang mit SPSS recht wichtig sind, sollen ihre wesentlichen Eigenschaften noch einmal wiederholt werden:

- Die Arbeitsdatei enthält die Datenmatrix und den Deklarationsteil.
- Mit dem Dateneditor können wir die Arbeitsdatei ansehen und modifizieren, auf dem Registerblatt **Datenansicht** die Datenmatrix und auf dem Registerblatt **Variablenansicht** den Deklarationsteil.
- Die Arbeitsdatei ist temporär, muss also nach einer (planvollen) Änderung in eine permanente SPSS-Datendatei gesichert werden (siehe unten).

### 3.2.2 Variablen definieren

Wie eben erwähnt, verwaltet SPSS für jede Variable zahlreiche verarbeitungsrelevante Merkmale (z.B. MD-Indikatoren). Diese werden im Deklarationsteil der Arbeitsdatei gespeichert und können vom Anwender festgelegt werden. Da SPSS für alle Attribute geeignete Voreinstellungen benutzt, setzt die Datenerfassung nicht unbedingt eine Variablendefinition voraus<sup>1</sup>, doch wird das Erfassen und die spätere Auswertungsarbeit z.B. durch benutzerdefinierte Variablennamen anstelle der automatisch generierten und wenig aussagekräftigen Namen VAR00001, VAR00002 usw. erleichtert. Daher liegt es nahe, dem SPSS-System die in unserem Kodierplan beschriebenen Variablen vor dem Eintragen der Daten bekannt zu machen.

#### 3.2.2.1 Das Datenfenster-Registerblatt Variablenansicht

Das Datenfenster besitzt *zwei* Registerblätter bzw. Tabellen:

- das Registerblatt **Datenansicht** zur Anzeige und Modifikation der (Fälle × Variablen)-Datenmatrix
- das Registerblatt **Variablenansicht** zur Anzeige und Modifikation der Variablenattribute

---

<sup>1</sup> Da in SPSS der Variablentyp *numerisch* voreingestellt ist, müssten wir vor dem Erfassen von Daten anderen Typs auf jeden Fall eine Variablendefinition vornehmen. Allerdings sind solche Variablen in unserem Kodierplan nicht vorgesehen.

In den Zeilen der **Variablenansicht** wird jeweils eine Variable beschrieben, wozu in den Spalten insgesamt zehn Attribute zur Verfügung stehen. Für unsere erste Variable (FNR) eignen sich z.B. folgende Angaben:

	Name	Typ	Spaltenformat	Dezimalstellen	Variablenlabel	Wertelabels	Fehlende Wert	Spalten	Ausrichtung	Maßniveau
1	fnr	Numerisch	2	0		Kein	Kein	3	Rechts	Nominal
2										

Um eine neue Variable anzulegen, trägt man ihren Namen in eine freie Zeile der Tabelle ein und ändert nach Bedarf die automatisch generierten Attributausprägungen. Darüber hinaus können auch Variablen eingefügt, gelöscht oder verschoben werden (siehe unten).

### 3.2.2.2 Die SPSS-Variablenattribute

Bevor wir die Variablen unserer KFA-Studie deklarieren, sollen vorab die SPSS-Variablenattribute erläutert werden:

- **Name**

Die wesentlichen Regeln für SPSS-Variablenamen wurden schon im Zusammenhang mit dem Kodierplan genannt (siehe Seite 23).

- **Typ**

Die wichtigsten SPSS-Variablentypen haben wir schon genannt: Numerisch, String und Datum (siehe Seite 18). In der Regel empfiehlt es sich, auch bei nominalskalierten Merkmalen eine numerische Kodierung vorzunehmen (siehe Abschnitt 1.4.3), so dass der voreingestellte numerische Variablentyp meist beibehalten werden kann.

- **Spaltenformat**

Bei einer *numerischen* Variablen beeinflusst dieses Attribut lediglich ihre voreingestellte Breite bei der Ausgabe in eine Textdatendatei über das Kommando WRITE (inkl. Vorzeichen und Dezimaltrennzeichen) und ist daher für die Arbeit mit dem Daten- und dem Ausgabefenster wenig relevant. Allerdings muss der Spaltenformatwert stets größer sein als die Anzahl der Dezimalstellen (s. u.).

Bei einer *alphanumerischen* Variablen legt das Spaltenformat die maximale Anzahl der gespeicherten Zeichen fest und ist folglich recht bedeutsam. So werden z.B. bei einer nachträglichen Reduktion der Spaltenzahl tatsächlich entsprechend viele Zeichen am rechten Rand gelöscht.

- **Dezimalstellen**

Bei einer *numerischen* Variablen können Sie festlegen, welche Anzahl von Dezimalstellen bei der Anzeige ihrer Werte im Datenfenster bzw. in der Ergebnisausgabe verwendet werden soll.

Diese Angabe betrifft *nicht* die Speichergenauigkeit im Datenfenster bzw. in der Arbeitsdatei, sondern nur die Anzeige.

Bei einer *alphanumerischen* Variablen ist das Attribut irrelevant und auf den Wert 0 fixiert.

- **Variablenlabel**

Hier können optional Variablenlabel mit einer maximalen Länge von 256 Zeichen vereinbart werden, die in Ergebnistabellen und Grafiken an Stelle der aus praktischen Erwägungen möglichst kurz gewählten und mit Syntaxrestriktionen belasteten Variablennamen (Verbot von Leerstellen) angezeigt werden sollen, z.B.:

Variablenname	Variablenlabel
FB	Fachbereich an der Universität Trier
AERGM	Ärger mit kontrafaktischer Alternative

Allerdings erscheinen die Labels in der Ausgabe mancher SPSS-Prozeduren nicht in voller Länge.

Umlaute und sonstige Sonderzeichen sind erlaubt, und die potentiellen Probleme bei Verwendung einer Datendatei unter alternativen Betriebssystemen sind weniger ernst, weil es sich beim Variablenlabel ja um ein optionales Attribut handelt.

Während wir die Variablennamen in SPSS der Einfachheit halber stets klein schreiben, ist bei den Variablenlabels eine publikationsreife Groß/Kleinschreibung angemessen.

Sind Variablenlabel vorhanden, werden diese auch in Dialogboxen zur Beschreibung der Variablen verwendet. Diese Voreinstellung kann aber über

**Bearbeiten > Optionen > Allgemein > Variablenlisten = Namen anzeigen**

abgeändert werden. Bei der in Dialogboxen üblichen Platzbeschränkung auf ca. 20 Stellen ist oft der abgeschnittene Anfang eines 50-stelligen Labels weniger informativ als der vollständige Name.

- **Wertelabels**

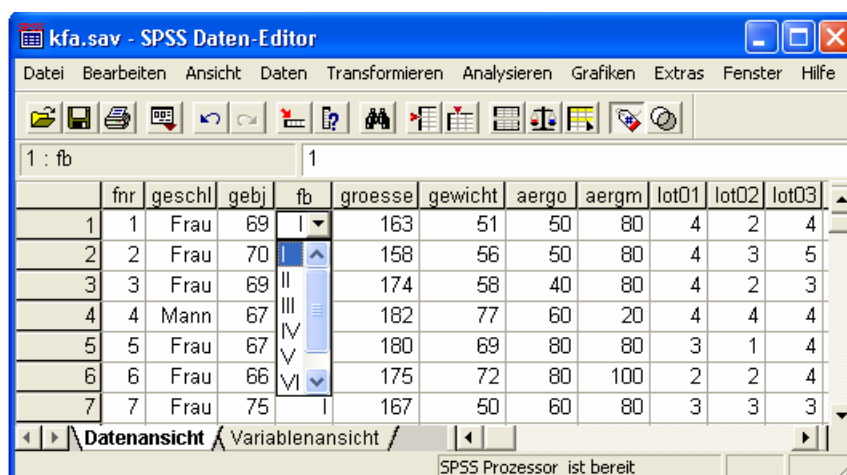
Hier können optional Wertelabels mit maximal 60 Zeichen zur Erläuterung von Variablenausprägungen vereinbart werden, was speziell bei numerisch kodierten nominalskalierten Merkmalen empfehlenswert ist, z.B.:

Variablenname	Werte	Wertelabels
GESCHL	1	Frau
	2	Mann

Diese Labels spielen bei Berechnungen keine Rolle, erscheinen aber in der Ergebnisausgabe und können deren Lesbarkeit verbessern.

Außerdem bietet auch die Datenansicht des Dateneditors nach dem Menübefehl **Ansicht > Wertelabels** einige Unterstützung für die Etiketten:

- Sie werden an Stelle der Werte angezeigt.
- Alternativ zur Werteingabe per Tastatur kann man per Drop-Down-Menü ein Label wählen:




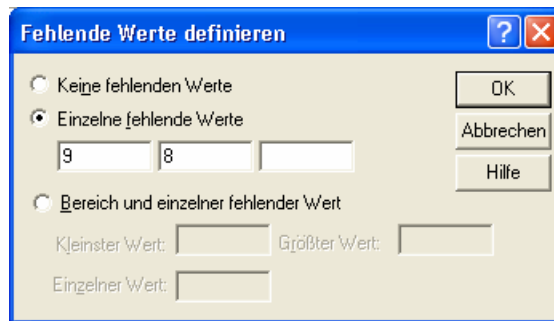
Die starke optische Präsenz der Wertelabels führt bei vielen SPSS-Anwendern zu einer Überschätzung ihrer Rolle bei der Datenverarbeitung: Es ist *nicht* möglich, durch Vergabe von Wertelabels die Menge der gültigen Werte einer Variablen zu definieren und eine Plausibilitätskontrolle für die Erfassung per Dateneditor einzurichten. Unter obiger Vereinbarung für die Variable GESCHL führt z.B. die Eingabe des Wertes 4711 zu keinerlei Protest durch SPSS.

- **Fehlende Werte**

Wenn Sie bei einer Variablen *benutzerdefinierte* MD-Indikatoren verwenden wollen, müssen Sie diese unbedingt deklarieren, weil sie sonst wie gültige Werte behandelt werden, z.B. bei einer Mittelwertbildung.

Da wir im KFA-Projekt laut Kodierplan ausschließlich System-Missing als MD-Indikator verwenden, müssen wir anschließend keine MD-Deklaration vornehmen (vgl. Abschnitt 1.4.3.5). Daher wird an dieser Stelle die simple Prozedur zum Deklarieren von benutzerdefinierten MD-Indikatoren beschrieben:

- Markieren Sie bei der betroffenen Variablen die Zelle zum Attribut **Fehlende Werte**.
- Nach einem Mausklick auf den nun vorhandenen Erweiterungsschalter  erscheint eine Dialogbox, in der man entweder bis zu drei Einzelwerte oder aber ein Intervall samt zusätzlichem Einzelwert als MD-Indikatoren vereinbaren kann, z.B.:



- **Spalten und Ausrichtung**

Wie breit soll die Spalte einer Variablen im Datenfenster sein? Wie sollen die Werte ausgerichtet werden (linksbündig, zentriert, rechtsbündig)? Die Attribute in dieser Subdialogbox wirken sich nur auf die Darstellung einer Variablen im Datenfenster aus.

Reicht der erlaubte Platz für die vollständige Darstellung eines Wertes nicht aus, erscheinen entsprechend viele Sternchen.

- **Messniveau**

Über die technischen Variablenattribute hinaus kann das Messniveau einer Variablen deklariert werden, wobei diese Vereinbarung bei der weiteren Arbeit mit SPSS allerdings bisher nur in wenigen Situationen relevant ist, z.B.:

- Bei der so genannten interaktiven Grafik (siehe unten)
- Bei der Verarbeitung von SPSS-Datendateien mit **Answer Tree**

In Zukunft werden wohl mehr SPSS-Prozeduren die Information über das Messniveau der Variablen ausnutzen.

Weil außerdem die Reflexion über dieses methodologisch wichtige Variablenattribut nicht schaden kann, wollen wir uns in diesem Kurs der Pflicht unterziehen, bei allen Variablen das korrekte Messniveau anzugeben.

### 3.2.2.3 Variablendefinition durchführen

Aktivieren Sie nun die **Variablenansicht** des Datenfensters, und tragen Sie für die erste Variable (zur Fallidentifikation) den Namen FNR ein. Nach dem Markieren der zugehörigen Zelle können Sie sofort mit dem Eintippen des Namens beginnen. Die Groß/Kleinschreibung ist dabei irrelevant. Im Manuskript werden Variablennamen nur aus darstellungstechnischen Gründen groß geschrieben.

Sobald Sie die Zelle mit dem Variablennamen verlassen (z.B. per Mausklick auf eine andere Zelle oder per Tabulatortaste) wird eine neue Variable mit dem gewünschten Namen in die Ar-




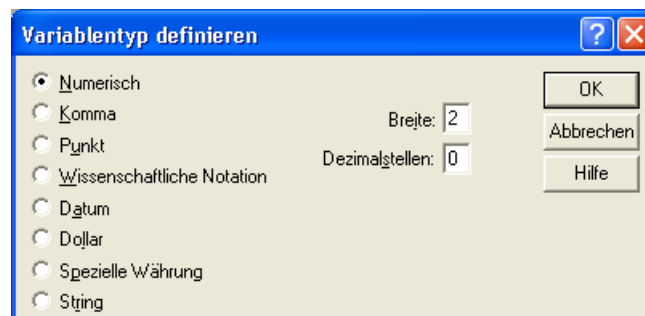
beitsdatei aufgenommen, sofern gegen den Variablennamen keine Einwände bestehen, und die restlichen Attribute der neuen Variablen werden mit Standardwerten versorgt.

Nach dem Markieren der Zelle **Dezimalstellen** kann man die gewünschte Anzahl von Dezimalstellen durch Eingabe einer Zahl oder per Up-Down - Regler wählen:

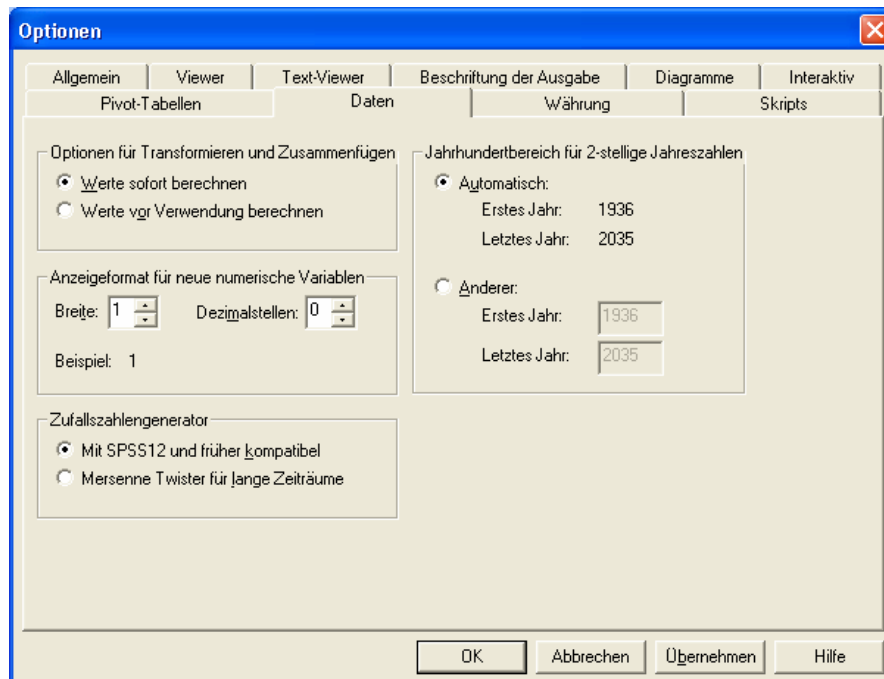


Analog wird auch das Attribut **Spaltenformat** festgelegt, das allerdings bei der von uns geplanten Arbeitsweise keine große Rolle spielt (siehe oben).

Eine alternative Möglichkeit zum Einstellen der Attribute **Dezimalstellen** und **Spaltenformat** findet sich in der (von uns eigentlich nicht benötigten) Dialogbox **Variablentyp definieren**, die nach einem Mausklick auf den Erweiterungsschalter  in der markierten **Typ**-Zelle erscheint:

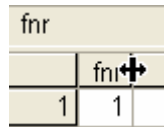


**Tipp:** Wenn in einem Projekt das voreingestellte Anzeigeformat für numerische Variablen (Breite = 8, Dezimalstellen = 2) häufig durch eine bestimmte Alternative ersetzt werden muss, kann zur Vereinfachung der Deklaration die Voreinstellung entsprechend geändert werden. Dazu öffnet man mit **Bearbeiten > Optionen** die Dialogbox **Optionen**, wechselt hier zum Registerblatt **Daten** und nimmt im Rahmen **Anzeigeformat für neue numerische Variablen** die gewünschten Einstellungen vor, z.B.:



Wenngleich die Variable FNR im Ausgabefenster nicht allzu oft auftauchen wird, tragen wir in die Zelle zum Attribut **Variablenlabel** den Text *Fallnummer* ein.

Statt die Breite der FNR-Spalte im Datenfenster über eine gut geschätzte **Spalten**-Angabe festzulegen, können Sie bei aktiviertem Datenfenster-Registerblatt **Datenansicht** auch folgendermaßen vorgehen: Setzen Sie den Mauszeiger auf den rechten Rand der Zelle mit dem Variablenamen, woraufhin der Zeiger eine neue Form und dementsprechend eine neue Funktion gewinnt:



Nun lässt sich der rechte Rand der aktuellen Spalte verschieben: Linke Maustaste drücken, ziehen und an der gewünschten Position wieder los lassen. Eine so festgelegte Spaltenbreite wird von SPSS als **Spalten**-Variablenattribut übernommen.

Klappen Sie in der markierten **Messniveau**-Zelle die versteckte Liste auf, um für die Fallnummer ein nominales Skalenniveau zu deklarieren:




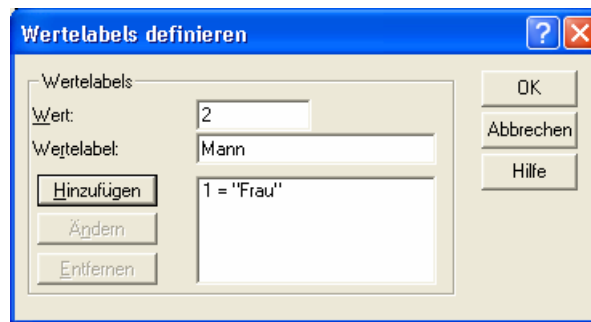
MD-Indikatoren müssen wir im KFA-Projekt generell nicht vereinbaren, Wertelabels sind bei der augenblicklich bearbeiteten Fallnummernvariablen irrelevant, und das Attribut **Ausrichtung** übernehmen wir stets unverändert. Daher können wir die Deklaration der Variablen FNR beenden:



Bei Bedarf sind Anpassungen jederzeit möglich.

Vereinbaren Sie nun in der zweiten Zeile der Variablenansicht für die Geschlechtsvariable den Namen GESCHL, eine einspaltige Anzeige ohne Dezimalstellen und das Variablenlabel *Geschlecht*.

Bei diesem numerisch kodierten nominalskalierten Merkmal ist es sinnvoll, die willkürliche Zuweisung von Zahlen zu den beiden Kategorien durch Wertelabels zu dokumentieren, damit wir bei der Lektüre von Ergebnisausgaben nicht rätseln müssen, welches Geschlecht die Nummer Eins ist. Öffnen Sie daher mit einem Mausklick auf den Erweiterungsschalter  in der markierten **Wertelabels**-Zelle die folgende Dialogbox:



Hier wird z.B. das weibliche Label folgendermaßen vereinbart:

- Tragen Sie den **Wert** 1 und das **Wertelabel** *Frau* ein.
- Drücken Sie auf den Schalter **Hinzufügen**.  
In der Schaltflächen-Beschriftung **Hinzufügen** signalisiert nach Betätigen der **Alt**-Taste das unterstrichene **H**, dass der umständlichen Mausklick auf die Schaltfläche durch die Tastenkombination **Alt+H** ersetzt werden kann.

Abschließend ist für GESCHL noch das nominale **Messniveau** zu deklarieren.

### 3.2.2.4 Übung

Definieren Sie alle Variablen zur ersten Seite unseres KFA-Fragebogens. Wie Sie nötigenfalls Variablen einfügen oder löschen können, erfahren Sie im nächsten Abschnitt.

## 3.2.3 Variablen einfügen, löschen oder verschieben

Bei der Variablendefinition kann sich leicht die Notwendigkeit ergeben, Variablen einzufügen oder zu löschen.

### 3.2.3.1 Variablen einfügen

Wenn Sie z.B. nach FNR und GESCHL die Variable FB definiert und folglich die Variable GEBJ vergessen haben, können Sie das Missgeschick in der Variablenansicht folgendermaßen korrigieren:

- Setzen Sie einen rechten Mausklick auf die Nummer der FB-Zeile (am linken Rand der Tabelle).
- Wählen Sie die Option **Variable einfügen** aus dem Kontextmenü.

Daraufhin stellt SPSS vor FB eine neue Variable mit voreingestellten Attributen zur Verfügung, die nun beliebig angepasst werden können:

	Name	Typ	Spaltenformat	Dezimalstellen	Variablenlabel	Wertelabels	Fehlende Werte	Spalten	Ausrichtung	Meßniveau
1	fnr	Numerisch	2	0	Fallnummer	Kein	Kein	3	Rechts	Nominal
2	geschl	Numerisch	1	0	Geschlecht	{1, Frau}...	Kein	7	Rechts	Nominal
3	VAR00005	Numerisch	8	2		Kein	Kein	8	Rechts	Metrisch
4	fb	Numerisch	1	0	Fachbereich an d	{1, I}...	Kein	3	Rechts	Metrisch
5										

Auf analoge Weise lässt sich eine neue Variable auch in der *Datenansicht* einfügen:

- Setzen Sie einen rechten Mausklick auf die Beschriftung der FB-Spalte im Kopfbereich der Tabelle,
- und wählen Sie die Option **Variable einfügen** aus dem Kontextmenü.

### 3.2.3.2 Variablen löschen

Gehen Sie in der Variablenansicht folgendermaßen vor, um eine Variable zu löschen:

- Setzen Sie einen rechten Mausklick auf die Zeilennummer der betroffenen Variablen (am linken Rand der Tabelle).
- Wählen Sie die aus dem Kontextmenü Option **Löschen**.

Auf analoge Weise lässt sich eine Variable auch in der Datenansicht löschen.

### 3.2.3.3 Variablen verschieben

Gehen Sie in der Variablenansicht folgendermaßen vor, um eine Variable per Drag & Drop (Ziehen und Ablegen) zu verschieben:

- Markieren Sie die zu verschiebende Variable durch einen Mausklick auf ihre Zeilennummer. Lassen Sie anschließend die Maustaste wieder los.
- Klicken Sie erneut auf die Nummer der zu verschiebenden Variablen, und halten Sie dabei die Maustaste gedrückt.
- Bewegen Sie bei gedrückter Maustaste den Mauszeiger zum Ziel der Verschiebungsaktion. Der aktuell anvisierte Zielort wird von SPSS durch eine rote Linie gekennzeichnet.
- Wenn Sie die Maustaste los lassen, erscheint die Variable am neuen Ort.

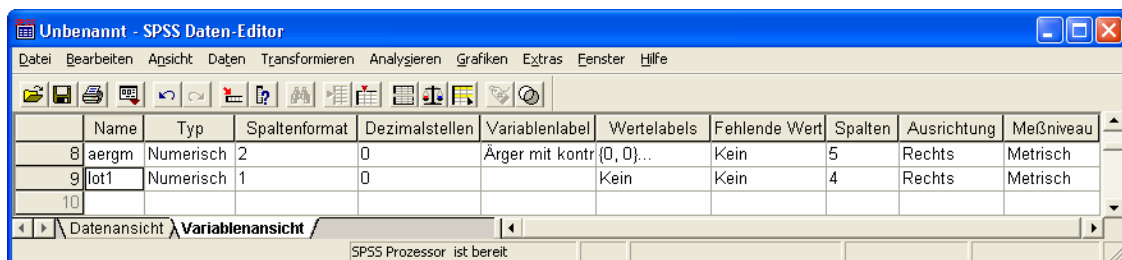
Auf analoge Weise lässt sich eine Variable auch in der Datenansicht verschieben.

## 3.2.4 Attribute auf andere Variablen übertragen

### 3.2.4.1 Variablendeklarationen vervielfältigen

Für unsere 12 LOT-Fragen sollen natürlich alle Variablenattribute mit Ausnahme des Namens identisch sein. Erfreulicherweise müssen wir die identische Variablendefinition nicht 12-mal wiederholen, sondern können nach einer ersten Definition die Attribute auf alle anderen Variablen übertragen. Mit der folgenden Vorgehensweise lässt sich sogar das Schreiben der restlichen Variablennamen automatisieren:

- Deklarieren Sie die Variable LOT1 mit geeigneten Attributen, z.B.:



The screenshot shows the SPSS Data Editor window with the Variable View selected. The table below represents the data shown in the screenshot:

	Name	Typ	Spaltenformat	Dezimalstellen	Variablenlabel	Wertelabels	Fehlende Wert	Spalten	Ausrichtung	Meßniveau
8	aergm	Numerisch	2	0	Ärger mit kontr	{0, 0}...	Kein	5	Rechts	Metrisch
9	lot1	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
10										

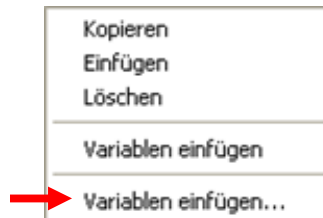
Das voreingestellte metrische Messniveau kann beibehalten werden, obwohl unsere fünf-stufigen Variablen LOT1 bis LOT12 wohl eher grobschlächlige Indikatoren für das angenommene latente Merkmal Optimismus sind. In den geplanten Auswertungen werden wir nicht die zwölf Rohvariablen selbst, sondern eine daraus abgeleitete Mittelwertvariable verwenden, für die ein approximativ metrisches Messniveau angenommen werden darf.

- Markieren Sie die komplette Variable LOT1 per Mausklick auf ihre Zeilennummer am linken Tabellenrand, und kopieren Sie alle Attribute mit **Strg+C** oder

### Bearbeiten > Kopieren

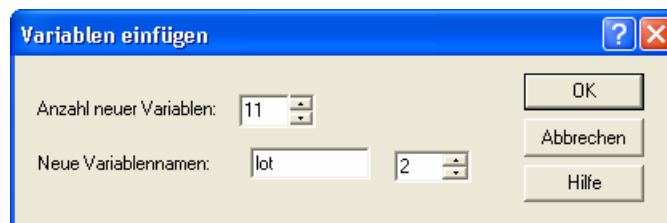
in die Zwischenablage.

- Setzen Sie einen rechten Mausklick auf die nächste freie Zeile der Variablenansicht und wählen Sie aus dem Kontextmenü die Option **Variablen einfügen** mit den drei Punkten am Ende der Beschriftung:



Diese Option ist nur verfügbar, wenn sich eine komplette Variablendeklaration in der Zwischenablage befindet.

- In der folgenden Dialogbox



können Sie nun festlegen, ...

- wie viele neue Variablen benötigt werden,
- welche gemeinsame Wurzel die neuen Variablennamen haben sollen,
- mit welchem Indexwert SPSS den Namen der ersten Variablen komplettieren soll.

Nach dem Quittieren der obigen Dialogbox entstehen elf neue Variablen mit den gewünschten Namen und Attributen:

	Name	Typ	Spaltenformat	Dezimalstellen	Variablenlabel	Wertelabels	Fehlende Wert	Spalten	Ausrichtung	Meßniveau
8	aergm	Numerisch	2	0	Ärger mit kontr	{0, 0}...	Kein	5	Rechts	Metrisch
9	lot1	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
10	lot2	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
11	lot3	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
12	lot4	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
13	lot5	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
14	lot6	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
15	lot7	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
16	lot8	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
17	lot9	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
18	lot10	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
19	lot11	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch
20	lot12	Numerisch	1	0		Kein	Kein	4	Rechts	Metrisch

### 3.2.4.2 *Alle Attribute einer Variablen übertragen*

Gehen Sie folgendermaßen vor, um alle Attribute einer Variablen (mit Ausnahme des Namens) auf andere, bereits vorhandene Variablen zu übertragen:

- Markieren Sie die Quellvariable per Mausklick auf ihre Zeilennummer am linken Tabellenrand, und kopieren Sie alle Attribute mit **Strg+C** oder

#### **Bearbeiten > Kopieren**

in die Zwischenablage.

- Markieren Sie *eine* Zielvariable per Mausklick auf ihre Zeilennummer oder eine Serie von Zielvariablen durch Mausklicks in Kombination mit der Umschalt- oder **Strg**-Taste.
- Übertragen Sie die in der Zwischenablage gespeicherten Attribute auf alle markierten Variablen mit **Strg+V** oder

#### **Bearbeiten > Einfügen**

### 3.2.4.3 *Einzelne Attribute einer Variablen übertragen*

Es ist auch möglich, ein *einzelnes* Attribut von einer Variablen auf andere zu übertragen:

- Quell-Attributzelle markieren
- Attribut mit **Strg+C** in die Zwischenablage kopieren
- Zu verändernde Attributzellen markieren
- Attribut mit **Strg+V** aus der Zwischenablage übernehmen

### 3.2.4.4 *Übung*

Definieren Sie die restlichen Variablen unserer KFA-Studie.

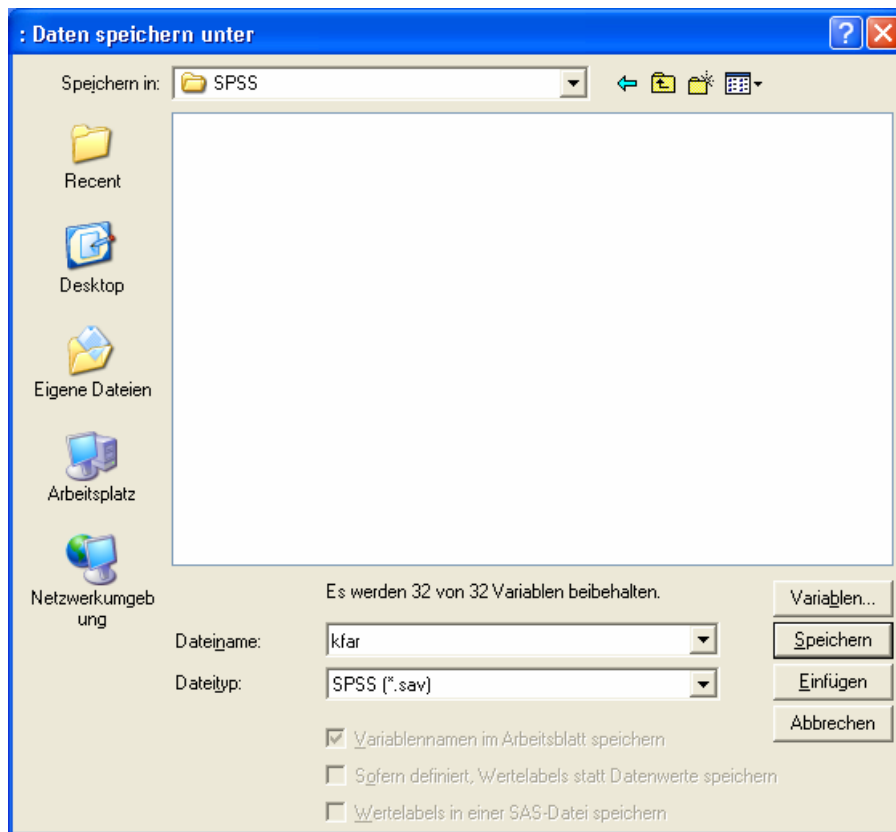
### 3.2.5 *Sichern der Arbeitsdatei als SPSS-Datendatei*

Wenn eine neu erstellte Arbeitsdatei über das Ende der Sitzung hinaus erhalten bleiben soll, muss sie explizit auf einen permanenten Datenträger gesichert werden. Dabei entsteht eine **SPSS-Datendatei**, früher auch als *SPSS-Systemdatei* bezeichnet. In späteren Sitzungen kann durch *Öffnen* dieser SPSS-Datendatei der gesicherte Zustand der Arbeitsdatei wieder hergestellt werden.

Zwar enthält Ihre aktuelle Arbeitsdatei noch keine Daten, aber im Deklarationsteil stehen bereits wertvolle Informationen, deren Verlust recht schmerzlich wäre. Daher sollten Sie schon jetzt die temporäre Arbeitsdatei in eine permanente SPSS-Datendatei sichern, indem Sie den folgenden Menübefehl wählen:

#### **Datei > Speichern unter...**

In der erscheinenden Dialogbox ist für die zu erzeugende SPSS-Datendatei ein Name, ein Verzeichnis und ein Laufwerk anzugeben:



Wenn Sie die für SPSS-Dateien vorgegebene Namenserweiterung **.sav** beibehalten, geht das spätere Öffnen besonders bequem.


Als Name für unsere Beispieldatei wird **kfar.sav** vorgeschlagen, verbunden mit der Versicherung, die Begründung für das *r* im nächsten Abschnitt nachzuliefern.

Wenn Sie an einem Pool-PC an der Universität Trier arbeiten, können Sie die Datei im Ordner **U:\Eigene Dateien\SPSS** speichern, der beim ersten SPSS-Einsatz automatisch angelegt wurde. Nach dem **Speichern** zeigt die Titelzeile des Datenfensters den Namen der nunmehr zugeordneten Datendatei, in unserem Fall also **kfar.sav**.

Beim Speichern einer Arbeitsdatei können auch alternative Dateiformate gewählt werden (z.B. MS-EXCEL, SAS, ASCII, dBase).

Zum späteren Sichern in eine bereits zugeordnete Datei dient der Befehl:

### **Datei > Speichern**

Alternativ können Sie mit der Maus auf das Symbol  klicken oder die Tastenkombination **Strg+S** benutzen.

### **3.2.6 Rohdatendatei, Transformationsprogramm und Fertigdatendatei**

Möglicherweise haben Sie sich beim Lesen des letzten Abschnitts gefragt, was das *r* im vorgeschlagenen Dateinamen **kfar.sav** bedeuten soll. Bei der Beantwortung dieser Frage sind leider einige Vorgriffe auf spätere Abschnitte nötig. Versuchen wir es trotzdem. Das *r* soll signalisieren, dass in dieser Datei die nach den Vorschriften des Kodierplans erfassten **Rohdaten** stehen. In **kfar.sav** sollen also ausschließlich folgende Arbeitsschritte einfließen:

- Variablendefinition gemäß Kodierplan
- Datenerfassung gemäß Kodierplan
- Nötigenfalls spätere Korrekturen von Erfassungsfehlern

Damit ist diese Datei für viele im Demoprojekt geplante Auswertungsschritte noch nicht geeignet. Es fehlt z.B. der Optimismus-Testwert, welcher aus den zwölf LOT-Fragen berechnet werden muss.

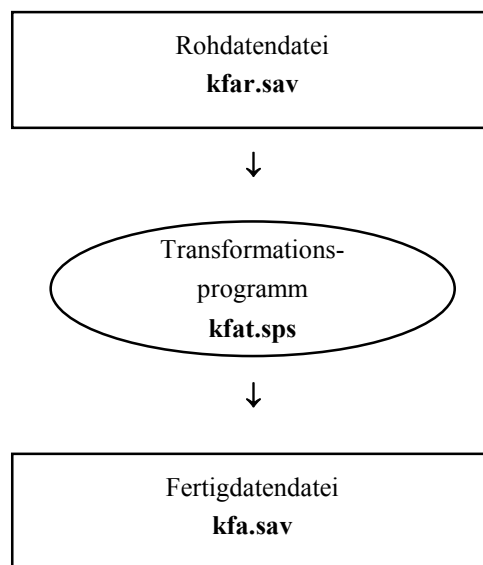
Aus der Rohdatendatei werden wir bald eine **Fertigdatendatei** herstellen, in die alle projektweit relevanten Variablenmodifikationen und -neuberechnungen einfließen sollen, so dass sie eine bequeme Datenbasis für alle statistischen und graphischen Analysen darstellt. In fast jedem Projekt sind Variablenmodifikationen und -neuberechnungen in erheblichem Umfang erforderlich.

Profis modellieren dabei nicht „per Hand“ so lange an der Rohdatei herum, bis die Fertigdatei entstanden ist, sondern sie erstellen, z.B. durch Konservieren von bearbeiteten Dialogboxen, ein so genanntes **SPSS-Programm** (siehe unten), das alle Transformationen erledigt und das bei Bedarf auch wiederholt ausgeführt werden kann.

Die zweistufige Projektdatenverwaltung mit Roh- und Fertigdatei verhindert in Kombination mit dem SPSS-Transformationsprogramm, dass bei jeder Änderung der Rohdaten die erwähnten Transformationen zur Fertigdatei „per Hand“ wiederholt werden müssen. Solche Änderungen der Rohdaten (z.B. durch Fehlerkorrekturen oder Stichprobenerweiterungen) sind eher die Regel als die Ausnahme.

Weil die Kommandos des Transformationsprogramms auch mit Hilfe von korrespondierenden Dialogboxen erstellt werden können, erfordert die professionelle Vorgehensweise kaum Programmierkenntnisse.

Es wird also folgende Struktur für die Verwaltung der Projektdaten vorgeschlagen:



Die Erläuterungen in diesem Abschnitt werden vermutlich erst dann voll verständlich, wenn Sie sich mit Variablentransformationen und SPSS-Programmen auskennen.

Nach diesem Vorausblick wenden wir uns wieder der aktuellen Aufgabe zu: Wir tragen die erhobenen Daten in die Rohdatendatei **kfar.sav** ein.



### 3.2.7 Dateneingabe

Wechseln Sie bei Bedarf zur Datenansicht, und geben Sie die Daten des ersten Falles ein:

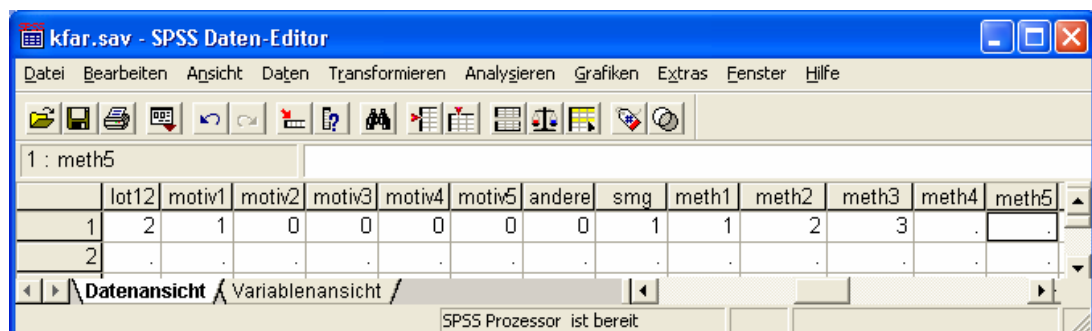
- Aktivieren Sie nötigenfalls die Zelle zur ersten Variablen des ersten Falles, und tippen Sie den zugehörigen Wert ein.
- Drücken Sie die Taste mit dem Rechtspfeil (→) oder die **Tabulator**-Taste (⇨), um den eingetippten Wert zu quittieren und die Zellenmarkierung um eine Spalte nach *rechts* zu verschieben (zur nächsten Variablen):



Auch die **Enter**-Taste quittiert den eingetippten Wert, bewegt jedoch anschließend die Zellenmarkierung um eine Zeile nach *unten* (zum nächsten Fall), was in unserer jetzigen Lage weniger praktisch ist.

Wenn Sie auf Abwege geraten sind, können Sie die Zellenmarkierung jederzeit per Mausklick neu positionieren.

- Sobald für einen neuen Fall die erste Variablenausprägung eingetragen und quittiert wurde, erhält er für die restlichen Variablen den Initialisierungswert SYSMIS (dargestellt durch einen Punkt).
- Wenn über den Menübefehl **Ansicht > Wertelabels** die Anzeige von Wertelabels aktiviert worden ist, erscheint z.B. in der markierten GESCHL-Zelle ein Drop-Down-Menü zur „Erleichterung“ der Werteingabe. Allerdings erscheint das Drop-Down-Menü nur bei bereits vorhandenen Fällen. Verzichten Sie durch einem erneuten Aufruf des Menübefehls auf die Wertelabels und die fragwürdigen Eingabehilfen.
- Tragen Sie die restlichen Werte des ersten Falles ein, jeweils quittiert mit der Tabulator-taste. So sieht der vollständig erfasste erste Fall unserer Stichprobe im Datenfenster aus (bei abgeschalteter Wertelabels-Anzeige):



- Wenn Sie den Wert der letzten Variablen mit der Tabulatortaste quittieren, setzt SPSS freundlicherweise die Zellenmarkierung gleich in die erste Datenzelle des nächsten Falles, so dass Sie die Dateneingabe unmittelbar fortsetzen können.

### 3.2.8 Daten korrigieren

#### 3.2.8.1 Wert in einer Zelle ändern

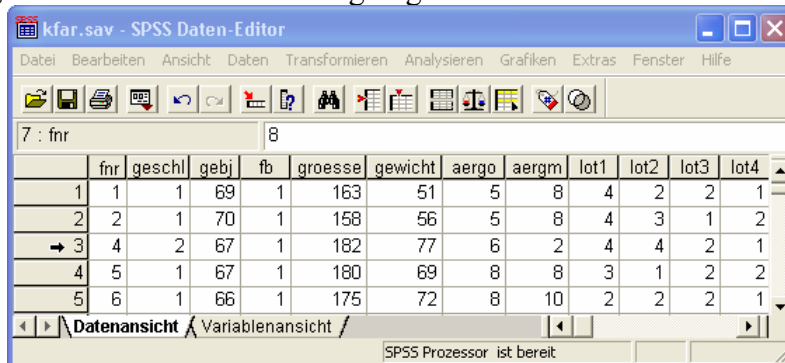
Natürlich können die Eintragungen in einer Zelle jederzeit korrigiert werden:

- Wert ersetzen:
  - Zelle markieren
  - neuen Wert eintippen, wobei der alte überschrieben wird
  - neuen Wert quittieren, z.B. mit **Enter**
- Wert editieren:
  - Doppelklick auf die Zelle
  - Wert editieren
  - neuen Wert quittieren, z.B. mit **Enter**

#### 3.2.8.2 Einen Fall einfügen

Gehen Sie folgendermaßen vor, um einen Fall einzufügen:

- Setzen Sie einen rechten Mausklick auf die (von SPSS gesetzte) Zeilennummer desjenigen Falles an, *vor* dem ein neuer Fall eingefügt werden soll.



Daraufhin wird die gesamte angeklickte Zeile markiert, und es erscheint ein Kontextmenü.

- Wählen Sie aus dem Kontextmenü die Option **Fälle einfügen**

Der neue Fall erhält bei allen Variablen den Wert System-Missing. Diese Initialisierungswerte können dann natürlich beliebig überschrieben werden.

#### 3.2.8.3 Einen Fall löschen

Gehen Sie folgendermaßen vor, um einen Fall, d.h. eine Zeile der Datenmatrix, komplett zu löschen:

- Setzen Sie einen rechten Mausklick die die (von SPSS gesetzte) Zeilennummer des überflüssigen Falles. Daraufhin wird die gesamte angeklickte Zeile markiert, und es erscheint ein Kontextmenü.
- Wählen Sie aus dem Kontextmenü die Option **Löschen**

### 3.2.8.4 *Einen Fall verschieben*

Gehen Sie folgendermaßen vor, um einen Fall per Drag & Drop (Ziehen und Ablegen) zu verschieben:

- Setzen Sie einen linken Mausklick auf die (von SPSS gesetzte) Zeilennummer. Daraufhin wird die gesamte Zeile markiert. Lassen Sie anschließend die Maustaste wieder los.
- Klicken Sie erneut auf die Zeilennummer, und halten Sie dabei die Maustaste gedrückt.
- Bewegen Sie bei gedrückter Maustaste den Mauszeiger zum Ziel der Verschiebungsaktion. Der augenblicklich eingestellte Zielort wird von SPSS durch eine rote Linie gekennzeichnet.
- Wenn Sie die Maustaste los lassen, erscheint der Fall am neuen Ort.

### 3.2.9 Weitere Möglichkeiten des Dateneditors

Über die beschriebenen Methoden hinaus bietet der Dateneditor u.a. die Möglichkeit, beliebige rechteckige Segmente der Datenmatrix auszuschneiden, zu kopieren und einzufügen.

Wer derartige, relativ fehleranfällige Umordnungsmaßnahmen vornimmt, wird gelegentlich von der Möglichkeit profitieren, mit:

#### **Bearbeiten > Rückgängig**

die letzte Änderung rückgängig machen zu können.

In Abschnitt 4.7 wird beschrieben, wie Sie im Datenfenster nach Variablenausprägungen suchen können.

### 3.2.10 Übung

Für die Teilnehmer(innen) des realen SPSS-Kurses steht nun die Erfassung der erhobenen Daten an. Geben Sie alle Fälle ein, und sichern Sie (auch zwischendurch) in die zugeordnete Datendatei, z.B. U:\Eigene Dateien\SPSS\kfar.sav.

Wer dem Vorschlag in diesem Manuskript folgend zur Erfassung der Antworten auf die offene Frage im Fragebogenteil 4b ein dynamisches und sparsames Set aus kategorialen Variablen vorgesehen hat (z.B. METH1 bis METH3), der muss nicht nur mechanisch Daten eintippen, sondern auch gelegentlich mit Kreativität und Ordnungssinn neue Methodenkategorien definieren und dokumentieren.

Beim Erfassen der Daten, die in diesem Manuskript analysiert werden, entstand folgende Liste:

Kategorie	Code
Faktorenanalyse	1
Regressionsanalyse	2
Korrelationsanalyse	3
Varianzanalyse	4
Strukturgleichungsanalyse	5
Clusteranalyse	6
Diskriminanzanalyse	7
Logistische Regression	8
Conjoint-Analyse	9

Diese Tabelle vervollständigt unseren Kodierplan (vgl. Abschnitt 1.4.3.5).

Es bietet sich an, die Definition der Variablen METH1 bis METH3 durch entsprechende Wertelabels zu vervollständigen (vgl. Abschnitt 3.2.2.3)

---

## 4 Univariate Verteilungs- und Fehleranalysen

In diesem Abschnitt werden Sie erfahren, wie schnell und bequem mit SPSS numerische und graphische Analysen durchgeführt werden können. Wir werden unsere Daten mit Hilfe deskriptiver Auswertungsmethoden sorgfältig auf Erfassungsfehler untersuchen. Dabei schlagen wir zwei Fliegen mit einer Klappe, denn eine sorgfältige Verteilungsanalyse aller Variablen gehört ohnehin zur Pflicht bei jeder empirischen Studie.

In manchen Projekten wird sich die Forschungsarbeit sogar auf die Beschreibung von univariaten Verteilungen beschränken (z.B. in der Meinungsforschung). Meist sind aber auch multivariate Zusammenhangsanalysen von Interesse.

### 4.1 Erfassungsfehler

Speziell bei der manuellen Datenerfassung sind Fehler praktisch unvermeidbar. Manche von ihnen sind als Verstöße gegen allgemeine Gültigkeitsregeln relativ leicht aufzuspüren:

Beispiel: Wenn bei der Variablen GESCHL nur die Werte 1 (für Frauen) und 2 (für Männer) erlaubt sind, dann ist z.B. der Wert 3 sofort als falsch erkennbar.

Weit schwieriger zu entdecken sind Fehler, die keine allgemeine Gültigkeitsregel verletzen:

Beispiel: Wenn unter der oben angegebenen GESCHL-Kodierungsvorschrift für den Untersuchungsteilnehmer Kurt Müller versehentlich der Wert 1 eingegeben wurde, dann kann dieser Fehler nur durch aufwändige Handarbeit gefunden werden.

Welcher Aufwand bei der Datenprüfung erforderlich bzw. sinnvoll ist, hängt wesentlich davon ab, wie die Daten erfasst worden sind (vgl. Abschnitt 3.1).

*Nobody is perfect* gilt übrigens nicht nur für Menschen, sondern auch für Maschinen. Daher sollte man vorsichtshalber auch bei Verwendung einer automatischen Erfassungsmethode stichprobenartig die Datenintegrität überprüfen.

Nach der Erfassung per Texteditor ist die Menge potentieller Fehler besonders groß. Deshalb wurde oben nachdrücklich von dieser veralteten Erfassungsmethode abgeraten. Konsequenterweise gehen wir auch im Abschnitt über Datenprüfung nicht auf die speziellen Probleme ein, die nach dem Erfassen per Texteditor auftreten können.

#### 4.1.1 Überprüfung von Gültigkeitsregeln

Wir beschränken uns auf die Suche nach Werten außerhalb der zulässigen Wertebereiche, wenngleich damit nicht alle Möglichkeiten zum Aufspüren von verletzten Gültigkeitsregeln ausgereizt sind.

Bei der Erfassung per Datenbankprogramm mit Plausibilitätskontrolle werden unzulässige Werte zurückgewiesen und folglich von der Datendatei fern gehalten. Bei der Erfassung mit dem SPSS-Dateneditor findet eine derartige Eingangskontrolle nicht statt. Eine so entstandene Datei muss daher systematisch nach Daten außerhalb der zulässigen Bereiche durchsucht werden. Dies kann allerdings ohne großen Zusatzaufwand im Rahmen der aus wissenschaftlichen Gründen ohnehin empfehlenswerten univariaten Verteilungsanalyse geschehen.

#### 4.1.2 Überprüfung von Einzelwerten

Fehler, die gegen keine Gültigkeitsregel verstoßen, lassen sich nur mit Fleißarbeit entdecken, wobei z.B. folgende Vorgehensweisen möglich sind:

- Man vergleicht die erfassten Daten Wert für Wert mit den schriftlichen Unterlagen.
- Manche Datenbankprogramme versuchen, die Erfasser durch Kontrollen und Sanktionen zu sorgfältiger Arbeit zu motivieren: INPUT II (siehe Abschnitt 3.1.2.2) erlaubt z.B. die Festlegung einer Kontrollwahrscheinlichkeit, mit der ein Erfasser einen Teil des letzten Datensatzes nochmals eingeben muss. Bei Erfolg sinkt die Kontrollwahrscheinlichkeit, bei Misserfolg werden die diskrepanten Daten präsentiert, und die Kontrollwahrscheinlichkeit steigt.

Eine aufwändige Prüfmethode ist *bei kleinen Stichproben* durchaus empfehlenswert, denn:

- Der Zeitaufwand ist erträglich.
- Erfassungsfehler wirken sich besonders stark aus.

Wir wollen exemplarisch den Effekt von Erfassungsfehlern auf die Varianz eines Stichprobenmittelwerts untersuchen und nehmen für  $n$  Beobachtungen  $X_i$  ( $i = 1, \dots, n$ ) an, dass sie jeweils mit einem Erfassungsfehler  $F_i$  belastet sind, wobei die Erfassungsfehler den Erwartungswert Null haben sowie untereinander und von den korrekten Beobachtungswerten  $T_i$  unabhängig sind:

$$X_i = T_i + F_i, \quad E(F_i) = 0, \quad E(X_i) = E(T_i) = \mu$$

$$\text{Var}(T_i) = \sigma^2, \quad \text{Var}(F_i) = \sigma_F^2$$

Für die Varianz des Mittelwertes aus den fehlerfrei erfassten Werten gilt:

$$\text{Var}(\bar{T}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n T_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(T_i) = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n}$$

Für die Varianz des Mittelwertes der fehlerhaft erfassten Werte erhalten wir:

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n (T_i + F_i)\right) = \frac{1}{n^2} \sum_{i=1}^n (\text{Var}(T_i) + \text{Var}(F_i)) = \frac{1}{n^2} n (\sigma^2 + \sigma_F^2) = \frac{\sigma^2}{n} + \frac{\sigma_F^2}{n}$$

Offenbar hängt der Präzisionsverlust im Stichprobenmittel, das als Schätzwert für den Erwartungswert in der Population dient, von der Erfassungsfehlervarianz  $\sigma_F^2$  und von der Stichprobengröße  $n$  ab. Während sich in einer großen Stichprobe der niedrige Ausgangswert  $\frac{\sigma^2}{n}$  des Schätzfehlers nur unwesentlich erhöht, kommt es in einer kleinen Stichprobe mit ihrem bereits ungünstigen Ausgangsniveau zu einem erheblichen Präzisionsverlust. Als unerwünschte Folgen stellen sich ein:

- Unpräzise Parameterschätzungen
- Reduzierte Power bei Hypothesentests

Obwohl bei unserer kleinen Stichprobe eine Einzelprüfung aller Werte angemessen ist, verzichten wir aus Zeitgründen darauf. Es gehört übrigens zu den lehrreichen Erfahrungen der realen SPSS-Kurse, dass die selbständig als Untersuchungsleiter agierenden Teilnehmer aus Kopien desselben Fragebogenstapels aufgrund individueller Erfassungsfehler recht unterschiedliche Ergebnisse ermitteln (auch bei den zentralen Hypothesentests).

## 4.2 Öffnen einer SPSS-Datendatei

Vermutlich haben Sie nach der anstrengenden Datenerfassung eine Pause eingelegt und SPSS verlassen, so dass wir jetzt offiziell die Fortsetzung einer unterbrochenen Projektarbeit üben können. Starten Sie SPSS, und öffnen Sie Ihre vorhandene Rohdatendatei **kfar.sav**, entweder mit Hilfe des Startassistenten oder über den Menübefehl

### Datei > Zuletzt verwendete Daten

Beim Öffnen einer Datendatei legt SPSS eine neue (temporäre) Arbeitsdatei an und kopiert die eingelesenen Daten samt Variablendeklarationen dorthin. Alle Veränderungen, die Sie in der Datenmatrix oder im Deklarationsteil vornehmen, wirken sich zunächst nur auf die temporäre Arbeitsdatei aus. Gegebenenfalls müssen Sie also diese Änderungen über den Menübefehl

### Datei > Speichern

in die permanente SPSS-Datendatei **kfar.sav** übernehmen.

## 4.3 Statistische Auswertungen durchführen: Häufigkeitsanalysen zur Prüfung der Variablen FNR

Da wir unsere Daten mit dem SPSS-Dateneditor erfasst haben, der keine Plausibilitätskontrolle bei der Eingabe vornimmt, müssen wir nach den Überlegungen aus Abschnitt 4.1 systematisch nach unzulässigen Werten suchen. Die meisten der dazu erforderlichen deskriptiven Datenanalysen wären im Rahmen der routinemäßigen Verteilungsuntersuchung ohnehin fällig.

Der erste Test dient allerdings ausschließlich zur Datenprüfung, weil dabei die Fallidentifikations-Variable FNR untersucht wird. Es ist sogar etwas zweifelhaft, ob man tatsächlich „der Vollständigkeit halber“ in die Überprüfung dieser administrativen Variablen Zeit investieren sollte.

Weil die Manuskript-Stichprobe den Umfang  $n = 31$  hat, und es keinen Grund für eine lückenhafte Nummerierung gab, müssen nach fehlerfreier Erfassung bei dieser Variablen die Werte 1, ..., 31 jeweils genau einmal auftreten. Daraus ergeben sich einige notwendige Bedingungen, die sich leicht nachprüfen lassen:

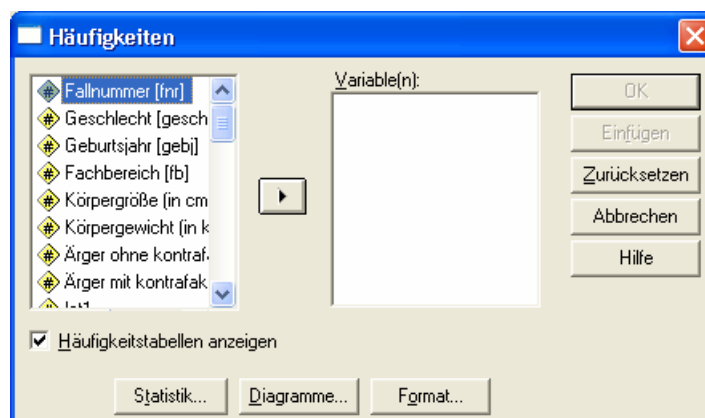
- In der Stichprobe müssen bei der Variablen FNR 31 gültige Werte vorliegen. (MD-Indikatoren sind hier nicht erlaubt.)
- Der kleinste Wert muss gleich 1, und der größte Wert muss gleich 31 sein.
- Jeder Wert darf höchstens einmal auftreten, d.h. der Stichproben-Modus muss die Häufigkeit 1 haben.

Zur Überprüfung der Bedingungen lassen wir in einer *Häufigkeitsanalyse* für die Variable FNR folgende Statistiken berechnen: Anzahl valider Fälle, Minimum, Maximum und Modus.

Mit dem Menübefehl

### Analysieren > Deskriptive Statistik > Häufigkeiten...

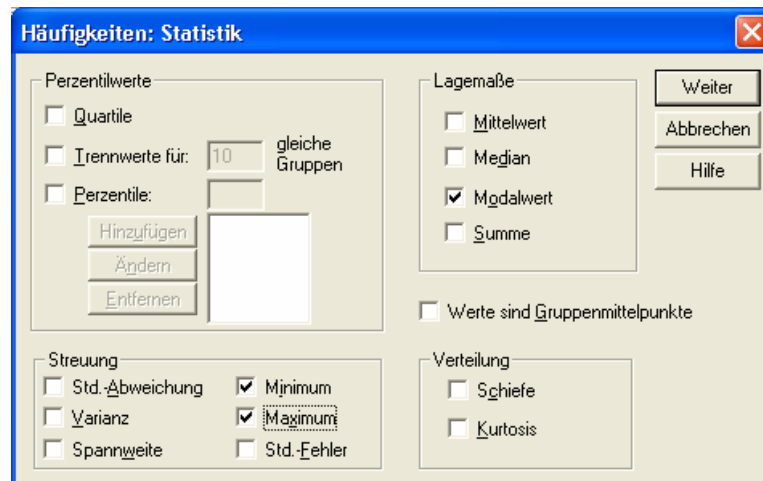
erhalten wir die folgende Dialogbox zur Anforderung von Häufigkeitsanalysen:



Zur bequemen Spezifikation der im aktuellen Prozeduraufruf zu analysierenden Variablen dienen die beiden Variablen-Auswahlbereiche. Links stehen alle Variablen der Arbeitsdatei, die derzeit *nicht* für die Analyse ausgewählt sind (*Anwärterliste*). Rechts daneben, im Bereich **Variable(n)**, stehen die Ausgewählten (*Teilnehmerliste*). Dazwischen befindet sich ein Transportschalter, mit dem sich links markierte Variablen nach rechts und rechts markierte Variablen nach

links verschieben lassen. Markieren Sie also links die Fallnummern-Variable FNR und drücken Sie auf den Transportschalter.

Zur Auswahl der gewünschten **Statistiken** müssen Sie die zuständige Subdialogbox per Knopfdruck aktivieren. Um eine der hier aufgelisteten Möglichkeiten zu wählen, ist das zugehörige Kontrollkästchen zu markieren:



Quittieren Sie die Subdialogbox mit **Weiter** und die Hauptdialogbox mit **OK**. Daraufhin führt SPSS die Berechnungen aus und präsentiert die Ergebnisse im Ausgabefenster (SPSS Viewer), das sich in den Vordergrund drängt.

Bei Anforderung einer Häufigkeitsanalyse produziert SPSS per Voreinstellung eine Tabelle, die für jeden aufgetretenen Wert eine Zeile mit folgenden Angaben enthält:

- Absolute Häufigkeit
- Prozentualer Anteil am Stichprobenumfang
- Prozentualer Anteil an den validen Werten (ohne MD-deklarierte Werte)
- kumulativer valider Prozentanteil (Anteil valider Werte, die nicht größer sind)

Außerdem berichtet SPSS unaufgefordert, wie viele Fälle einen validen Wert bzw. einen MD-deklarierten Wert haben. Weitere Leistungen müssen explizit angefordert werden.

Obige Dialogbox liefert folgende Statistiken:

#### Statistiken

Fallnummer		
N	Gültig	31
	Fehlend	0
Modus		1 <sup>a</sup>
Minimum		1
Maximum		31

a. Mehrere Modi vorhanden. Der kleinste Wert wird angezeigt.

Indizien für Erfassungsfehler finden sich nicht: Alle 31 Personen haben einen validen Wert, das Minimum ist 1, das Maximum ist 31.

Laut Häufigkeitstabelle (hier verkürzt wiedergegeben) hat der (natürlich nicht eindeutige) Modalwert die Häufigkeit 1:

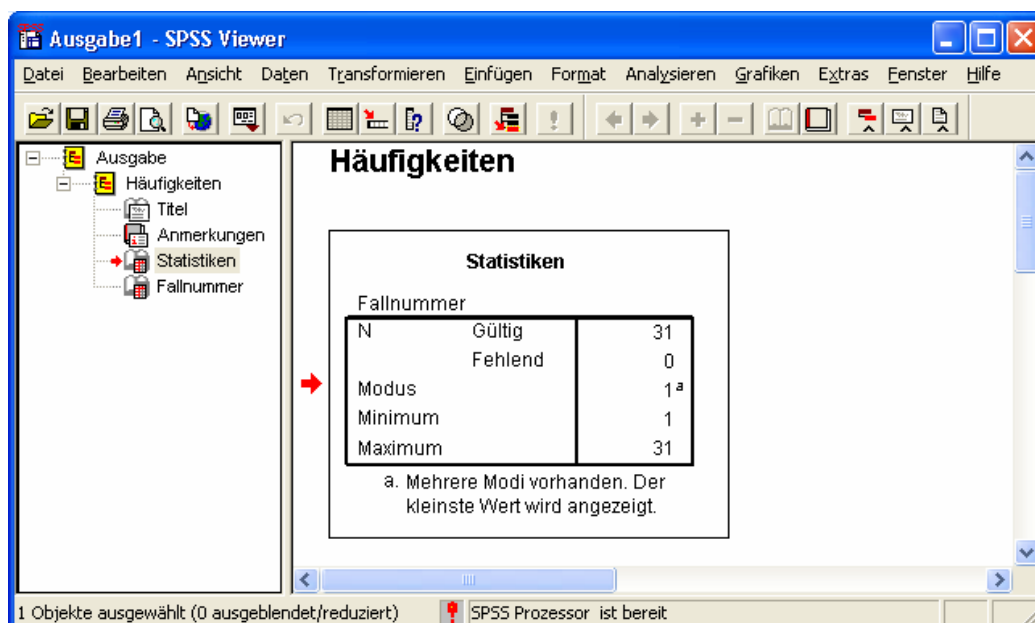
		Fallnummer			
		Häufigkeit	Prozent	Gültige Prozente	Kumulierte Prozente
Gültig	1	1	3,2	3,2	3,2
	2	1	3,2	3,2	6,5
	3	1	3,2	3,2	9,7
	4	1	3,2	3,2	12,9
	,	,	,	,	,
	,	,	,	,	,
	,	,	,	,	,
	30	1	3,2	3,2	96,8
	31	1	3,2	3,2	100,0
	Gesamt	31	100,0	100,0	

Wir haben uns bei der FNR-Prüfung auf einige *notwendige* Bedingungen beschränkt, weil momentan nur elementare SPSS-Operationen benutzt werden sollen. Eine perfekte Kontrolle ist bei dieser administrativen Variablen ohnehin nicht erforderlich.

Die obigen SPSS-Ausgaben wurden übrigens aus dem Ausgabefenster via Windows-Zwischenablage in Microsoft Word<sup>®</sup> übertragen. Mit dieser Form des Datenaustauschs und mit anderen Möglichkeiten beim Arbeiten mit dem Ausgabefenster (Viewer) beschäftigen wir uns im nächsten Abschnitt.

#### 4.4 Arbeiten mit dem Ausgabefenster (Teil I)

In seiner voreingestellten Variante ist das SPSS-Ausgabefenster, das auch als **Viewer** bezeichnet wird, zweigeteilt in den Navigationsbereich (die Gliederungsansicht) am linken Rand und den eigentlichen Inhaltsbereich:



So soll ein schnelles Navigieren zwischen den verschiedenen Ausgabebestandteilen ermöglicht werden.

Die Aufteilung des verfügbaren Platzes auf die beiden Teile des Viewers kann per Maus beliebig verändert werden: Trennlinie anklicken und bei gedrückter Maustaste horizontal verschieben. Wesentliche Bestandteile des Inhaltsbereichs sind Pivot-Tabellen, Grafiken und Textausgaben. Zu ihrer Nachbearbeitung steht jeweils ein spezieller Editor zur Verfügung, der per Doppelklick



auf das Objekt gestartet wird. Außerdem können in einem Viewer-Dokument noch protokollierte SPSS-Anweisungen, Warnungen, Anmerkungen und Titelzeilen auftreten.

#### 4.4.1 Arbeiten im Navigationsbereich

Die meisten der anschließend beschriebenen Aktionen im Navigationsbereich wirken sich synchron auch auf den Inhaltsbereich aus.

##### 4.4.1.1 Fokus positionieren

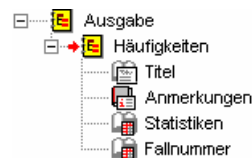
Ein kleiner roter Pfeil zeigt im Gliederungsbereich auf die Bezeichnung derjenigen Ausgabe, die im Inhaltsbereich gerade privilegiert dargestellt wird. Per Mausklick auf eine andere Ausgabenbeschriftung kann dieser Fokus beliebig verschoben werden.


##### 4.4.1.2 Ausgabeblöcke bzw. Teilausgaben aus- oder einblenden

Ein *Block* mit zusammengehörigen Ausgaben (in der Regel entstanden aus einer Analyseanforderung) wird ...

- ausgeblendet: per Mausklick auf das Minus-Zeichen neben dem Block-Symbol  oder per Doppelklick auf das Block-Symbol.

Beispiel:



- eingeblendet: per Mausklick auf das Plus-Zeichen neben Block-Symbol  oder per Doppelklick auf das Block-Symbol.


Beispiel:



Eine *Teilausgabe* innerhalb eines Blockes wird per Doppelklick auf das zugehörige Buchsymbol aus- bzw. eingeblendet. Das Buchsymbol erscheint dementsprechend zugeklappt (im Beispiel: Anmerkungen) oder aufgeklappt (im Beispiel: Statistiken).

##### 4.4.1.3 Ausgabeblöcke oder -teile markieren

Im Navigationsbereich können Sie auf windows-übliche Weise Ausgabeblöcke und/oder Teilausgaben markieren:

- Einen Ausgabeblock: Per Mausklick auf das Block-Symbol oder auf die Beschriftung
- Eine Teilausgabe: Per Mausklick auf das Buchsymbol oder auf die Beschriftung
- Mehrere Blöcke bzw. Teile: Per -Mausklick bzw. **Strg**-Mausklick

Sie können markierte Blöcke bzw. Teilausgaben z.B. mit der **Entf**-Taste löschen oder in die Windows-Zwischenablage befördern (siehe Abschnitt 4.4.4).

#### 4.4.2 Viewer-Dokumente drucken

Über den Standardbefehl

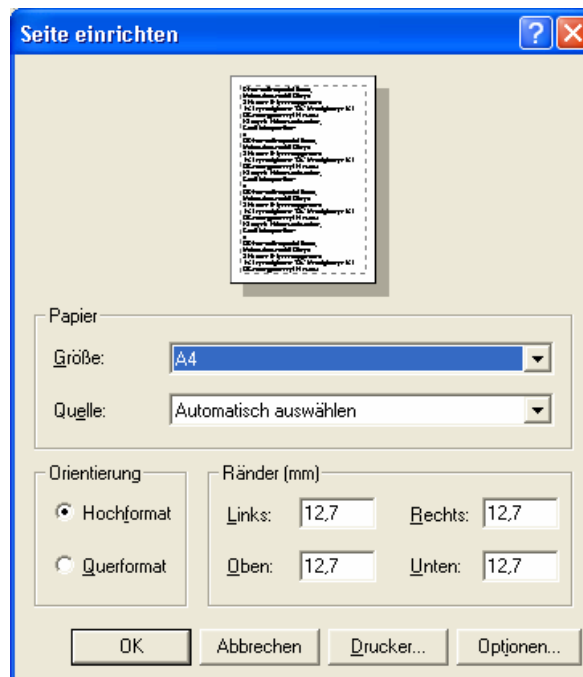
**Datei > Drucken**




können Sie alle angezeigten oder alle markierten Ausgabebestandteile drucken.

Zur Gestaltung der Ausgabe finden sich nach

## Datei > Seite einrichten

in der folgenden Dialogbox einige Möglichkeiten:



In der **Optionen**-Subdialogbox können u.a. Kopf- und Fußzeilen festgelegt werden, z.B. unter Verwendung von Standardelementen wie **Überschrift erster Stufe** , **Datum**  und **Uhrzeit** .

Den Erfolg Ihrer Bemühungen können Sie über **Datei > Seitenansicht** auch schon vor dem Ausdruck begutachten.

Auf den Pool-PCs an der Universität Trier können Sie den Inhalt des Ausgabefensters als PDF-Datei exportieren, indem Sie im Druckdialog den Drucker mit dem Namen **Rumborak PDF Writer Plus** wählen.

### 4.4.3 Ausgaben sichern und öffnen

Zum Speichern eines Viewer-Dokuments dienen die Menübefehle **Datei > Speichern unter** bzw. **Datei > Speichern**. Dabei entstehen Viewer-Dateien, die üblicherweise durch die Namensendung **.spo** gekennzeichnet werden. SPSS-Ausgaben sollten z.B. *dann* in elektronischer Form gespeichert werden, wenn sie (auszugsweise) in Dokumente anderer Programme eingegangen sind, z.B. in MS-Word - Dateien. Mit SPSS ist eine nachträgliche Modifikation dieser Ausgaben leicht möglich, mit den Fremdprogrammen aber kaum.

Zum Öffnen eines Viewer-Dokuments mit den Befehlen **Datei > Öffnen > Ausgabe** oder **Datei > Zuletzt geöffnete Dateien** gibt es nichts Ungewöhnliches zu berichten.

### 4.4.4 Objekte via Zwischenablage in andere Anwendungen übertragen

Mit der Tastenkombination **Strg+C** oder mit dem Menübefehl

#### **Bearbeiten > Kopieren**

fordert man SPSS auf, ein markiertes Ausgabe-Objekt (z.B. Tabelle oder Diagramm) in die Zwischenablage zu befördern. Zum Einfügen in der Zielanwendung kann man den Menübefehl

#### **Bearbeiten > Einfügen**

bzw. die Tastenkombination **Strg+V** verwenden.

SPSS legt die Daten in mehreren Formaten in der Zwischenablage ab, und je nach Zielanwendung kann es sinnvoll sein, über den Menübefehl

#### **Bearbeiten > Inhalte Einfügen**


auf das entnommene Format Einfluss zu nehmen. Wenn Sie beim Einfügen einer Tabelle das Format **Grafik (Windows-Metadatei)** wählen, erhalten Sie in der Zielanwendung ein Grafik-Implantat mit dem Original-Design aus dem SPSS-Viewer, das nur noch Größen- und Positionsänderungen erlaubt. So wurden z.B. die in Abschnitt 4.3 wiedergegebenen Tabellen übertragen. Zum selben Ergebnis gelangt, wer im SPSS-Viewer Tabellen mit der Tastenkombination **Strg+K** oder mit dem Menübefehl

#### **Bearbeiten > Objekte Kopieren**

in die Zwischenablage befördert und in der Zielanwendung mit **Bearbeiten > Einfügen** bzw. **Strg+V** entnimmt.

Über **Bearbeiten > Objekte Kopieren** lassen auch *mehrere* markierte Tabellen gemeinsam aus dem Viewer in die Zwischenablage übertragen.

#### **4.4.5 Übungen**

- 1) Markieren Sie den Ausgabeblock mit der Häufigkeitsanalyse, und löschen Sie ihn mit der **Entf**-Taste.
- 2) Steigen Sie erneut in die Dialogbox zur Häufigkeitsanalyse ein. Statt den zugehörigen Menübefehl zu wiederholen, können Sie einfacher mit dem Symbol  eine Liste der zuletzt benutzten Dialogboxen aufrufen und daraus per Mausklick den Eintrag **Häufigkeiten** wählen. Die Dialogbox ist noch im selben Zustand, den Sie eben verlassen haben. Dies gilt selbstverständlich generell in SPSS, so dass Sie bei der sukzessiven Modifikation einer Anforderung innerhalb einer Sitzung jeweils auf dem letzten Stand weitermachen können.
- 3) Schalten Sie die Häufigkeitstabelle über das zugehörige Kontrollkästchen aus, und lassen Sie die Häufigkeitsanalyse erneut ausführen.

#### **4.5 Graphische Darstellungen in Statistik-Dialogboxen anfordern: Häufigkeits- bzw. Fehleranalyse für die Variablen GESCHL und FB**

Nun wollen wir weitere Variablen untersuchen und dabei auch graphische Verteilungsdarstellungen verwenden. Dazu rufen wir erneut die Dialogbox zur Häufigkeitsanalyse auf und beseitigen alle alten Festlegungen (auch in den Subdialogboxen) mit dem Schalter **Zurücksetzen**.

Dann transportieren wir nacheinander die Variablen GESCHL und FB aus der Anwärterliste (links) in die Teilnehmerliste (rechts).

Anschließend begeben wir uns in die Subdialogbox **Diagramme** und entscheiden uns im Rahmen **Diagrammtyp** für **Balkendiagramme**, weil die Merkmale Geschlecht und Fachbereich nominalskaliert sind. Wer nicht mehr genau weiß, wozu man Balkendiagramme und Histogramme verwendet, kann sich mit der kontextsensitiven **Hilfe** Aufklärung verschaffen.

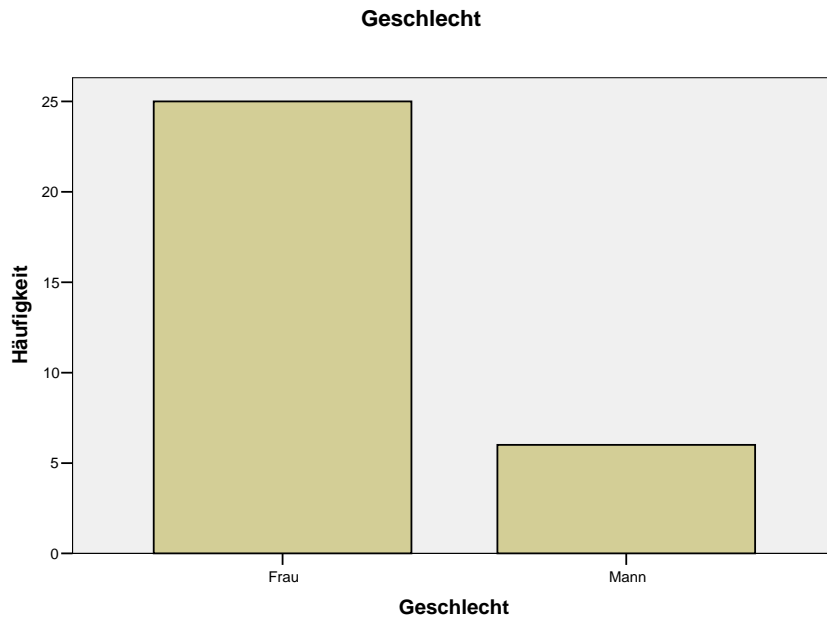
Jetzt starten wir die neue Analyse, indem wir die Subdialogbox mit **Weiter** und die Hauptdialogbox mit **OK** quittieren.

Im Viewer erhalten wir für die Variable GESCHL die Häufigkeitstabelle

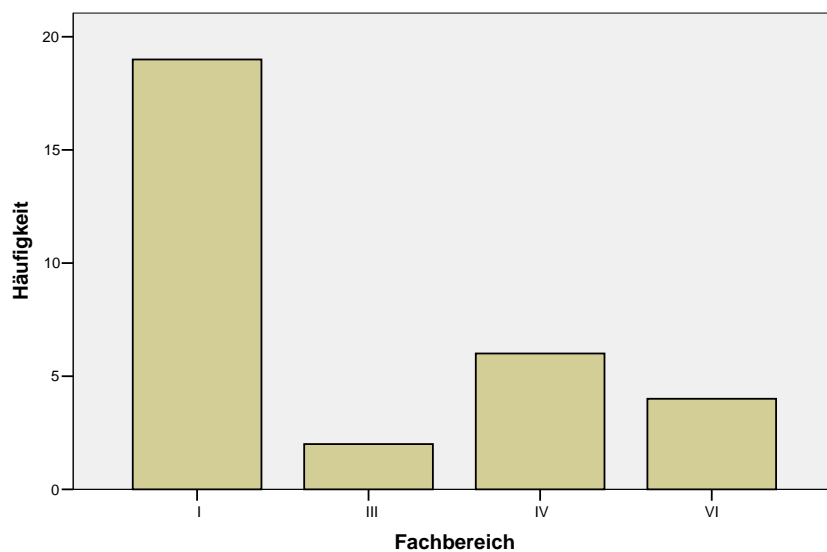
**Geschlecht**

	Häufigkeit	Prozent	Gültige Prozente	Kumulierte Prozente
Gültig Frau	25	80,6	80,6	80,6
Mann	6	19,4	19,4	100,0
Gesamt	31	100,0	100,0	

und das folgende Balkendiagramm:



Zunächst beobachten wir, dass bei der Variablen GESCHL kein unzulässiger Wert vorliegt. Bei der Geschlechtsverteilung stellen wir einen sehr hohen Frauenanteil fest, der als wesentliches Merkmal unserer Stichprobe berichtet werden muss. Bei potentiell geschlechtsabhängigen Ergebnissen müssen wir besonders vorsichtig interpretieren und generalisieren. Erste Hinweise zur Ursache der hohen Frauenquote liefert die empirische Verteilung der Fachbereichs-Variablen:

**Fachbereich**

Wir sehen, dass im SPSS-Kurs, der die Manuskript-Daten geliefert hat, der Fachbereich I sehr stark vertreten war, was mit dem Kurstermin zusammenhängen mag. Im Fachbereich I der Universität Trier (Fächer: Philosophie, Pädagogik, Psychologie) ist aber der Frauenanteil sehr hoch. Obige Abbildungen wurden übrigens mit der in Abschnitt 4.4.4 beschriebenen Methode vom SPSS-Viewer in MS-Word übertragen.

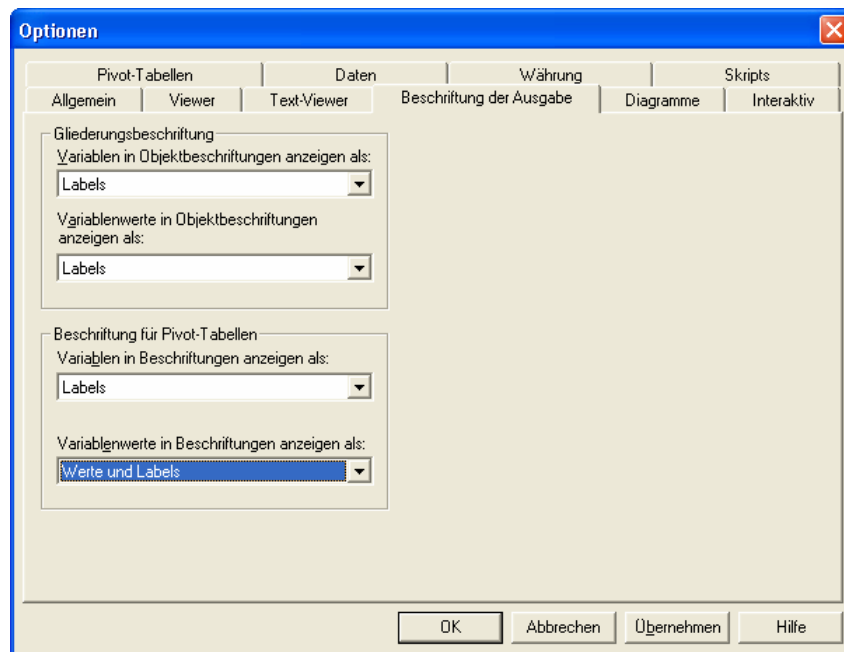
Der aktuelle Abschnitt sollte nur einen ersten Eindruck von den Grafikmöglichkeiten des SPSS-Systems vermitteln. Wir haben eine integrierte Grafik-Option der Dialogbox zur Häufigkeitsanalyse benutzt. Die meisten graphischen Darstellungsmöglichkeiten bietet SPSS über das Hauptmenü **Grafiken** an, mit dessen Optionen wir uns später befassen werden.

## 4.6 Häufigkeits- bzw. Fehleranalysen für die restlichen Projektvariablen

### 4.6.1 Übung

Mittlerweile verfügen Sie über genügend SPSS-Kenntnisse, um die restlichen Häufigkeits- bzw. Fehleranalysen zu unserem Projekt selbstständig durchführen zu können:

- 1) Die Merkmale Geburtsjahr, Größe, Gewicht und die beiden Ärgermaße können näherungsweise als metrisch angesehen werden. Lassen Sie sich daher für die zugehörigen Variablen ausgeben:
  - keine Häufigkeitstabellen  
Das für Tabellen zuständige Kontrollkästchen in der Dialogbox **Häufigkeiten** ist per Voreinstellung markiert. Sie müssen also die Markierung durch Anklicken beseitigen.
  - Histogramme mit eingezeichneter Normalverteilungsdichte
  - folgende Statistiken: Mittelwert, Median, Modalwert, Standardabweichung, Varianz, Minimum, Maximum, Schiefe, Kurtosis (Exzeß)
- 2) Lassen Sie sich für die LOT-Variablen ausgeben:
  - Häufigkeitstabellen
  - keine Grafiken
  - folgende Statistiken: Mittelwert, Median, Modalwert, Standardabweichung, Varianz, Minimum, Maximum
- 3) Lassen Sie sich für die Variablen MOTIV1 bis MOTIV5, ANDERE, SMG und METH1 bis METH3 ausgeben:
  - Häufigkeitstabellen
  - keine Grafiken
  - keine Statistiken
- 4) Prüfen Sie für alle Variablen nach, ob unzulässige Werte vorliegen.  
Sorgen Sie vorsichtshalber nach  
**Bearbeiten > Optionen > Beschriftung der Ausgabe**  
dafür, dass in Häufigkeitstabellen neben den eventuell definierten Labels auf jeden Fall auch die eigentlichen Werte angezeigt werden:



Anderenfalls ist der unglückliche Fall denkbar, dass ein falscher Wert aufgrund eines korrekten Labels unentdeckt bleibt, z.B.:

Fachbereich					
		Häufigkeit	Prozent	Gültige Prozente	Kumulierte Prozente
Gültig	0 I	19	61,3	61,3	61,3
	3 III	2	6,5	6,5	67,7
	4 IV	6	19,4	19,4	87,1
	6 VI	4	12,9	12,9	100,0
	Gesamt	31	100,0	100,0	

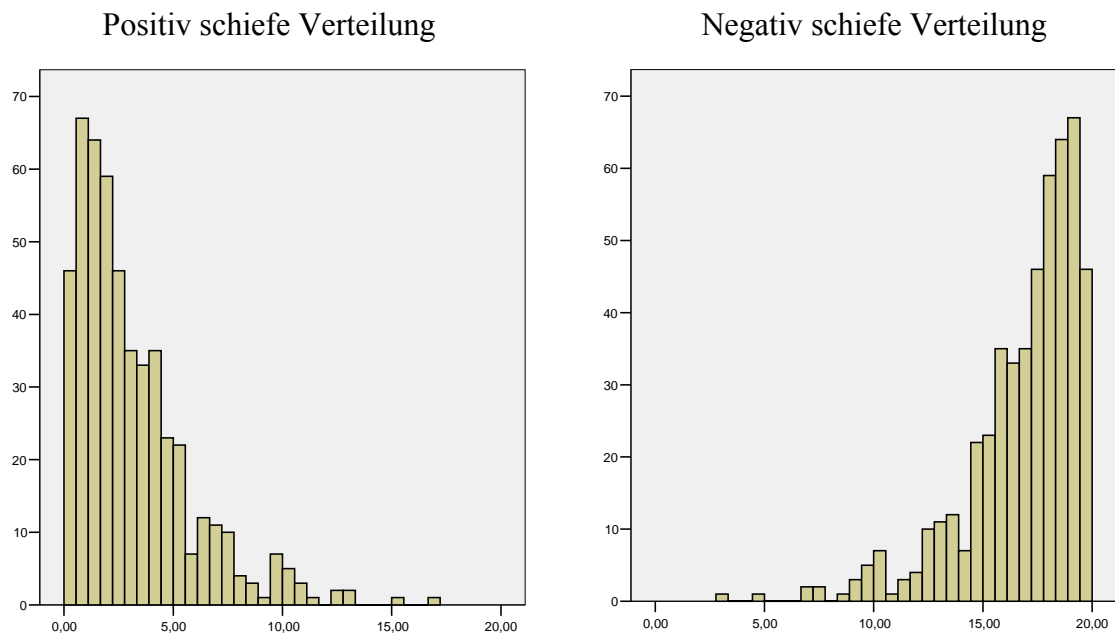
- 5) Untersuchen Sie bei den metrischen Variablen GROESSE, GEWICHT, AERGO und AERGM zusätzlich, ob diese annähernd normal verteilt sind. Beziehen Sie in Ihr Urteil die Statistiken Schiefe und Kurtosis sowie deren Standardfehler ein.

Die Vergleiche mit der Normalverteilung erfolgen hier aus purem Interesse an den Verteilungen der betrachteten Variablen, ohne dabei bereits an die *Verteilungsvoraussetzungen* irgendwelcher Testverfahren zu denken. Diese Voraussetzungen beziehen sich ohnehin häufig nicht auf die momentan von uns analysierten univariaten Verteilungen, sondern z.B. auf bedingte Verteilungen bzw. auf die Verteilungen der Residuen eines bestimmten statistischen Modells. Nähere Aussagen sind nur im Zusammenhang mit konkreten Testverfahren möglich.

Hinweise zu den Statistiken Schiefe und Kurtosis:

### Schiefe

Bei symmetrischen Variablen ist die Schiefe Statistik gerade gleich 0. Sie wird positiv bei linkssteil (bzw. rechtsschief) verteilten Variablen, wenn also die Verteilungsmasse am linken Rand konzentriert ist, und negativ bei rechtssteil (bzw. linksschief) verteilten Variablen, z.B.:



Zur Stichprobenschiefe wird auch der zugehörige Standardfehler ausgegeben, mit dessen Hilfe wir Tests zur Populationsschiefe veranstalten können. Diese sind allerdings nur approximativ gültig und vor allem in kleineren Stichproben mit Vorsicht zu genießen. Ihr Vorzug gegenüber später den vorzustellenden Normalverteilungs-Anpassungstests besteht darin, dass sie gezielt auf Verletzungen der Verteilungs-Symmetrie ansprechen.

Bei einem  $\alpha$ -Fehlerrisiko von 5 % ist die *zweiseitige* Nullhypothese, dass die Schiefe in der Population gerade gleich Null sei, zu verwerfen, falls:

$$\frac{|\text{Schiefe}|}{\text{SF}(\text{Schiefe})} > 1,96$$

Beim selben  $\alpha$ -Niveau entscheidet sich der Test zum gerichteten Hypothesenpaar:

$$H_0: \text{Schiefe} \geq 0 \quad \text{versus} \quad H_1: \text{Schiefe} < 0$$

gegen die Nullhypothese, falls:

$$\frac{\text{Schiefe}}{\text{SF}(\text{Schiefe})} < -1,65$$

Analog lässt sich natürlich auch die einseitige Nullhypothese mit umgekehrtem Vorzeichen prüfen.<sup>1</sup>

### Kurtosis (Exzeß)

Der Exzeß (synonym: Kurtosis, Breitipfligkeit, Wölbung) ist bei normalverteilten Variablen gleich Null. Er wird negativ bei breiteren und positiv bei schlankeren Verteilungen.

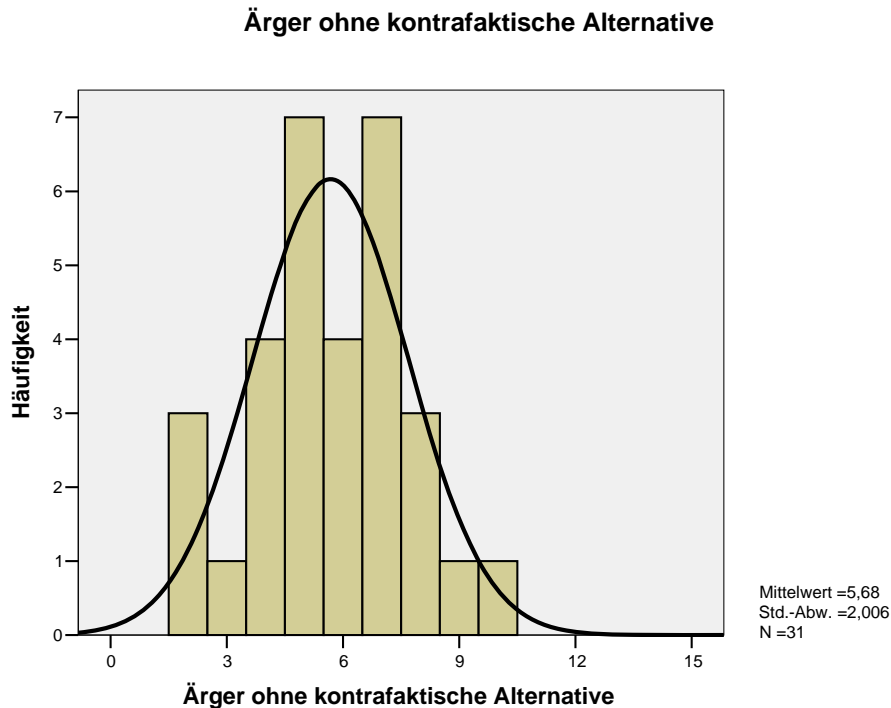
Mit Hilfe des Standardfehlers können analog zum Vorgehen bei der Schiefe-Statistik (siehe oben) „quick-and-dirty-Tests“ zum Exzeß in der Population durchgeführt werden.

<sup>1</sup> Wer in seinem Gedächtnis nicht mehr genügend Kenntnisse zur Inferenzstatistik reaktivieren konnte, der sei auf den Abschnitt 7.1 vertröstet.

#### 4.6.2 Diskussion ausgewählter Ergebnisse

##### a) Die Verteilungen der zentralen KFA-Variablen (AERGO, AERGM)

Bei den zentralen KFA-Variablen (AERGO, AERGM) finden sich keine irregulären Werte. Die Verteilungen fallen unterschiedlich aus. Einen recht normalen Eindruck macht die Verteilung der Ärgermessung in der Situation *ohne* kontrafaktische Alternative (AERGO):



Die Verteilungskennwerte Schiefe (= -0,08) und Kurtosis (= -0,277) sind nach den oben angegebenen Tests nicht signifikant von Null verschieden:

**Statistiken**

		Ärger ohne kontrafaktische Alternative	Ärger mit kontra- faktischer Alternative
N	Gültig	31	31
	Fehlend	0	0
Mittelwert		5,68	7,68
Median		6,00	8,00
Modus		5(a)	8
Standardabweichung		2,006	2,271
Varianz		4,026	5,159
Schiefe		-,080	-1,451
Standardfehler der Schiefe		,421	,421
Kurtosis		-,277	2,013
Standardfehler der Kurtosis		,821	,821
Minimum		2	1
Maximum		10	10

a Mehrere Modi vorhanden. Der kleinste Wert wird angezeigt.

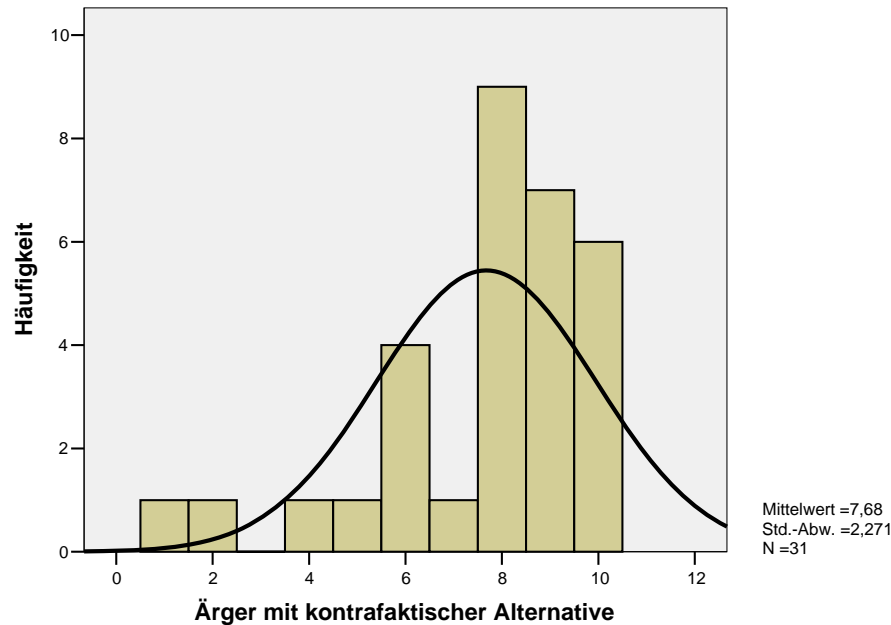
Wir sind nun sehr gespannt auf die Verteilung der Ärgermessung in der Situation *mit* kontrafaktischer Alternative (AERGM), weil sich ein KFA-Effekt in der erwarteten Richtung hier deutlich abzeichnen sollte. Es ist generell zu empfehlen, sich mit möglichst einfachen Grafiken und Sta-



tistiken ein präzises Bild von der Effektlage zu verschaffen, statt einem Signifikanztest blind zu vertrauen, der eventuell durch technische Fehler belastet ist.

Im Vergleich zur „neutralen“ Ärgerverteilung von AERGO (mit dem Mittelwert 5,68) zeigt sich bei AERGM eine dramatisch andere Verteilung (mit dem Mittelwert 7,68):

#### Ärger mit kontrafaktischer Alternative



Wir sehen einen mittleren Ärgeranstieg um 20° (bei Rückübersetzung in die Celsius-Skala des Fragebogens). Außerdem ist die AERGM-Verteilung am rechten Rand konzentriert und deutlich verschieden von einer Normalverteilung, was sich auch in signifikanten Ergebnissen der Tests zu Schiefe und Kurtosis widerspiegelt:

$$\frac{|\text{Schiefe}|}{\text{SF}(\text{Schiefe})} = 3,447 > 1,96$$

$$\frac{|\text{Kurtosis}|}{\text{SF}(\text{Kurtosis})} = 2,451 > 1,96$$

Hier sind *zweiseitige* Tests durchzuführen, weil keine gerichteten Hypothesen vorlagen. Wir haben zwar eine explizite Hypothese über die Richtung des KFA-Effekts (vgl. Abschnitt 1.3.2), doch muss die Verschiebung einer Verteilung nach rechts keinesfalls zu einer negativen Schiefe führen. Offenbar ist aber der KFA-Effekt so stark, dass er die Ärgerverteilung an die „Decke“ geschoben und damit rechtssteil (negativ schief) gemacht hat.

#### b) Ergebnis der Fehleranalyse

Unsere Fehleranalyse liefert nur einen „Treffer“. In der Häufigkeitstabelle zur Variablen LOT10 entdecken wir den verbotenen Wert Null:


LOT10

	Häufigkeit	Prozent	Gültige Prozente	Kumulierte Prozente
Gültig 0	1	3,2	3,2	3,2
1	4	12,9	12,9	16,1
2	10	32,3	32,3	48,4
3	9	29,0	29,0	77,4
4	7	22,6	22,6	100,0
Gesamt	31	100,0	100,0	

Diese Fehlerquote kann als erfreulich niedrig eingestuft werden.

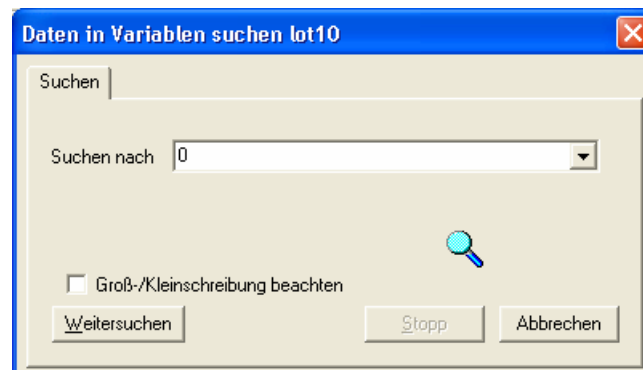
#### 4.7 Suche nach Daten

In der Häufigkeitstabelle zu LOT10 haben wir den unzulässigen Wert Null (mit Häufigkeit 1) entdeckt. Nun möchten wir natürlich sofort wissen, bei welchem Fall dieser Wert auftritt, um geeignete Korrekturen vornehmen zu können. Der betroffene Fall ist sehr leicht zu ermitteln:

- Holen Sie nötigenfalls das Datenfenster in den Vordergrund.
- Markieren Sie in der **Datenansicht** eine beliebige Zelle der Variablen LOT10.
- Klicken Sie auf das Symbol , oder wählen Sie den Menübefehl:

##### Bearbeiten > Suchen...

Dann erscheint die folgende Dialogbox:



- Tragen Sie den zu suchenden Wert ein, und klicken Sie auf den Schalter **Weitersuchen**. Für die Suche nach SYSMIS ist ein Leerzeichen einzutragen.
- Daraufhin markiert SPSS die erste Trefferzelle, und Sie kennen den Fall mit fehlerhaftem LOT10-Wert: Es ist zufällig der erste Fall (FNR = 1), dessen ausgefüllter Fragebogen im Manuskript wiedergegeben ist (siehe Seite 25), so dass Sie den korrekten Wert ablesen und im Datenfenster eintragen können. Nach dieser Datenkorrektur sollten Sie die Arbeitsdatei sichern und damit die SPSS-Datendatei **kfar.sav** auf den neuen Stand bringen.

#### 4.8 Arbeiten mit dem Ausgabefenster (Teil II)

Weil es sich beim SPSS Viewer um eine komplexe Anwendung handelt, wird ihre umfangreiche Funktionalität in mehreren Portionen präsentiert.

##### 4.8.1 Nachbearbeitung von Tabellen

Sie werden noch sehr flexible Möglichkeiten zum Umstrukturieren („Pivotieren“) von Tabellen mit dem so genannten Pivot-Editor kennen lernen (z.B. Zeilen- und Spaltendimension vertauschen). Zunächst beschränken wir uns auf Gestaltungsmöglichkeiten, die das Erscheinungsbild

einer Tabelle beeinflussen, ohne ihre Grundstruktur zu verändern. Auch für solche Nachbearbeitungen ist der Pivot-Editor zuständig.

Als Beispiel soll im Folgenden die Häufigkeitstabelle zur Fachbereichsvariablen verwendet werden:

Fachbereich an der Universität Trier

	Häufigkeit	Prozent	Gültige Prozente	Kumulierte Prozente
Gültig I	19	61,3	61,3	61,3
III	2	6,5	6,5	67,7
IV	6	19,4	19,4	87,1
VI	4	12,9	12,9	100,0
Gesamt	31	100,0	100,0	

#### 4.8.1.1 Pivot-Editor starten

Um das Editieren einer Tabelle zu beginnen, können Sie einen Doppelklick darauf setzen oder die Option **Objekt: SPSS Pivot-Tabelle** aus ihrem Kontextmenü wählen. Bei der letztgenannten Methode bietet ein Untermenü die Auswahl zwischen dem **Bearbeiten** innerhalb des Viewers (*in-place-editing*) und dem **Öffnen** eines separaten Fensters für das Editieren der Tabelle.

#### 4.8.1.2 Modifikation von Zellinhalten

##### a) Text editieren

Bei aktivem Pivot-Editor können Sie nach einem Doppelklick auf eine Zelle den enthaltenen Text beliebig ändern. Wir wollen den Titel und die Spaltenbeschriftungen ändern sowie das Wort *Gültig* am linken Rand der Tabelle löschen:

Fachbereiche im SPSS-Kurs

	n	%	gültige %	kum %
I	19	61,3	61,3	61,3
III	2	6,5	6,5	67,7
IV	6	19,4	19,4	87,1
VI	4	12,9	12,9	100,0
Gesamt	31	100,0	100,0	

Mit der Pivot-Funktion **Gruppierung aufheben** werden wir übrigens später eine Möglichkeit kennen lernen, die überflüssige Zelle mit der Beschriftung „Gültig“ komplett zu entfernen.

##### b) Zellen zur weiteren Bearbeitung markieren

Mit dem Menübefehl **Bearbeiten > Auswählen** lassen sich Tabellenbestandteile (z.B. Tabellenkorpus, Datenzellen) zur weiteren Bearbeitung markieren. Außerdem stehen die windows-üblichen Markierungsmethoden per Maus und Tastatur zur Verfügung.

##### c) Schriftmerkmale

Für eine oder mehrere markierte Zellen kann man nach **Format > Schriftart...** diverse Schriftmerkmale ändern.

##### d) Zelleneigenschaften

Nach **Format > Zelleneigenschaften** können zahlreiche Attribute der markierten Zellen beeinflusst werden, z.B.:

- Zahlenformate, Anzahl der Dezimalstellen
- Ausrichtung der Zellinhalte
- Randabstände der Zellinhalte
- Schattierung

Mit zentrierten Werten, zwei Dezimalstellen bei den Prozentangaben und rechtsbündig gesetzten Fachbereichsbezeichnungen sieht unsere Beispieltabelle folgendermaßen aus:

Fachbereiche im SPSS-Kurs

	n	%	gültige %	kum %
I	19	61,29	61,29	61,29
III	2	6,45	6,45	67,74
IV	6	19,35	19,35	87,10
VI	4	12,90	12,90	100,00
Gesamt	31	100,00	100,00	

### e) Spaltenbreite

Wenn sich der Mauszeiger über dem rechten Rand einer Spalte befindet, ändert er seine Form zu einem doppelseitigen Pfeil. Jetzt können Sie durch Klicken und Ziehen bei gedrückter linker Maustaste die rechte Spaltenbegrenzung verschieben und somit die Spaltenbreite ändern.

Der Menübefehl

#### Ansicht > Gitterlinien

blendet Hilfslinien an der Stelle unsichtbarer Zellenbegrenzungen ein (bzw. aus) und erleichtert damit die Anpassung der Spaltenbreiten.

In unserer Beispieltabelle kann die erste Spalte eine Schlankeitskurve vertragen:

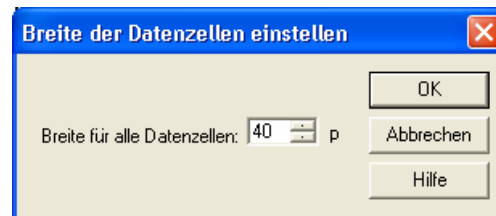
Fachbereiche im SPSS-Kurs

	n	%	gültige %	kum %
I	19	61,29	61,29	61,29
III	2	6,45	6,45	67,74
IV	6	19,35	19,35	87,10
VI	4	12,90	12,90	100,00
Gesamt	31	100,00	100,00	

Über den Menübefehl

#### Format > Breite der Datenzellen...

lässt sich die Breite sämtlicher Datenzellen einer Tabelle numerisch spezifizieren, z.B.:



Nach missratenen Gestaltungsbemühungen bringt eventuell

#### Format > Automatisch anpassen

wieder ein akzeptables Ergebnis zu Stande.

### 4.8.1.3 Tabellenvorlagen

Für eine Pivot-Tabelle kann nach **Format > Tabellenvorlagen...** das Design einer Tabellenvorlage übernommen werden. So sieht unser Beispiel nach Anwendung der Vorlage **Akademisch** aus:

Fachbereiche im SPSS-Kurs

	n	%	gültige %	kum %
I	19	61,29	61,29	61,29
III	2	6,45	6,45	67,74
IV	6	19,35	19,35	87,10
VI	4	12,90	12,90	100,00
Gesamt	31	100,00	100,00	

## 4.8.2 Weitere Gestaltungsmöglichkeiten im Navigationsbereich

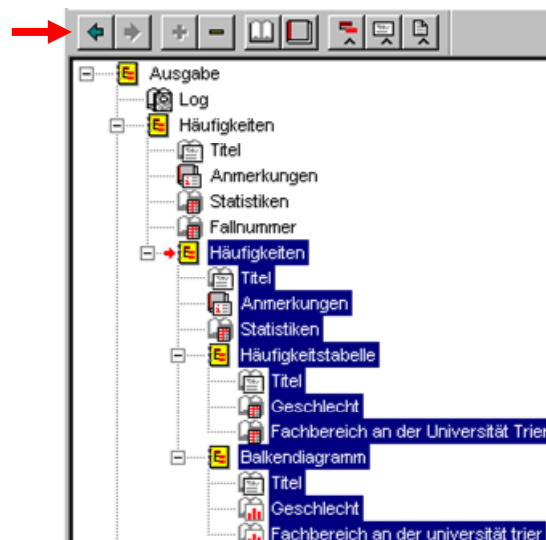
### 4.8.2.1 Blöcke bzw. Teilausgaben kopieren, verschieben oder löschen

Sie können markierte Blöcke bzw. Teilausgaben ...

- Löschen: mit der **Entf**-Taste
  - Kopieren bzw. Verschieben: mit der Maus: Ziehen und Ablegen, beim *Kopieren* zusätzlich *nach* Beginn der Bewegung die **Strg**-Taste drücken mit den Items aus dem Menü **Bearbeiten** oder den äquivalenten Tastenkombinationen: **Kopieren** bzw. **Ausschneiden**, Ziel markieren und **Einfügen**
- via Zwischenablage:

### 4.8.2.2 Befördern und Degradieren

Wenn kopierte oder verschobene Ausgabeblöcke versehentlich auf einer unerwünschten Gliederungsebene gelandet sind, können sie mit den Pfeiltasten oberhalb der Navigationszone „befördert“ oder „degradiert“ werden, z.B.:



Die Ausgabeblöcke in einem Viewer-Dokument müssen nicht unbedingt nebeneinander auf derselben Gliederungsebene liegen, sondern können baumartig angeordnet werden. Von dieser Strukturierungsmöglichkeit macht z.B. auch die SPSS-Prozedur zur Häufigkeitsanalyse Gebrauch.

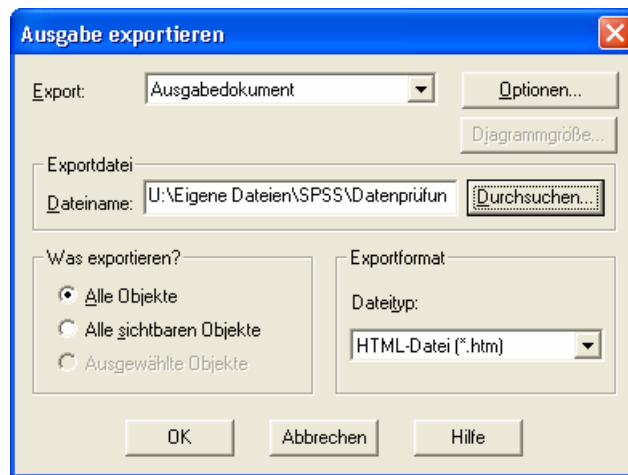
## 4.8.3 Ausgaben exportieren

Pivot-Tabellen, Diagramme und sonstige Ausgaben können in diversen Formaten (z.B. HTML, MS-Word/RTF, Text) exportiert werden. So lassen sich z.B. Ergebnispakete in elektronischer Form an Mitglieder einer Arbeitsgruppe übergeben, die über keine passende SPSS-Version zum Öffnen der Ausgabedateien (Namenserweiterung **spo**) verfügen.

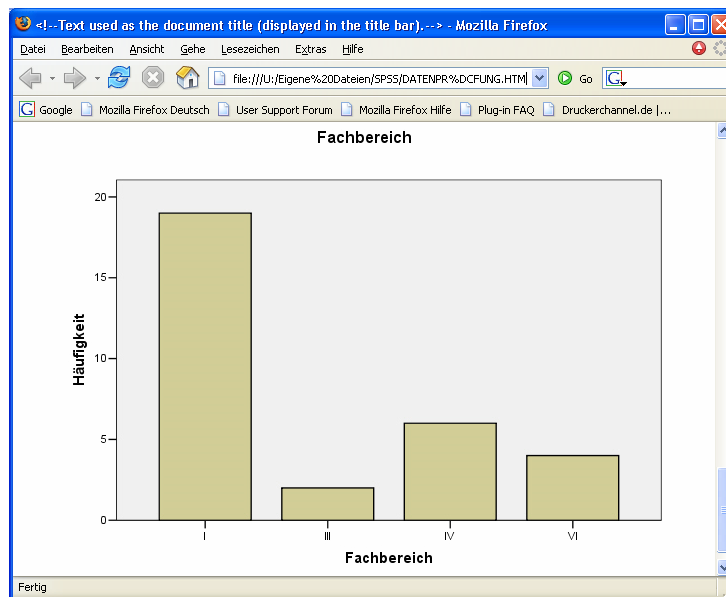
Der Export wird angefordert mit

**Datei > Exportieren...**

Mit folgender Dialogbox wird z.B. das gesamte Viewer-Dokument im HTML-Format exportiert:



So sieht das FB-Balkendiagramm nach dem HTML-Export im Firefox-Browser aus:



Beim Exportumfang gibt es folgende Alternativen:

- **Ausgabedokument**
- **Ausgabedokument (ohne Diagramme)**
- **Nur Diagramme**

Dann sind folgende **Dateitypen** zulässig: EMF, CGM, JPG, PCT, PNG, EPS, TIF, BMP, WMF

Für jedes zu exportierende Diagramm wird eine eigene Datei erstellt. Beim Exportumfang **Ausgabedokument** können die oben genannten Dateiformate (CGM, JPG etc.) in der **Optionen**-Subdialogbox eingestellt werden. In Abhängigkeit vom gewählten Grafik-Dateityp sind für den Export von Diagrammen weitere Optionen vorhanden, z.B. zur Größe und Farbumsetzung.

Beim Export für MS-Word erhält man Tabellen im Format dieses Textverarbeitungsprogramms, die also in Word beliebig modifiziert werden können (vgl. Abschnitt 4.4.4).

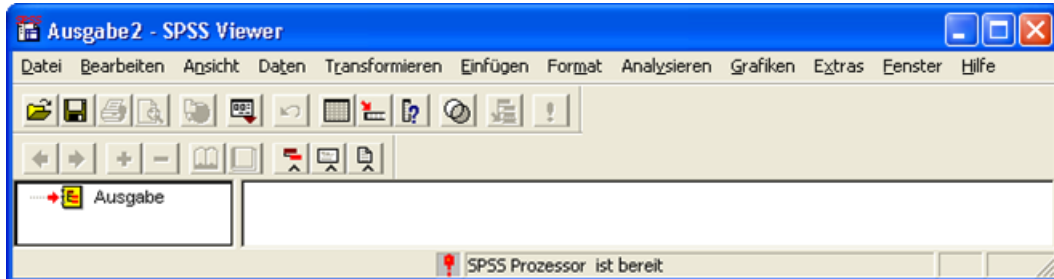
Auf den Pool-PCs an der Universität Trier können Sie den Inhalt des Ausgabefensters als PDF-Datei exportieren, indem Sie nach **Datei > Drucken** den Drucker mit dem Namen **Rumborak PDF Writer Plus** wählen.


#### 4.8.4 Mehrere Ausgabefenster verwenden

Bislang war immer von *dem* Ausgabefenster die Rede. Im Verlauf einer längeren Auswertungsarbeit kann es der Übersichtlichkeit halber sinnvoll sein, ein zusätzliches Ausgabefenster anzufordern. Dazu dient der Menübefehl:

##### **Datei > Neu > Ausgabe**

Wenn mehrere Ausgabefenster vorhanden sind, muss geregelt werden, in welches Fenster SPSS zukünftige Ausgaben schreiben soll. Daher ist stets ein *Haupt*ausgabefenster festgelegt, das durch ein Ausrufezeichen in seiner Statuszeile gekennzeichnet ist, z.B.:



Außerdem ist der Ausrufezeichen-Schalter  in der *Symbolleiste* des Hauptfensters notwendigerweise inaktiv. Dieser Schalter dient nämlich ggf. dazu, ein Ausgabefenster zum Hauptfenster zu ernennen.

Um ein bestimmtes Ausgabefenster in den Vordergrund zu holen, können Sie es anklicken oder das **Fenster**-Menü eines beliebigen SPSS-Fensters benutzen.

Jedes Ausgabefenster kann auf windows-übliche Weise geschlossen werden, z.B. indem Sie es in den Vordergrund holen und dann anordnen:

##### **Datei > Schließen**

---

## 5 Speichern der SPSS-Kommandos zu wichtigen Anweisungsfolgen

### 5.1 Zur Motivation

Eventuell möchten Sie nach zahlreichen Datenkorrekturen alle Testprozeduren erneut durchführen, um ein beruhigendes Ergebnis *Null Fehler* zu sehen. Leider müssen dazu zahlreiche Dialogboxen erneut ausgefüllt und abgeschickt werden. Eventuell erhalten Sie nach Abschluss der Fehlerkontrolle noch weitere bearbeitete Fragebögen. Sie freuen sich natürlich über die Stichprobenerweiterung und erfassen sofort die neuen Fälle. Dann allerdings fällt Ihnen ein, dass nun alle Kontrollanalysen nochmals wiederholt werden müssen.

Um solchen Frust zu vermeiden, brauchen wir eine Möglichkeit, aufwändige und potentiell mehrfach benötigte Anweisungssequenzen zur späteren Wiederverwendung abzuspeichern. In SPSS eignen sich dazu in natürlicher Weise die **Kommandos**, die den einzelnen Dialogboxen zugrunde liegen, und die von SPSS stets im Hintergrund erzeugt und ausgeführt werden, wenn wir eine ausgefüllte Dialogbox mit **OK** abschicken.

In diesem Zusammenhang lohnt ein kurzer Blick auf die Architektur des SPSS-Systems, das aus den beiden folgenden Komponenten besteht:

- **Benutzerschnittstelle**

Wir interagieren mit der Benutzerschnittstelle, die unsere Anweisungen entgegennimmt und die Ergebnisse präsentiert. Wir können der Benutzerschnittstelle unsere Anweisungen in Form von ausgefüllten Dialogboxen oder als Folge von SPSS-Kommandos übergeben.

- **SPSS-Prozessor**

Die Benutzerschnittstelle gibt unsere Anweisungen in jedem Fall in Form von SPSS-Kommandos an den Prozessor weiter, der im Hintergrund arbeitet. Wir erfahren übrigens in der Statuszeile der SPSS-Fenster, was der Prozessor gerade treibt. Da wir den Prozessor bislang nur minimal belastet haben, hat die Statuszeile meistens angezeigt: **SPSS Prozessor bereit**. Während der Prozessor arbeitet, wird in der Statuszeile protokolliert, mit welchem SPSS-Kommando er gerade beschäftigt ist. Nach dem Abschicken einer Häufigkeitsdialogbox erscheint z.B. **Ausführen: FREQUENCIES**, bei unserem kleinen Datensatz allerdings nur sehr kurz. Wenn wir eine ausgefüllte Häufigkeitsdialogbox mit **OK** quittieren, führt der SPSS-Prozessor also im Hintergrund das korrespondierende FREQUENCIES-Kommando aus.

In fast allen SPSS-Dialogboxen kann man über die Standardschaltfläche **Einfügen** die zugrunde liegenden SPSS-Kommandos produzieren lassen. Diese werden dann *nicht* ausgeführt, sondern in ein so genanntes **Syntaxfenster** übertragen, das weitgehend analog zu einem Texteditor funktioniert. Hier kann man alle Kommandos zu einer Sequenz ansammeln, nach Bedarf einzeln oder geschlossen ausführen lassen und schließlich in einer Datei abspeichern. Später kann man die Kommandos aus dieser Datei wieder laden und, eventuell nach manueller Überarbeitung, erneut ausführen lassen. Das genaue Vorgehen wird in Abschnitt 5.2 an einem konkreten Beispiel geübt.

Eine Folge von SPSS-Kommandos kann man (leicht hochstaplerisch) als **SPSS-Programm** bezeichnen. In fast jedem Projekt sollte es mindestens *ein* SPSS-Programm geben, nämlich das bereits in Abschnitt 3.2.6 vorgeschlagene Transformationsprogramm, das aus der Rohdatendatei durch diverse Transformationen die Fertigdatendatei des Projektes erstellt. Wir werden für unser KFA-Projekt ein solches Programm in Abschnitt 6 erstellen.



Ob sich bei einer konkreten Anweisungssequenz das Abspeichern als SPSS-Programm lohnt, muss von Fall zu Fall entschieden werden. Bei kurzen, simplen Sequenzen mit geringer Wiederholungswahrscheinlichkeit ist ein Konservieren unrentabel.

Es soll nicht verschwiegen werden, dass die Ausführung einer Anweisungssequenz mit dem Umweg über ein Syntaxfenster geringfügig mehr SPSS-Kenntnisse erfordert als die direkte Ausführung durch Quittieren der Dialogboxen mit **OK**. Wer sich beim Umgang mit SPSS-Kommandos unsicher fühlt, bei seinem relativ kleinen Projekt eventuell erforderliche Wiederholungen von Dialogbox-Sequenzen nicht scheut und das Risiko inkonsistenter Datenzustände durch große Sorgfalt kontrolliert, der kann auf das Erzeugen und Abspeichern von SPSS-Kommandos verzichten.

Für ambitionierte SPSS-Anwender(innen) muss noch klargestellt werden, dass die Erstellung, Überarbeitung und Ausführung von Programmen in einem Syntaxfenster eine eigenständige Methode der SPSS-Benutzung darstellt, über die fast alle Analyse-Funktionen erreichbar sind. Viele SPSS-Leistungen stehen sogar *ausschließlich* über die Syntax zur Verfügung, z.B.:

- Conjoint-Analyse
- Kontrollstrukturen wie z.B. DO REPEAT - Schleifen, mit denen man komplexe Datentransformationen auf effiziente Weise durchführen kann.
- Die MATRIX-Programmiersprache, mit der man eigene Statistikprozeduren programmieren kann.

Der Hersteller SPSS Inc. meint im Hilfesystem zu der Debatte „Dialogbox kontra Programm“:

„Erfahrene SPSS-Anwender bevorzugen möglicherweise die rationellere Befehlssprache.“

Im aktuellen Abschnitt 5 werden der Einfachheit halber nur sehr oberflächliche Hinweise zur Kommandosprache gegeben. Diese sollten genügen für Anwender, die nicht frei programmieren, sondern nur gelegentlich ein von SPSS automatisch erzeugtes Kommando modifizieren wollen. Der Anhang enthält für ambitionierte SPSS-Anwender eine ausführlichere Beschreibung der Kommandosprache. Eine vollständige Dokumentation auf ca. 2000 Seiten finden Sie als PDF-Dokument im Hilfesystem von SPSS 13 über

### Hilfe > Command Syntax Reference

Wie schon erwähnt, sind die Dialogboxen beim Erstellen eines SPSS-Programms sehr nützlich. Mit Hilfe der bislang ignorierten Standardschaltfläche **Einfügen** kann nämlich die zu einer Dialogbox-Bearbeitung äquivalente Kommandofolge in ein Syntaxfenster übertragen werden. Sie müssen sich also nicht zwischen zwei unvereinbaren SPSS-Bediensystemen entscheiden, sondern sollten eine möglichst effiziente Kombination beider Methoden verwenden.

## 5.2 Dialogunterstützte Erstellung von SPSS-Programmen

Das folgende SPSS-Programm führt für unser KFA-Projekt die Häufigkeitsanalysen zur Fehlersuche bei den Variablen FNR, GESCHL und FB durch (siehe Abschnitt 4):

```

GET
  FILE='U:\Eigene Dateien\SPSS\kfar.sav'.
FREQUENCIES
  VARIABLES=fnr
  /STATISTICS=MINIMUM MAXIMUM MODE
  /ORDER= ANALYSIS.
FREQUENCIES
  VARIABLES=geschl fb
  /BARCHART  FREQ
  /ORDER= ANALYSIS.

```

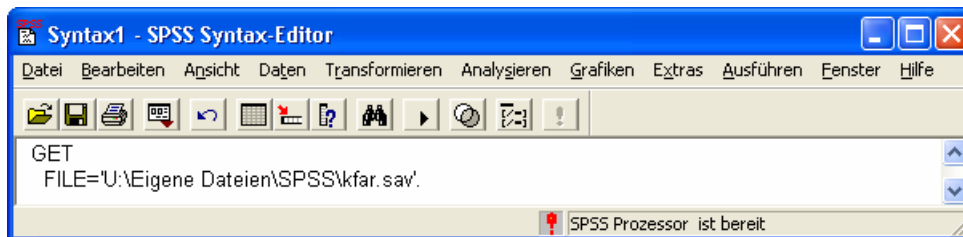
Wir werden dieses Programm gleich „vollautomatisch“ mit drei Mausklicks auf **Einfügen**-Schalter produzieren und dabei auch seine Bestandteile kurz erläutern. Als Ausgangssituation für die anschließenden Erläuterungen wird eine aktive SPSS-Sitzung mit leerem Datenfenster angenommen. Starten Sie nötigenfalls SPSS bzw. entleeren Sie das Datenfenster mit:

### Datei > Neu > Daten

Rufen Sie die Dialogbox zum Öffnen einer Datendatei auf:


### Datei > Öffnen > Daten

Schreiben oder klicken Sie den Namen Ihrer Rohdatendatei in das entsprechende Feld, und betätigen Sie dann den Schalter **Einfügen**. Daraufhin beginnt SPSS *nicht* damit, aus der angegebenen Datendatei eine neue Arbeitsdatei zu erstellen und diese im Datenfenster anzuzeigen, sondern SPSS schreibt das für diese Aktion zuständige GET-Kommando in ein Syntaxfenster mit dem Titel **Syntax1**:



Der Aufbau des GET-Kommandos ist sehr einfach:

- Es beginnt mit dem Kommandonamen GET.
- Im FILE-Subkommando wird die zu öffnende Datei spezifiziert.
- Am Ende muss wie bei jedem SPSS-Kommando ein **Punkt** stehen.

Weil das Datenfenster momentan noch leer ist, stehen die Menübefehle zum Anfordern von Statistik- und Grafikprozeduren nicht zur Verfügung. Daher wollen wir jetzt das GET-Kommando ausführen lassen, um die Daten einzulesen. Setzen Sie dazu die Schreibmarke an eine beliebige Position *innerhalb des GET-Kommandos*, und klicken Sie auf das Symbol . Nun erstellt SPSS eine neue Arbeitsdatei mit den Rohdaten und zeigt diese im Datenfenster an.


Spezifizieren Sie jetzt mit Hilfe der zuständigen Dialogbox dieselbe Häufigkeitsanalyse zur FNR-Variablen wie in Abschnitt 4.3. Verlassen Sie die Dialogbox jedoch nicht mit **OK**, sondern mit **Einfügen**. Daraufhin erscheint am Ende des Syntaxfensters ein FREQUENCIES-Kommando (siehe oben):

- Es beginnt mit dem Kommandonamen FREQUENCIES.
- Im VARIABLES-Subkommando ist angegeben, welche Variable analysiert werden soll.
- Im STATISTICS-Subkommando ist angegeben, welche Verteilungskennwerte berechnet werden sollen.

- Das (im vorliegenden Fall irrelevante) ORDER-Subkommando entscheidet bei der Analyse *mehrerer* Variablen darüber, ob die Statistiken für jede Variable in einer eigenen Tabelle oder für alle Variablen in einer gemeinsamen Tabelle erscheinen sollen. Um diese Entscheidung in der **Häufigkeiten**-Dialogbox zu treffen, müssen Sie übrigens die **Format**-Subdialogbox öffnen und im Rahmen **Mehrere Variablen** die passende Option wählen.
- Das FREQUENCIES-Kommando wird wie jedes SPSS-Kommando durch einen **Punkt** abgeschlossen.

Produzieren Sie als nächstes die Syntax zu der in Abschnitt 4.5 durchgeführten Häufigkeitsanalyse für die Variablen GESCHL und FB.

Nun sollte Ihr Syntaxfenster den zu Beginn des Abschnitts wiedergegebenen Inhalt haben. Das GET-Kommando ist schon gelaufen, folglich müssen Sie noch die beiden FREQUENCIES-Kommandos ausführen lassen. Weil es sich um *zwei* Kommandos handelt, müssen Sie folgendermaßen vorgehen:

- Markieren Sie zunächst per Maus *die beiden* auszuführenden Kommandos.
- Klicken Sie dann auf das Symbol , oder drücken Sie die Tastenkombination **Strg+R**. Daraufhin werden alle Kommandos im Syntaxfenster ausgeführt, die (zumindest teilweise) markiert sind.

Im Ausgabefenster protokolliert SPSS übrigens zu jeder Analyseanforderung in der zunächst zugeklappten Teilausgabe **Anmerkungen** u.a. die zugrunde liegende Syntax, z.B.:

#### Anmerkungen

Ausgabe erstellt		18-OCT-2005 18:00:28
Kommentare		
Eingabe	Daten	U:\Eigene Dateien\SPSS\kfar.sav
	Filter	<keine>
	Gewichtung	<keine>
	Aufgeteilte Datei	<keine>
	Anzahl der Zeilen in der Arbeitsdatei	31
Behandlung fehlender Werte	Definition von fehlenden Werten	Benutzerdefinierte fehlende Werte werden als fehlend behandelt.
	Verwendete Fälle	Statistik basiert auf allen Fällen mit gültigen Daten.
Syntax		FREQUENCIES VARIABLES=geschl fb /BARChart FREQ /ORDER= ANALYSIS .
Ressourcen	Verstrichene Zeit	0:00:00,59
	Zugelassene Werte	224841

Damit sich durch spätere Wiederverwendung der SPSS-Kommandos der Rationalisierungseffekt der programm-orientierten Arbeitsweise einstellen kann, müssen Sie Ihr SPSS-Programm sichern.

Wechseln Sie dazu nötigenfalls in das Syntaxfenster, und wählen Sie den Menübefehl:


#### **Datei > Speichen unter...**

Verwenden Sie im Dateinamen die vorgeschlagene Erweiterung **sps**, indem Sie *keine* Erweiterung angeben.

Wenn Sie später dieselbe Auswertung nochmals benötigen, dann müssen Sie lediglich das vorhandene Programm mit dem Menübefehl:

### **Datei > Öffnen > Syntax**

laden und ausführen lassen. Um die Ausführung *sämtlicher* Kommandos in einem Syntaxfenster anzuordnen, haben Sie folgende Möglichkeiten:



- Menübefehl **Ausführen > Alles**
- Alle Kommandos markieren (z.B. mit **Strg+A**) und die Ausführung anfordern (z.B. per Mausklick auf das Symbol  oder mit der Tastenkombination **Strg+R**)

## **5.3 Arbeiten mit dem Syntax-Fenster**

Das Syntaxfenster bietet die Funktionalität eines einfachen Texteditors, so dass man automatisch erstellte SPSS-Kommandos leicht modifizieren kann, um z.B. die in einer Statistikprozedur zu analysierenden Variablen auszutauschen.

Man kann ein neues Syntaxfenster auch unabhängig vom **Einfügen**-Schalter direkt anfordern mit:

### **Datei > Neu > Syntax**

Wenn *mehrere* Syntaxfenster vorhanden sind, muss geregelt werden, in welches Fenster SPSS die per **Einfügen**-Schalter automatisch erzeugten Kommandos übertragen soll. Dies geschieht genauso wie bei den Ausgabefenstern: Ein Mausklick auf den aktiven Schalter  in seiner Symbolleiste macht ein Syntaxfenster zum Hauptfenster. Ein passiver (nicht verwendbarer) Schalter  signalisiert ebenso wie ein Ausrufezeichen in der *Statuszeile*, dass es sich um das Hauptfenster handelt.

Um ein bestimmtes Syntaxfenster in den Vordergrund zu holen, können Sie es anklicken oder das **Fenster**-Menü eines beliebigen SPSS-Fensters benutzen.

Jedes Syntaxfenster kann auf windows-übliche Weise geschlossen werden, z.B. indem Sie es in den Vordergrund holen und dann anordnen:

### **Datei > Schließen**

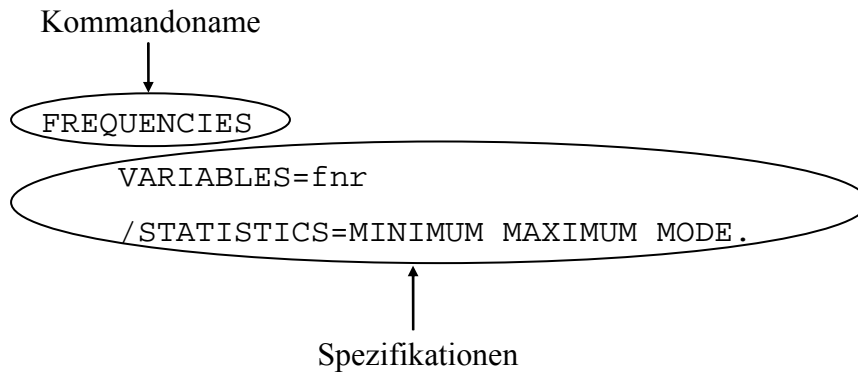
Wenn Sie längere Zeit mit SPSS arbeiten, wird sich vermutlich Ihr Umgang mit SPSS-Syntax in folgenden Stufen weiterentwickeln:

- Kommandos automatisch erzeugen lassen und später wieder verwenden  
Bei dieser Arbeitsweise müssen Sie nur wissen, wie man SPSS-Kommandos per Dialogbox in ein Syntaxfenster befördert, und wie man überflüssige Kommandos löscht.
- Automatisch erzeugte Kommandos modifizieren  
Es zeigt sich, dass SPSS-Kommandos meist leicht zu durchschauen und zu modifizieren sind.
- Freies Programmieren

## **5.4 Elementare Regeln zur SPSS-Syntax**

Für den im Kurs vorgeschlagenen Einsatz von SPSS-Kommandos sollte die Kenntnis der folgenden Regeln genügen:

- Ein Kommando besteht aus seinem Namen und den Spezifikationen, die sich aus Schlüsselwörtern (z.B. VARIABLES, STATISTICS), Variablennamen usw. zusammensetzen, z.B.:



- Zwei Elemente der Kommandosprache sind durch mindestens ein Leerzeichen oder durch einen Zeilenwechsel voneinander zu trennen. Manche Zeichen mit festgelegter Bedeutung wie z.B. "=", "/", "(", "+", ">" sind aber selbstbegrenzend, d.h. davor und danach sind keine Leerzeichen nötig (aber erlaubt).
- Ein Kommando kann sich über beliebig viele Fortsetzungszeilen erstrecken, dabei dürfen aber *innerhalb* des Kommandos keine Leerzeilen auftreten. Diese signalisieren nämlich per Voreinstellung (wie der Punkt) das Ende des Kommandos.
- Zwischen zwei Kommandos dürfen beliebig viele Leerzeilen stehen, was eine übersichtliche Gestaltung von SPSS-Programmen erlaubt.
- **Jedes Kommando muss in einer neuen Zeile beginnen und mit einem Punkt enden.**

Gut kommentierte Programme sind später leichter zu verstehen. Die SPSS-Syntax bietet zum Kommentieren das Kommando COMMENT, dessen Name durch ein Sternchen ersetzt werden darf:

```
COMMENT kommentar.
```

```
* kommentar.
```

Beachten Sie beim Kommentar-Kommando:

- Es darf sich über beliebig viele Fortsetzungszeilen erstrecken, wobei innerhalb des Kommandos keine Leerzeilen erlaubt sind.
- **Jedes Kommentar-Kommando muss mit einem Punkt abgeschlossen werden.** Wenn Sie den Punkt am Ende vergessen, dann betrachtet SPSS den folgenden Programmtext bis zum nächsten Punkt (oder zur nächsten Leerzeile) als Teil des Kommentars!
- Endet eine Kommentarzeile mit einem Punkt, so betrachtet SPSS das Kommentar-Kommando als abgeschlossen. Wenn Sie einen Punkt als *Satzzeichen* ans Ende einer Kommentarzeile gesetzt haben, dann müssen Sie die nächste Kommentarzeile wieder mit COMMENT oder \* einleiten.
- Punkte innerhalb einer Kommentarzeile sind kein Problem.

Beispiel:     \* Mit diesem Programm wird die Rohdatendatei KFAR.SAV  
              auf Erfassungsfehler untersucht.  
GET  
              FILE='U:\Eigene Dateien\SPSS\KFAR.SAV'.

---

## 6 Datentransformation

### 6.1 Vorbemerkungen

Die zur Untersuchung unserer differentialpsychologischen Hypothese benötigte Optimismus-Variable existiert noch nicht, sondern muss erst aus den 12 LOT-Variablen berechnet werden. Vor dieser Berechnung müssen allerdings die aus messtechnischen Gründen umgepolten (negativ formulierten) LOT-Fragen geeignet rekodiert werden (z.B. Frage 3). Es ist typisch für empirische Studien, dass vor der eigentlichen Auswertung aus den Rohvariablen mit zahlreichen Datentransformationen neue oder modifizierte Fertigvariablen erstellt werden müssen.

In diesem Abschnitt werden Sie häufig benötigte SPSS-Befehle zur Datentransformation kennen lernen. Diese wirken sich auf die Datenmatrix der Arbeitsdatei aus, wo entweder neue Variablen aufgenommen oder vorhandene Variablen verändert werden. Per Voreinstellung werden dabei *alle Fälle gleichermaßen* behandelt.

Man kann die Ausführung einer Datentransformation aber auch von einer Bedingung abhängig machen, so dass nicht mehr alle Fälle davon betroffen sind. Diese Möglichkeit werden wir z.B. dazu verwenden, die MD-Behandlung bei den Motiv-Variablen in Ordnung zu bringen, indem wir genau für *die* Fälle mit

$$\text{MOTIV1} = \text{MOTIV2} = \dots = \text{ANDERE} = 0$$

bei allen genannten Variablen die 0 in SYSMIS umkodieren.

SPSS unterstützt Transformationen für Variablen beliebigen Typs. Wir beschränken uns jedoch auf die besonders wichtigen numerischer Variablen.

#### 6.1.1 Rohdatendatei, Transformationsprogramm und Fertigdatendatei

In Abschnitt 3.2.6 wurde vorgeschlagen, zu jedem Projekt ein SPSS-Transformationsprogramm zu erstellen, dessen Aufgabe darin besteht, ausgehend von der Rohdatendatei alle Fertigvariablen zu entwickeln, die im weiteren Verlauf routinemäßig benötigt werden. *Alle* potentiell relevanten Variablen (roh oder fertig) sollen in einer erweiterten Datendatei gesichert werden, die sich für alle Auswertungsarbeiten eignet<sup>1</sup>. Mit Rücksicht auf diese Idee haben wir die bislang existierende Datendatei mit **kfar.sav** (*r* für *roh*) bezeichnet. Im Namen der Projekt-Fertigdatendatei werden wir das *r* dann weglassen.

Wir werden im Verlauf des aktuellen Abschnitts 6 das SPSS-Transformationsprogramm zu unserem KFA-Projekt sukzessive mit Hilfe verschiedener Dialogboxen erstellen. Dabei ist eine besondere Sorgfalt erforderlich, weil fehlerhafte Anweisungen im Transformationsprogramm schwerwiegende Konsequenzen für die weitere Arbeit haben können.

Weil das Transformationsprogramm eventuell wiederholt benötigt wird, z.B. nach einer Stichprobenerweiterung oder nach einer Fehlerkorrektur in den Rohdaten, muss es ebenso sorgfältig gesichert werden wie die Rohdatendatei. Als Dateinamen wollen wir **kfat.sps** wählen.

Wie in Abschnitt 5.1 ausführlich diskutiert, können Sie alle erforderlichen Transformationen auch durch direkte Ausführung von Dialogboxen (Schalter **OK**) erledigen. Diese Arbeitsweise ist zweifellos für Anfänger leichter zu handhaben als die programmorientierte Methode, bei der

---

<sup>1</sup> Unter gewissen, am ehesten in großen Projekten anzutreffenden Umständen kann es sinnvoll sein, die auszuwertenden Daten in *mehreren* Dateien bereitzuhalten. In der Regel führt das Verteilen der Variablen oder Fälle auf mehrere Dateien früher oder später zu dem Problem, dass sich die in einer Analyse zu vergleichenden Fälle oder Variablen in verschiedenen Dateien befinden. Daher ist unreflektierte Anwendung der allgemeinen Lebensregel „Teile und herrsche!“ auf die Dateiorganisation eines Forschungsprojektes *nicht* zu empfehlen.

mit Hilfe von Dialogboxen zunächst mehrere SPSS-Kommandos in ein Syntaxfenster befördert werden (Schalter **Einfügen**), um sie anschließend ausführen zu lassen. Die direkte Arbeitsweise hat aber folgende Nachteile:

- Beim sukzessiven manuellen Modifizieren der Datendatei geht bei größeren Projekten leicht der Überblick verloren. Z.B. weiß irgendwann von einer bestimmten Variablen niemand mehr, in welchen Zwischenschritten sie aus welchen anderen Variablen berechnet worden ist. Spätestens nach dem Auftreten unerwarteter Ergebnisse muss die *tatsächlich* angewendete Berechnungsvorschrift als mögliche Fehlerquelle überprüft werden. Bei der Verwendung eines Transformationsprogramms ist die Herkunft der abgeleiteten Variablen jedoch stets dokumentiert.
- Sind Wiederholungen von Datentransformationen erforderlich, müssen diese komplett neu spezifiziert werden. Solche Wiederholungen sind z.B. nach einer Datenkorrektur fällig, weil SPSS abgeleitete Variablen **nicht** automatisch anpasst, wenn sich Werte der Ursprungsvariablen ändern. Nach Korrekturen bei den Rohvariablen müssen Sie also alle Datentransformationen wiederholen, in die diese Rohvariablen eingehen. Ein weiterer potentieller Anlass für die Wiederholungen von Datentransformationen ist die Erweiterung der Stichprobe.

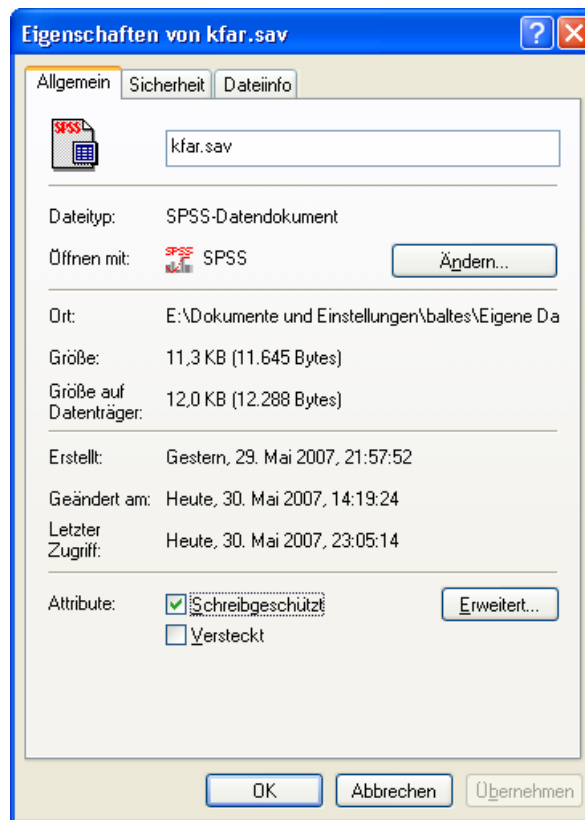
Die für ein Projekt erforderlichen Datentransformationen in Form von SPSS-Anweisungen zu konservieren, lohnt sich meist, denn:

- Die einzelnen Anweisungen sind relativ komplex und damit ebenso fehleranfällig wie zeitaufwändig.
- Es ist relativ wahrscheinlich, dass die gesamte Anweisungsfolge wiederholt durchgeführt werden muss (bei entdeckten Fehlern in den Rohvariablen oder bei einer Stichprobenerweiterung).
- Die Anweisungen zur Datentransformation sind dokumentationspflichtig.

### 6.1.2 Hinweise zum Thema Datensicherheit

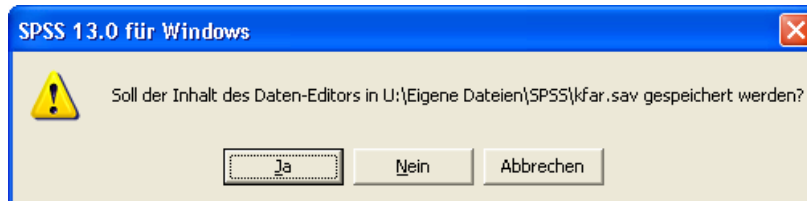
Ihre Rohdaten können nach der sorgfältigen Datenerfassung und -prüfung vorerst als korrekt gelten. Sichern Sie den erreichten Stand, indem Sie die Rohdaten in mindestens **zwei** Dateien speichern (möglichst auf verschiedenen Datenträgern).

Es ist sinnvoll, für beide Dateien das Schreibschutzattribut mit dem Windows-Explorer zu setzen, z.B.:



Vor der geplanten Änderung einer Datei muss das Schreibschutzattribut natürlich wieder aufgehoben werden. Ähnlich sorgfältig sollten Sie nach seiner Fertigstellung das Transformationsprogramm aufbewahren.

Wenn Sie beim Verlassen von SPSS gefragt werden, ob Sie das Daten- oder ein Syntaxfenster sichern wollen, sollten Sie sehr sorgfältig prüfen, ob bei dem entsprechenden Objekt während der Sitzung tatsächlich nur geplante Veränderungen stattgefunden haben.



Antworten Sie im Zweifelsfall mit **Nein**. Möglicherweise haben Sie durch unbeabsichtigte Tastendrucke Daten gelöscht oder verändert. Diese Fehler sollten dann auf keinen Fall auf die Festplatte geschrieben werden.

### 6.1.3 Initialisierung neuer numerischer Variablen

Wenn Sie in einer Datenmodifikationsanweisung die Erstellung einer *neuen* numerischen Variablen anfordern, dann wird die (Fälle  $\times$  Variablen)-Datenmatrix in der Arbeitsdatei um eine Spalte erweitert. SPSS **initialisiert** dabei zunächst die neue Variable, indem es für alle Fälle den globalen MD-Indikator System-Missing als Wert einträgt. Gelingt anschließend die Ermittlung der neuen Variablenausprägung für einen Fall, so wird der Initialwert entsprechend ersetzt. Anderenfalls bleibt System-Missing stehen, so dass der betroffene Fall bei allen Berechnungen mit der neuen Variablen ausgeschlossen wird.



## 6.2 Alte Werte einer Variablen auf neue abbilden (Umkodieren)

Mit dem Befehl **Umkodieren** aus dem Menü **Transformieren** bzw. mit dem äquivalenten RECODE-Kommando können die Werte einer bestehenden Variablen in neue Werte überführt werden. Man kann die Ausgangsvariable verändern oder eine neue Variable mit dem rekodierten Wertevektor erstellen.

### 6.2.1 Das praktische Vorgehen am Beispiel einer künstlichen Gruppenbildung

Da wir im Abschnitt 6 das KFA-Transformationsprogramm sukzessive aufbauen wollen, öffnen wir zunächst unsere Rohdatendatei **kfar.sav**.

Um das Umkodieren zu üben, wählen wir ein mäßig sinnvolles Beispiel aus unserer Studie: Wir konstruieren unter dem Namen DEKADE eine vergrößerte Variante der JahrgangsvARIABLEN, bei der alle in den 60'er Jahren geborenen Personen den Wert 1 und alle in den 70'er Jahren geborenen Personen den Wert 2 erhalten sollen. Wie man sich anhand der Häufigkeitstabelle zur Variablen GEBJ überzeugen kann, ist damit für alle Fälle in unserer Stichprobe ein DEKADE-Wert definiert.

Mit Hilfe der neuen Variablen kann man z.B. den Einfluss des Geburtsjahrzehnts auf diverse abhängige Variablen untersuchen, wobei man sich von der Informationsreduktion (im Vergleich zu GEBJ) keinen allzu großen Nutzen versprechen sollte.

Bei der geplanten Rekodierung wird die (Fälle  $\times$  Variablen)-Datenmatrix um eine neue Variable erweitert, die folgendermaßen aus der vorhandenen Variablen GEBJ entsteht:

GEBJ			DEKADE	
69		→		1
70		→		2
69		→		1
67		→		1
.			.	
.			.	
.			.	
72				2
68		→		1
67		→		1
67		→		1

Wählen Sie den Menübefehl:

#### **Transformieren > Umkodieren**

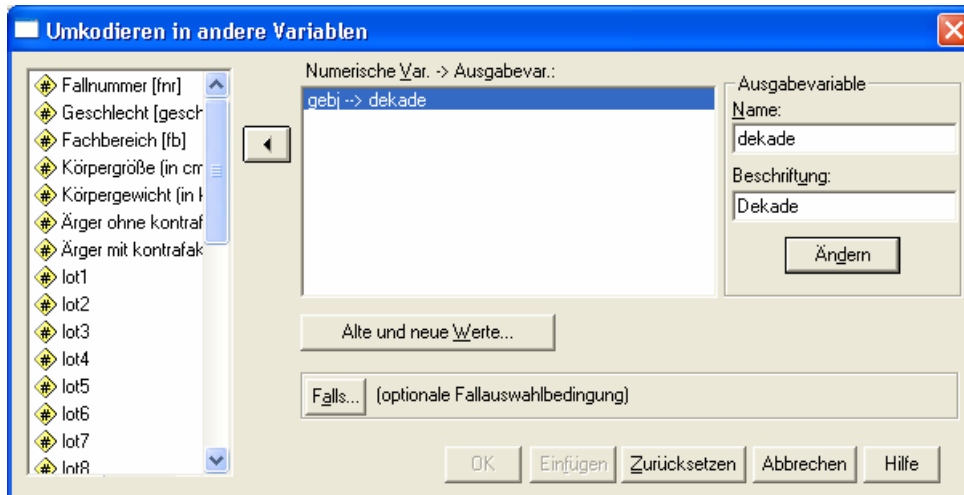
Daraufhin erscheint das folgende Untermenü:

<p><b>In dieselben Variablen...</b>  <b>In andere Variablen...</b></p>
--

Da wir eine neue Variable erzeugen wollen, ist die zweite Alternative zu wählen. Machen Sie folgendermaßen weiter:

- Befördern Sie in der nun erscheinenden Dialogbox **Umkodieren in andere Variablen** die Variable GEBJ in das Feld **Numerische Var. -> Ausgabevar.**
- Tragen Sie im Bereich **Ausgabevariable** den gewünschten Namen DEKADE der neu zu erzeugenden Variablen ein.
- Ergänzen Sie ein Variablenlabel.
- Klicken Sie auf **Ändern**.

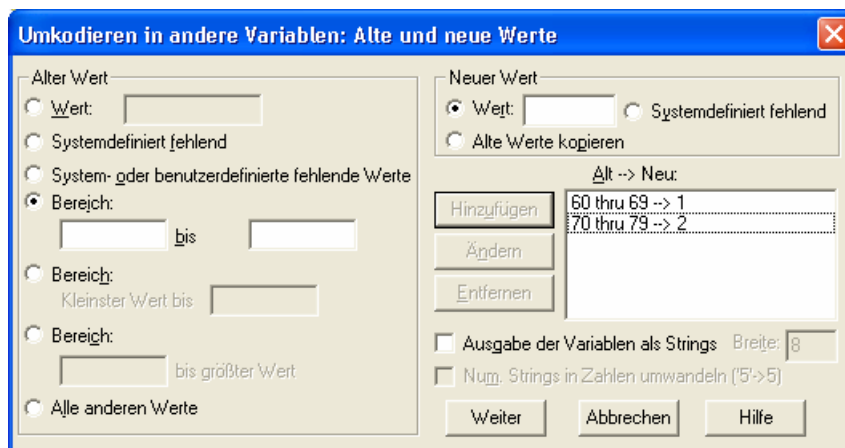
Danach müsste Ihre Dialogbox ungefähr so aussehen:



Legen Sie nun die Abbildungsregeln fest:

- Aktivieren Sie mit dem Schalter **Alte und neue Werte** die Subdialogbox **Umkodieren in andere Variablen: Alte und neue Werte**.
- Geben Sie unter **Alter Wert** den Bereich von 60 bis 69 an, und wählen Sie als zugehörigen **Neuen Wert** die 1.
- Beenden Sie die Definition der ersten Abbildungsvorschrift mit **Hinzufügen**.
- Vereinbaren Sie analog die Zuordnungsvorschrift: „70 bis 79 → 2“.

Jetzt müssten Sie dieses Bild sehen:



Damit ist die Rekodierung vollständig spezifiziert. Quittieren Sie die Subdialogbox mit **Weiter**. Da wir das KFA-Transformationsprogramm sukzessive aufbauen wollen, müssen Sie nun in der Dialogbox **Umkodieren in andere Variablen** auf den Schalter **Einfügen** klicken, um die implizit definierten Kommandos zu produzieren. Wir erhalten ein Syntaxfenster mit folgendem Inhalt:

```

RECODE
  gebj
  (60 thru 69=1) (70 thru 79=2) INTO dekade .
VARIABLE LABELS dekade "Dekade".
EXECUTE .

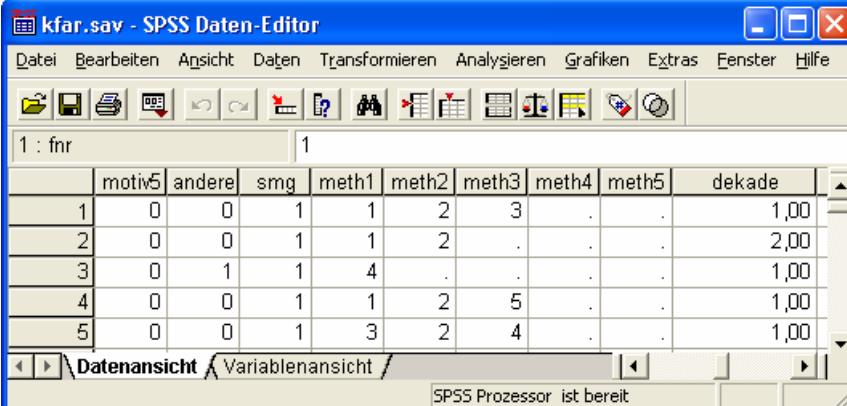
```

Neben dem RECODE-Kommando, das die eigentliche Umkodierung bewirkt, hat SPSS noch zwei weitere Kommandos erzeugt:

- **VARIABLE LABELS**  
Mit diesem Kommando wird das Label für die neue Variable in den Deklarationsteil der Arbeitsdatei eingetragen.
- **EXECUTE**  
Die Rolle dieses Kommandos wird in Abschnitt 6.3 erläutert.

Offenbar hat SPSS unsere Angaben nur in leicht verständliche, englischsprachige Formulierungen übersetzt, so dass Sie es eigentlich wagen können, die Kommandos bei Bedarf auch in abgeänderter Form zu verwenden.

Unabhängig von den guten Argumenten für das Transformationsprogramm gibt es in Ihrer aktuellen Lernphase einen Grund, die obige **Umkodieren**-Dialogbox per **OK**-Schalter zu quittieren oder die zugehörigen Kommandos jetzt schon ausführen zu lassen: Sie können den Effekt auf die Arbeitsdatei sofort beobachten, statt bis zum Abschicken des kompletten Transformationsprogramms warten zu müssen. Weil keine Konflikte mit unserer langfristigen Strategie zu befürchten sind, kehren wir zur **Umkodieren**-Dialogbox zurück und quittieren Sie mit **OK**. Anschließend befindet sich am rechten Rand der Arbeitsdatei die neue Variable DEKADE:



The screenshot shows the SPSS Daten-Editor window for 'kfar.sav'. The data table is displayed in 'Datenansicht' (Data View) mode. The table has 10 columns: 'fnr', 'motiv5', 'andere', 'smg', 'meth1', 'meth2', 'meth3', 'meth4', 'meth5', and 'dekade'. The 'dekade' column contains values 1, 2, 1, 1, and 1 for rows 1 through 5, respectively. The status bar at the bottom indicates 'SPSS Prozessor ist bereit'.

fnr	motiv5	andere	smg	meth1	meth2	meth3	meth4	meth5	dekade
1	0	0	1	1	2	3	.	.	1,00
2	0	0	1	1	2	.	.	.	2,00
3	0	1	1	4	.	.	.	.	1,00
4	0	0	1	1	2	5	.	.	1,00
5	0	0	1	3	2	4	.	.	1,00

## 6.2.2 Technische Details

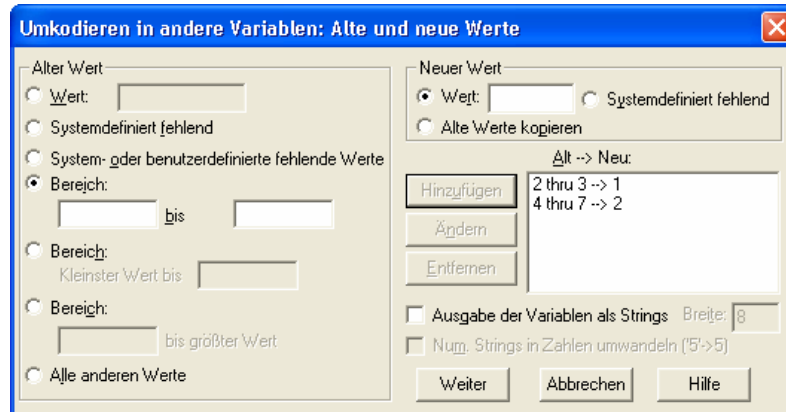
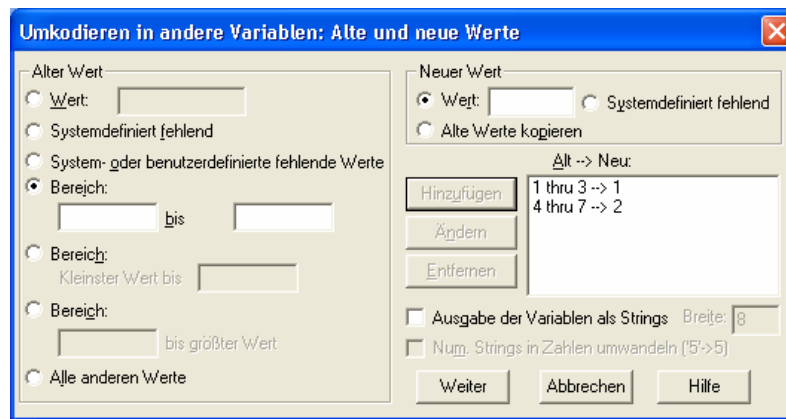
Obwohl das Umkodieren eine sehr simple Datentransformation ist, sind bei der praktischen Anwendung doch einige technische Details zu beachten:

- Sie können bei einem Einsatz der Dialogbox **Umkodieren in andere Variablen** beliebig viele Variablen gleichzeitig umkodieren.
- Bei der Spezifikation der alten Werte, die auf einen neuen Wert abgebildet werden sollen, können Sie angeben:
  - Einen einzelnen **Wert**
  - **Systemdefiniert fehlend** (SYSMIS )  
So ist es also möglich, den automatischen MD-Indikator auf einen anderen Wert umzusetzen.
  - **System- oder benutzerdefinierte fehlende Werte**  
Alle MD-Indikatoren werden umgesetzt.

- Den **Bereich** von einem ersten Wert bis zu einem zweiten Wert (inklusive Grenzwerte)
  - Den **Bereich** vom kleinsten Wert in der Stichprobe bis zu einem bestimmten Wert (inklusive Grenzwert)
  - Den **Bereich** von einem bestimmten Wert bis zum größten Wert in der Stichprobe (inklusive Grenzwert)
  - **Alle anderen Werte**  
Damit sind alle in keiner anderen Ersetzungsvorschrift genannten Werte angesprochen (inklusive MD-Indikatoren, auch System-Missing). **Alle anderen Werte** kann nur in *einer* Ersetzungsvorschrift angegeben werden. Diese wird von SPSS in der Liste aller Ersetzungsvorschriften automatisch an die letzte Stelle gesetzt und damit bei der Kommando-Ausführung zuletzt abgearbeitet.
- Als neuen Wert, auf den die alten Werte einer Ersetzungsvorschrift abgebildet werden sollen, können Sie angeben:
    - Einen **Wert**
    - **Systemdefiniert fehlend** (SYSMIS )  
Dann werden alle zugehörigen alten Werte auf SYSMIS umgesetzt.
    - **Alte Werte kopieren**  
Diese Möglichkeit steht nur beim Umkodieren in *andere* Variablen zur Verfügung und bewirkt für die zugehörigen alten Werte eine unveränderte Übernahme. Dies ist besonders nützlich, wenn die alten Werte mit **Alle anderen Werte** spezifiziert worden sind.
  - Sie können beliebig viele Ersetzungsvorschriften festlegen. SPSS bringt diese automatisch in eine sinnvolle Ordnung.
  - Wenn beim Umkodieren in andere Variablen eine *neue* Variable entsteht, so wird diese zunächst initialisiert, d.h. für alle Fälle wird in der neuen Spalte der Arbeitsdatei der Wert System-Missing eingetragen (vgl. Abschnitt 6.1.3). Durch die *erste zutreffende* Übersetzungsregel wird bei einem Fall der Initialisierungswert durch den zugehörigen neuen Wert überschrieben. Wird der alte Wert eines Falles in keiner Übersetzungsregel angesprochen, dann bleibt bei der neuen Variablen der Initialisierungswert System-Missing stehen! Dies würde in obigem Beispiel etwa einem 1980 geborenen Untersuchungsteilnehmer passieren.
  - **Benutzerdefinierte MD-Indikatoren werden wie gültige Werte behandelt!**  
Ist z.B. für eine Variable der Wert 99 als benutzerdefinierter MD-Indikator deklariert, und wird die 99 rekodiert zur 98, dann **bleibt** die 99 ein MD-Indikator der Variablen, und die 98 wird **nicht** zum MD-Indikator. Eventuell muss also nach der Rekodierung die Variablen-deklaration angepasst werden.
  - Jeder Fall wird **nur einmal** umkodiert, und zwar gemäß der **ersten zutreffenden** Ersetzungsregel (bei Anordnung von oben nach unten).

### 6.2.3 Übungen

- 1) In den beiden folgenden Dialogboxen, die wir allerdings in unserem Projekt *nicht* ausführen wollen, wird jeweils eine Umkodierung der Fachbereichs-Variablen (FB) in eine andere (neue) Variable spezifiziert. Hätten die beiden Dialogboxen denselben Effekt?



2) Bei unserem LOT-Fragebogen wurden die Fragen 3, 4, 5, und 12 aus messtechnischen Gründen umgepolt. Indem eine optimistische Antwort abwechselnd durch Zustimmung oder Ablehnung zum Ausdruck kommt, wird vermieden, dass systematische Ja- oder Neinsager einen extremen Optimismuswert erhalten. Bevor wir einen Mittelwert aus den LOT-Fragen als Optimismus-Schätzwert errechnen können, müssen die negativ gepolten Variablen folgendermaßen umkodiert werden:

5	→	1
4	→	2
2	→	4
1	→	5

Arbeiten Sie mit der **Umkodieren**-Dialogbox, aber quittieren Sie Ihre Eintragungen nicht mit **OK**, sondern mit **Einfügen**, damit das zugehörige RECODE-Kommando in das Syntaxfenster eingetragen wird, in dem wir gerade unser Transformationsprogramm aufbauen.

Machen Sie sich klar, warum die Abbildungsvorschrift „3 → 3“ beim Umkodieren **In dieselben Variablen** überflüssig ist, beim Umkodieren in andere (neue) Variablen aber unbedingt erforderlich wäre.

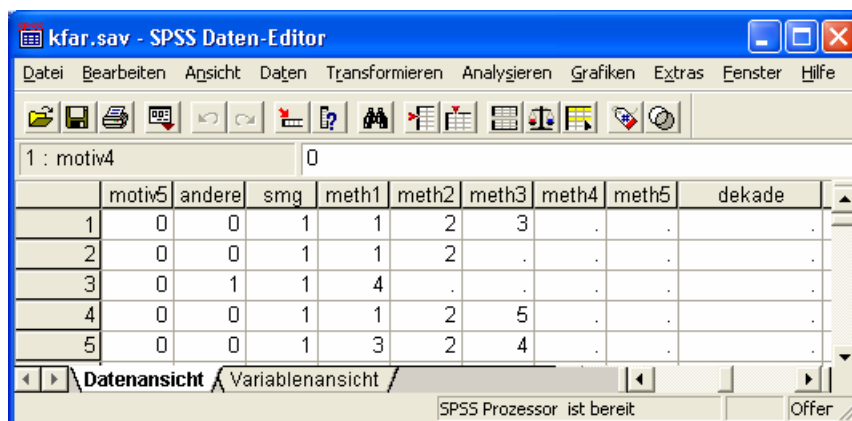
### 6.3 Zur Rolle des EXECUTE-Kommandos

Wenn Sie eine **Umkodieren**-Dialogbox mit **OK** quittieren, dann führt SPSS per Voreinstellung die angeforderte Rekodierung sofort in der Arbeitsdatei aus. Obwohl dieses Verhalten sehr sinnvoll und natürlich erscheint, gibt es eine erwägenswerte Alternative. Zum Rekodieren muss SPSS nämlich die Arbeitsdatei vollständig durchlaufen, was bei einer großen Stichprobe durchaus einige Zeit in Anspruch nehmen kann. Bei einer nächsten und übernächsten Transformationsanweisung (z.B. Rekodierung oder Neuberechnung) ist jeweils ein weiterer Durchlauf fällig. Dabei könnte SPSS zeitsparend *alle* Transformationen in einer *einzig* Datenpassage erledigen. Diese könnte so lange aufgeschoben werden, bis durch die Anforderung einer Statistikprozedur das Durchhackern der Daten unvermeidlich wird. Genau in dem zuletzt beschrieb-

nen, ökonomischen Sinn funktionieren seit jeher die SPSS-Transformationskommandos: Sie werden erst bei der nächsten Prozedur ausgeführt. Allerdings kann dieses zeitoptimierte Verhalten SPSS-Neulinge verwirren. Daher setzt SPSS für Windows hinter jedes per Dialogbox implizit (bei Quittieren mit **OK**) oder explizit (bei Quittieren mit **Einfügen**) produzierte Transformationskommando ein EXECUTE-Kommando, welches die *sofortige Ausführung* aller noch offenen Transformationen erzwingt. Wenn wir z.B. eine **Umkodieren**-Dialogbox mit **OK** quittieren, verarbeitet der SPSS-Prozessor im Hintergrund ein RECODE- und ein EXECUTE-Kommando. Das erste bewirkt nur eine Arbeitsvorbereitung, das zweite erzwingt die Ausführung der vorbereiteten Arbeit. Quittieren wir dieselbe Dialogbox mit **Einfügen**, erscheinen die beiden Kommandos im Syntaxfenster (siehe oben).<sup>1</sup>

Bei der in diesem Manuskript vorgestellten Arbeitsweise sind die von SPSS produzierten EXECUTE-Kommandos in der Regel überflüssig. Aufgrund der heute verfügbaren Rechenleistung lohnt es sich allerdings nur bei einer sehr großen Arbeitsdatei, diese Kommandos aus einem automatisch produzierten Programm zu entfernen.

Beim Arbeiten mit dem Syntaxfenster kann es zu dem folgenden, recht frustrierenden Erlebnis kommen: Sie lassen wohlgeformte Transformationskommandos ausführen, doch im Datenfenster stellt sich nur ein partieller Erfolg ein. Zwar erscheinen die neu anzulegenden Variablen, doch alle Fälle haben den Wert SYSMIS, z.B.:



	motiv5	andere	smg	meth1	meth2	meth3	meth4	meth5	dekade
1	0	0	1	1	2	3	.	.	.
2	0	0	1	1	2	.	.	.	.
3	0	1	1	4	.	.	.	.	.
4	0	0	1	1	2	5	.	.	.
5	0	0	1	3	2	4	.	.	.

Die Ursache ist dann meist: Sie haben nach den Transformationskommandos noch kein Prozedur- bzw. EXECUTE-Kommando ausführen lassen, so dass SPSS zwar die neue Variablen initialisiert, aber noch keine Werte ermittelt hat.

In dieser Situation wird in der Statuszeile angezeigt, dass **Offene Transformationen** zur Bearbeitung anstehen. Sie können deren Ausführung erzwingen, indem Sie im Syntaxfenster ein EXECUTE-Kommando abschicken oder folgenden Menübefehl wählen:

### **Transformieren > Offene Transformationen ausführen**

Es soll nicht verschwiegen werden, dass hier für SPSS-Neulinge Schwierigkeiten auftauchen, die bei rein dialogbox-orientierter Arbeitsweise nicht entstehen können.

<sup>1</sup> Man kann nach

#### **Bearbeiten > Optionen > Daten**

im Rahmen **Optionen für Transformieren und Zusammenfügen** mit der Option **Werte vor Verwendung berechnen** die voreingestellte EXECUTE-Inflation abstellen. Dann zeigt SPSS das oben beschriebene zeitoptimierte Verhalten, führt also z.B. nach dem Quittieren einer **Umkodieren**-Dialogbox mit **OK** das zugrunde liegende RECODE-Kommando zunächst noch nicht aus, sondern reiht es in die Warteschlange der offenen Transformationen ein. Diese werden vom SPSS-Prozessor erst dann ausgeführt, wenn er ein Prozedur- oder ein EXECUTE-Kommando erhält.

Für angehende SPSS-Profis möchte ich noch erwähnen, dass EXECUTE-Kommandos *innerhalb eines Blocks von Transformationsanweisungen* durchaus bedeutsam sein können. In dem folgenden (manuell erstellten) Beispiel wird mit Hilfe des Transformationskommandos SELECT IF jeder zweite Fall aus der Arbeitsdatei entfernt:

```
compute nr = $casenum.
execute.
select if (mod(nr,2) = 1).
execute.
```

Lässt man das erste EXECUTE weg, entfernt das Programm *alle* Fälle mit Ausnahme des ersten.

## 6.4 Berechnung von Variablen nach mathematischen Formeln

In der Dialogbox **Variable Berechnen** bzw. im äquivalenten COMPUTE-Kommando wird ein numerischer Ausdruck (z.B. GROESSE - 100) definiert und einer Ergebnisvariablen zugewiesen. Dabei kann man eine *neue* Variable erzeugen oder eine vorhandene verändern.

### 6.4.1 Beispiel

Sie sollen später anhand unserer Stichprobe untersuchen, ob die Trierer Studierenden im Mittel wenigstens das folgende Idealgewicht auf die Waage bringen (Nullhypothese)

$$\text{Gewicht (in kg)} = \overset{!}{\text{Größe(in cm)}} - 100$$

oder ob sie relativ zu dieser Formel zu leicht sind (Alternativhypothese). Zur Prüfung dieser Frage mit einem t-Test für gepaarte Stichproben muss die Arbeitsdatei um eine neue Variable, z.B. IDGEW genannt, erweitert werden, deren Werte nach obiger Formel aus der Körpergröße zu berechnen sind. Anschließend enthält die (Fälle × Variablen)-Datenmatrix in der Arbeitsdatei u.a. die beiden folgenden Variablen:

GROESSE	IDGEW
163	63
158	58
174	74
182	82
.	.
.	.
.	.
176	76
176	76
170	70
169	69

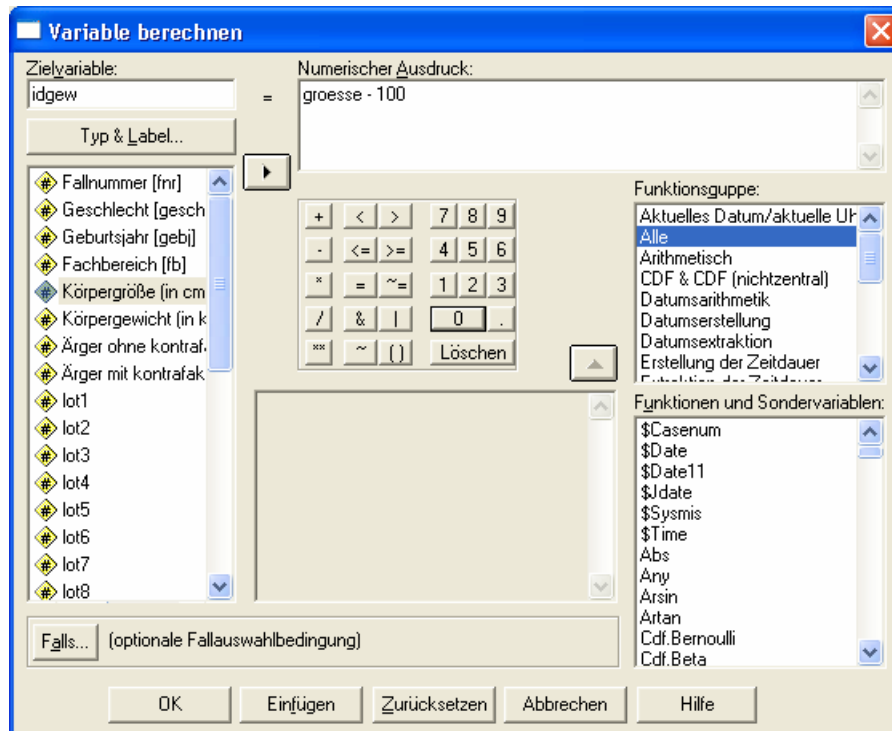
Starten Sie zum Definieren der neuen Variablen die Dialogbox **Variable berechnen** mit:

### Transformieren > Berechnen...

Tragen Sie zunächst im Feld **Zielvariable** den Namen für die neu in die Arbeitsdatei aufzunehmende Variable ein (IDGEW), und schreiben Sie dann in das Feld **Numerischer Ausdruck** die Definitionsvorschrift (GROESSE - 100), wobei einige Schreibhilfen zur Verfügung stehen:

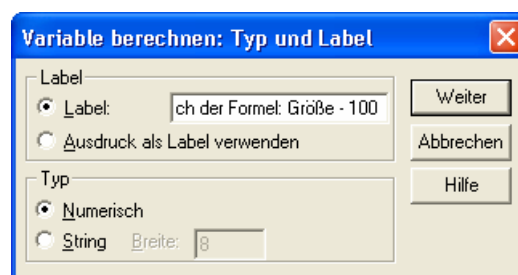
- Der Variablenname kann aus einer Liste per Transportschalter oder Doppelklick übernommen werden.
- Mit Hilfe einer virtuellen Tastatur können Sie das Minuszeichen und die Zahl 100 auch per Maus eingeben.

Anschließend sollte Ihre Dialogbox ungefähr so aussehen:



Die Dialogbox bietet über unsere momentanen Bedürfnisse hinausgehend auch die in SPSS verfügbaren Funktionen (siehe unten) und Sondervariablen (z.B. **\$Casenum** für die fortlaufende Fallnummer in der Arbeitsdatei) in Gruppen geordnet zum Transport in das Feld **Numerischer Ausdruck** an, so dass man bei der Verwendung von Funktionen nicht nachschlagen muss und Tippfehler vermeidet.

Rufen Sie nun mit dem gleichnamigen Schalter die Subdialogbox **Typ und Label** auf, und tragen Sie dort für IDGEW das Etikett *Idealgewicht nach der Formel: Größe - 100* ein:



Quittieren Sie die Subdialogbox mit **Weiter** und die Hauptdialogbox mit **Einfügen**. Daraufhin erhalten Sie im Syntaxfenster ein COMPUTE- und ein VARIABLE LABELS - Kommando:

```
COMPUTE idgew = groesse - 100 .
VARIABLE LABELS idgew 'Idealgewicht nach der Formel: Größe - 100' .
EXECUTE .
```



## 6.4.2 Technische Details

### 6.4.2.1 Numerischer Ausdruck

Im Texteingabefeld **Numerischer Ausdruck** der Dialogbox **Variable berechnen** sind wir trotz der SPSS-Scheibhilfen im Wesentlichen wieder in das „Syntaxzeitalter“ zurückgeworfen: Auf der weißen Fläche ist ein sprachlicher Ausdruck nach gewissen Syntaxregeln zu formulieren. Zum Glück sind uns aber numerische Ausdrücke aus der Schule wohlbekannt<sup>1</sup>.

Konkret darf ein numerischer Ausdruck im Sinne von SPSS folgende Bestandteile enthalten:

- Bereits definierte Variablen
- Zahlen
- arithmetische Operatoren:
  - Addition (+)
  - Subtraktion (-)
  - Multiplikation (\*)
  - Division (/)
  - Potenzfunktion (\*\*)
- Klammern
- Funktionen

#### 6.4.2.1.1 Numerische Funktionen

In numerischen Ausdrücken können Sie zahlreiche Funktionen verwenden, die numerische Variablen oder Zahlen als Argumente (in den folgenden Syntaxdarstellungen vertreten durch den Platzhalter *arg*) verarbeiten.<sup>2</sup> Diese Funktionen lassen sich in mehrere Gruppen einteilen, aus denen jeweils einige wichtige Vertreter genannt werden sollen:

- **Arithmetische Funktionen, z.B.:**

- |                                    |   |
|------------------------------------|---|
| - ABS( <i>arg</i> )                | Absoluter Wert  |
| - EXP( <i>arg</i> )                | Exponentialfunktion   |
| - LG10( <i>arg</i> )               | Dekadischer Logarithmus   |
| - LN( <i>arg</i> )                 | Natürlicher Logarithmus   |
| - MOD( <i>arg1</i> , <i>arg2</i> ) | Rest aus der Division von <i>arg1</i> durch <i>arg2</i> , z.B.<br>mod(1.3, 1) = 0.3 |
| - RND( <i>arg</i> )                | Auf eine ganze Zahl gerundeter Wert   |
| - SQRT( <i>arg</i> )               | Quadratwurzel   |

Beispiel: `compute logi = exp(3+1.2*x) / (1+exp(3+1.2*x)) .`  
Hier wird eine spezielle logistische Funktion der Variablen X definiert.

- **Statistische Funktionen, z.B.:**

- |  |  |
|--|--|
| - MEAN[.n]( <i>arg1</i> , <i>arg2</i> [, ...]) | Arithmetisches Mittel<br>Voreinstellung für <i>n</i> : 1 |
| - MAX[.n]( <i>arg1</i> , <i>arg2</i> [, ...])  | Maximum<br>Voreinstellung für <i>n</i> : 1               |
| - MIN[.n]( <i>arg1</i> , <i>arg2</i> [, ...])  | Minimum<br>Voreinstellung für <i>n</i> : 1               |

<sup>1</sup> Zwar gibt es gewisse Unterschiede zwischen mathematischen *Gleichungen* (z.B.  $y = a + b \cdot x$ ) und EDV-sprachlichen *Zuweisungen* (z.B. `compute x = x + 2.`), doch sind die Regeln für die numerischen Ausdrücke auf den *rechten* Seiten weitgehend identisch.

<sup>2</sup> SPSS kennt auch zahlreiche Funktionen für String- und Datums-Variablen, die aber aus Zeitgründen in diesem Kurs nicht behandelt werden. Informieren Sie sich bei Bedarf im Hilfesystem, z.B. über eine Indexsuche nach dem Stichwort *Funktionen*.

- SD[.n](arg1,arg2[, ...])      Standardabweichung  
Voreinstellung für n: 2
- SUM[.n](arg1,arg2[, ...])      Summe  
Voreinstellung für n: 1

- Regeln:
- Die eckigen Klammern schließen optionale Angaben ein.
  - Der Funktionsparameter *n* hat folgende Bedeutung: Wenn bei einem Fall mindestens *n* valide Argumente vorliegen, wird der Funktionswert berechnet. Ansonsten wird dem Fall der Wert SYSMIS zugewiesen. Wenn Sie mit der sehr liberalen Voreinstellung für *n* nicht einverstanden sind, können Sie einen alternativen Wert festlegen.
  - Mit „[, ...]“ wird zum Ausdruck gebracht, dass die Liste der Argumente optional beliebig verlängert werden darf.
  - Sie können eine Serie von Variablen, *die in der Arbeitsdatei hintereinander stehen*, bequem auf folgende Weise in einer Argumentenliste angeben:

*erste TO letzte*

Es kommt nicht auf die alphanumerische Ordnung der Variablennamen an, sondern auf die Reihenfolge der Variablen in der Arbeitsdatei.

- Beispiel:      `compute mfrei = mean.45(sport to angeln).`  
 Wenn für einen Fall bei den Variablen SPORT bis ANGELN, die in der Arbeitsdatei hintereinander stehen, mindestens 45 valide Argumente vorliegen, wird deren Mittelwert der neuen Variablen zugewiesen, ansonsten wird der MD-Indikator System-Missing zugewiesen.

Beachten Sie den wesentlichen Unterschied zwischen den gerade beschriebenen statistischen *Funktionen* und den Statistik*prozeduren*, mit denen wir z.B. die Verteilungsanalysen durchgeführt haben:

- Wenn wir in der Dialogbox **Häufigkeiten** (erreichbar über **Analysieren > Deskriptive Statistiken > Häufigkeiten**) z.B. den Mittelwert der Variablen GEWICHT anfordern, werden die (validen) Gewichtsangaben aller Fälle in der Stichprobe gemittelt. Es werden also die Ausprägungen *einer Variablen* über *alle Fälle* gemittelt. SPSS arbeitet sich *senkrecht* durch eine komplette Variable bzw. Spalte der Arbeitsdatei. Es resultiert ein einziger Stichprobenkennwert, welcher im Ausgabefenster erscheint.
- Mit der statistischen Funktion MEAN können wir für *jede einzelne Person* z.B. den Mittelwert über *mehrere LOT-Variablen* berechnen lassen. SPSS geht *waagerecht* vor, wobei dasselbe Verfahren *auf jeden Fall, d.h. auf jede Zeile* der Datenmatrix angewendet wird. Die statistische Funktion MEAN erzeugt (oder modifiziert) eine Variable, d.h. eine komplette Spalte im Datenfenster, in die für jeden Fall sein eigenes Berechnungsergebnis eingetragen wird.

• **Funktionen für fehlende Werte, z.B.:**

- NMISS(arg1[, ...])      Anzahl fehlender Werte bei den aufgelisteten Variablen
- VALUE(arg)      Es wird der Wert der Variablen *arg* geliefert, wobei *benutzerdefinierte* MD-Deklarationen ignoriert werden.

- Regeln:
- Mit „[, ...]“ wird zum Ausdruck gebracht, dass die Liste der zu untersuchenden Variablen optional beliebig verlängert werden darf.
  - Mit dem Schlüsselwort TO können bequem Serien von Variablen angegeben werden (siehe 1. Beispiel und obige Erläuterungen zu den statistischen Funktionen).

Beispiele: - `compute nmfrei = nmiss(sport to angeln) .`  
 Der numerische Ausdruck liefert die Anzahl der fehlenden Werte (SYMIS oder benutzerdefiniert) bei den Variablen SPORT bis ANGELN, die in der Arbeitsdatei hintereinander stehen.

- `compute vala = value(a) .`  
 Diese Funktion liefert auch dann den Wert der Variablen A, wenn es sich um einen benutzerdefinierten MD-Indikator handelt.

- **Pseudozufallszahlengeneratoren, z.B.:**

- `NORMAL(arg)` Die Funktion liefert normalverteilte Zufallszahlen mit Mittelwert 0 und Standardabweichung *arg*.
- `UNIFORM(arg)` Die Funktion liefert gleichverteilte Zufallszahlen im Intervall von 0 bis *arg*.

Beispiel: `COMPUTE av = NORMAL(1) .`  
`EXECUTE .`  
`T-TEST`  
`GROUPS=geschl(1 2)`  
`/MISSING=ANALYSIS`  
`/VARIABLES=av`  
`/CRITERIA=CIN(.95) .`

Die Kommandos aus diesem Beispiel wurden mit Hilfe von Dialogboxen erzeugt (Schalter **Einfügen**). Im COMPUTE-Kommando wird die normalverteilte Zufallsvariable AV erstellt. Es ist klar, dass Frauen und Männer denselben Erwartungswert (Populationsmittelwert) bei AV haben. Damit können wir ausprobieren, wie sich der t-Test für unabhängige Stichproben bei Gültigkeit der Nullhypothese identischer Erwartungswerte verhält. Die Dialogbox zu diesem t-Test erhält man mit **Analysieren > Mittelwerte vergleichen > t-Test bei unabhängigen Stichproben**.

Wenn Ihnen die Erläuterungen zu diesem Beispiel „spanisch“ vorkommen, hilft Ihnen vielleicht der Abschnitt 7.1 weiter, wo einige Grundprinzipien der Inferenzstatistik erläutert werden.

Hinweis: Bei NORMAL und UNIFORM wird ein Pseudozufallszahlengenerator verwendet, der per Voreinstellung mit dem festen Wert 2000000 startet und damit stets dieselben Zahlen liefert. Ein alternativer Startwert, der andere Zufallszahlen zur Folge hat, kann gewählt werden:

- mit dem Menübefehl:  
**Transformieren > Zufallszahlengeneratoren**
- oder mit dem SPSS-Kommando:  
`SET SEED=n.`

#### 6.4.2.1.2 Regeln für die Bildung numerischer Ausdrücke

Auch bei Verwendung der Dialogbox **Variable berechnen** müssen wir die numerischen Ausdrücke im Wesentlichen selbst formulieren. Dabei sind folgende Regeln zu beachten:

- Die **Auswertungsreihenfolge** hängt von der Priorität der Operatoren ab. Es gilt folgende Rangordnung:
  - Priorität 1: Funktionen
  - Priorität 2: Potenzfunktion (\*\*)
  - Priorität 3: Multiplikation (\*), Division (/) und Vorzeichen-Minus (z.B.: "-b")
  - Priorität 4: Addition (+), Subtraktion (-)

Bei gleicher Priorität erfolgt die Auswertung von links nach rechts. Eine alternative Auswertungsreihenfolge kann durch Klammern erzwungen werden: Klammersausdrücke werden zuerst ausgewertet. Bei geschachtelten Klammern erfolgt die Auswertung von innen nach außen.

- Bei Funktionen mit mehreren Argumenten müssen die einzelnen Argumente **durch jeweils genau ein Komma** (optional ergänzt durch Leerzeichen) getrennt werden.

Beispiel: `compute mabc = mean(a, b, c) .`

- Obwohl SPSS im Daten- und im Ausgabefenster das ländertypische Dezimaltrennzeichen benutzt, bei uns also das Komma, müssen in numerischen Ausdrücken gebrochene Zahlen generell mit Dezimalpunkt geschrieben werden.

Richtig: `2.75`

Falsch: `2,75`

Dies gilt sowohl für das Feld **Numerischer Ausdruck** der Dialogbox **Variable berechnen** als auch für das COMPUTE-Kommando in einem Syntaxfenster.

Es kann also durchaus passieren, dass Sie ein und dieselbe gebrochene Zahl im Datenfenster (als Wert eines Falles für eine bestimmte Variable) mit Dezimalkomma und in der Dialogbox **Variable berechnen** (z.B. als Konstante in einer Berechnungsanweisung) mit Dezimalpunkt schreiben müssen.

- In der Regel sind numerische Ausdrücke als Argumente von Funktionen zugelassen.

Beispiel: `compute albmax = max(a, ln(b)) .`

Das zweite Argument der Funktion MAX ist der numerische Ausdruck `ln(b)`.

#### 6.4.2.2 Sonstige Hinweise

##### SYSMIS als Ergebnis eines numerischen Ausdrucks

Durch eine Berechnungsanweisung wird der Wert des numerischen Ausdrucks auch dann der Zielvariablen zugewiesen, wenn dieser Wert gleich SYSMIS ist (z.B. bei fehlenden Argumenten). Dieses Vorgehen ist kompatibel mit dem in Abschnitt 6.1.3 beschriebenen Initialisierungsprinzip für neue numerische Variablen. Ist die Zielvariable bereits vorhanden, bleibt bei missglückter Berechnung des numerischen Ausdrucks keinesfalls der alte Wert bestehen, sondern es wird sinnvollerweise SYSMIS zugewiesen.

##### Rechnen mit fehlenden Werten

Wenn bei einem Fall eine Variable aus dem numerischen Ausdruck keinen validen Wert hat, dann erhält die Ergebnisvariable den Wert SYSMIS. Ausnahmen sind die folgenden SPSS-eigenen Regeln für das „Rechnen“ mit fehlenden Werten:

- $0 * \text{unbekannt} = 0$

Diese Regel ist schlau, denn für beliebige reelle Zahlen  $x$  gilt:

$$0 \cdot x = 0$$

- $0 / \text{unbekannt} = 0$

Diese Regel ist kritisierbar, denn:

$$\frac{0}{x} = \begin{cases} 0 & \text{für } x \neq 0 \\ \text{undefiniert} & \text{für } x = 0 \end{cases}$$

### 6.4.3 Übungen

1) Welche Werte haben die folgenden numerischen Ausdrücke?

$$(3 + 4) / 2$$

$$3 + 4 / 2$$

$$(3**2 / 2) + 4$$

$$3**2 / 2 + 4$$

2) Erstellen Sie im KFA-Projekt die Variablen, auf die sich unsere zentralen Hypothesen beziehen (vgl. Abschnitt 1.3):

- Berechnen Sie die Variable LOT als arithmetisches Mittel der (nötigenfalls rekodierten!) LOT-Variablen 1, 3, 4, 5, 8, 9, 11 und 12. Die restlichen Fragen dienen nicht zur Messung von Optimismus, sondern sollen verhindern, dass der Zweck des Fragebogens deutlich wird. Dies könnte das Antwortverhalten verzerren. Tolerieren Sie bei der Berechnung des Mittelwertes bis zu zwei fehlende Werte.
- Berechnen Sie die Variable AERGAM als arithmetisches Mittel der beiden Ärgervariablen und die Variable AERGZ als Ärgerzuwachs auf Grund der kontrafaktischen Alternative.

AERGAM benötigen wir zum Testen der differentialpsychologischen Hypothese. Beim geplanten Test der allgemeinspsychologischen Hypothese wird letztlich mit einem Einstichproben-t-Test geprüft, ob der Erwartungswert (Populationsmittelwert) der Variablen AERGZ signifikant größer als Null ist. Man kann den Test zwar bequem mit der SPSS-Prozedur zum t-Test für gepaarte Stichproben durchführen, ohne die Variable AERGZ explizit berechnen zu müssen, doch bietet diese Prozedur keine Möglichkeit, die Normalverteilungsvoraussetzung des Tests (vgl. Abschnitt 7.1) zu prüfen. Daher berechnen wir AERGZ explizit und prüfen die Verteilungsvoraussetzung mit der Prozedur zur explorativen Datenanalyse (siehe Abschnitt 7.3).

Rufen Sie jeweils mit dem Menübefehl:

#### **Transformieren > Berechnen...**

die zuständige Dialogbox auf, quittieren Sie aber Ihre Eintragungen nicht mit **OK**, sondern mit **Einfügen**, damit die zugehörigen COMPUTE-Kommandos in das Syntaxfenster eingetragen werden, in dem gerade das Transformationsprogramm entsteht.

Weil SPSS eine Folge von mehreren Kommandos stets in der natürlichen Reihenfolge abarbeitet, wird beim späteren Ablauf unseres Transformationsprogramms z.B. die für einige Items angeordnete Rekodierung bereits erledigt sein, wenn das COMPUTE-Kommando zur LOT-Berechnung ausgeführt wird.

### 6.5 Bedingte Datentransformation

Häufig ist es erforderlich, eine Datenmodifikation auf diejenigen Fälle zu beschränken, die eine bestimmte Bedingung erfüllen. Wir benötigen z.B. im KFA-Projekt eine solche Möglichkeit, um bei den Motivations- und Methodenvariablen das bisher vertagte Problem der fehlenden Werte adäquat behandeln zu können (siehe Abschnitt 1.4.3.2).

Manchmal ist es angebracht, für mehrere disjunkte Teilmengen der Gesamtstichprobe jeweils spezifische Modifikationen durchzuführen (Fallunterscheidung). Z.B. könnte man im Rahmen einer Untersuchung zum Essverhalten bei der Berechnung der neuen Variablen Idealgewicht aus der bereits vorhandenen Variablen Körpergröße bei Frauen und Männern unterschiedliche Formeln anwenden.

In den SPSS - Transformations-Dialogboxen erreichen Sie über den Schalter **Falls** eine Subdialogbox zur Definition einer Bedingung, unter der die Transformation ausgeführt werden soll. Sie

können z.B. eine bedingte Umkodierung (vgl. Abschnitt 6.2), Berechnung (vgl. Abschnitt 6.4) oder Wertauszählung (vgl. Abschnitt 6.6) vornehmen.

Wenn unter ein und derselben Bedingung gleich *mehrere* Transformationen vorgenommen werden sollen, muss diese Bedingung in allen benötigten Transformations-Dialogboxen, wiederholt werden. Ähnlich umständlich ist die Realisation von Fallunterscheidungen mit Hilfe der Transformations-Dialogboxen. Für solche Aufgaben bietet die SPSS-Kommandosprache mit der DO IF - ELSE IF - END IF - Kontrollstruktur bessere Lösungen. Diese lassen sich jedoch nicht komplett mit Dialogboxen generieren, so dass sie in diesem Kurs aus Zeitgründen nicht behandelt werden.

### 6.5.1 Beispiel

In diesem Abschnitt soll endlich das MD-Problem bei den Motivationsvariablen gelöst werden. Wir haben bei den Variablen MOTIV1 bis MOTIV5 und ANDERE systematisch die angekreuzten Kästchen mit 1 und die leeren Kästchen mit 0 kodiert, um während der Erfassung möglichst wenige zeitraubende und fehleranfällige Entscheidungen treffen zu müssen. Ein Fall mit Nullen bei MOTIV1 bis MOTIV5 *und* ANDERE hat aber offenbar den Fragebogenteil 4a komplett ausgelassen. Daher müssen für genau diese Fälle die Nullen bei den Variablen MOTIV1 bis MOTIV5 und ANDERE in SYSMIS umkodiert werden. Gehen Sie folgendermaßen vor:

- Wählen Sie den Menübefehl:

#### **Transformieren > Umkodieren > in dieselben Variablen...**

- Transportieren Sie die Variablennamen MOTIV1 bis MOTIV5 und ANDERE in die Teilnehmerliste der **Umkodieren**-Dialogbox.
- Legen Sie in der Subdialogbox **Alte und neue Werte** die benötigte Abbildungsvorschrift fest.
- Öffnen Sie die **Falls**-Subdialogbox, markieren Sie die Option **Fall einschließen, wenn Bedingung erfüllt ist**, und tragen Sie in das darunter liegende Textfeld eine geeignete Bedingung ein, z.B.:



Aufgrund unserer Datenüberprüfung können wir uns darauf verlassen, dass bei den Variablen MOTIV1 bis MOTIV5 und KEINE ausschließlich die Werte Null und Eins vorliegen. Daher ist die Summe dieser Variablen genau dann gleich Null, wenn jede einzelne Variable gleich Null ist.

Die obige Eintragung im Bedingungsfeld kann „semiautomatisch“ z.B. folgendermaßen erzeugt werden:

- Markieren Sie in der Funktionenliste **SUM(numausdr,numausdr,...)** und klicken Sie auf den zugehörigen Transportschalter.
- Transportieren Sie aus der Variablenliste MOTIV1 in das Bedingungsfeld.
- Schreiben Sie den Rest der Einfachheit halber per Hand.
- Machen Sie **Weiter**, und quittieren Sie die Hauptdialogbox mit **Einfügen**.

Daraufhin wird Ihr Transformationsprogramm um die folgende Sequenz erweitert:

```
DO IF (SUM(motiv1 to andere) = 0) .
RECODE
  motiv1 motiv2 motiv3 motiv4 motiv5 andere (0=SYSMIS) .
END IF .
EXECUTE .
```

Wenn Sie diese Kommandos ausführen lassen, gleichgültig ob direkt per **OK** in der **Umkodieren**-Dialogbox oder indirekt via Syntaxfenster, passiert bei jedem einzelnen Fall in der Stichprobe folgendes:

- SPSS prüft die Bedingung, die wir auch als **logischen Ausdruck** bezeichnen wollen.
- Ist bei einem Fall die Bedingung erfüllt, dann wird umkodiert, anderenfalls passiert nichts.

Weil die Variablen MOTIV1 bis MOTIV5 und ANDERE vor der Rekodierung garantiert nur Nullen oder Einsen als Werte aufweisen, hat unser logischer Ausdruck übrigens die Eigenschaft, in jedem Fall entweder wahr oder falsch zu sein. Das erscheint nach dem aussagenlogischen Axiom vom ausgeschlossenen Dritten als selbstverständlich, ist es aber in der empirischen Forschung z.B. wegen des nahezu allgegenwärtigen Problems fehlender Werte keineswegs. Für die Fälle in unserer Stichprobe kann z.B. der logische Ausdruck „GESCHL = 1“ folgende Wahrheitswerte annehmen:

- wahr            ⇔    Der GESCHL-Wert ist gleich Eins.
- falsch         ⇔    Der GESCHL-Wert ist eine von Eins verschiedene Zahl.
- unbestimmt ⇔    Der GESCHL-Wert fehlt (ist gleich SYSMIS).

Komplexere logische Ausdrücke (z.B. „LN(ML)/ANZ > 1“) können auch wegen undefinierter Funktionswerte unbestimmt sein (bei  $ML \leq 0$  oder  $ANZ = 0$ ).

Wenn Sie eine bedingte Transformationsanweisung verwenden, sollten Sie beachten, wie SPSS auf bestimmte und unbestimmte logische Ausdrücke reagiert:

- Ist der logische Ausdruck **wahr**, dann wird die Transformation ausgeführt.  
Im Fall einer bedingten Berechnung wird der Ergebnisvariablen also der Wert des numerischen Ausdrucks zugewiesen. Die Zuweisung erfolgt auch dann, wenn der numerische Ausdruck den Wert SYSMIS hat.
- Ist der logische Ausdruck **falsch oder unbestimmt**, dann passiert **nichts**, d.h.:
  - Eine bereits vorhandene Ergebnisvariable behält für den betroffenen Fall ihren bisherigen Wert.
  - Bei einer neu definierten Variablen behält der betroffene Fall den Initialisierungswert SYSMIS.

### 6.5.2 Bedingungen formulieren

Der in obigem Beispiel aufgetretene logische Ausdruck war recht einfach aufgebaut, weil er nur aus einem einzigen Vergleich bestand. Obwohl Ihnen auch komplexere Exemplare (z.B. aus der Schule) wohlvertraut sein dürften, soll der Begriff *logischer Ausdruck* zur Klärung einiger Spe-

zialprobleme etwas genauer beschrieben werden. Zunächst wird der einfachere Begriff *Vergleich* eingeführt.

### 6.5.2.1 Vergleich

Ein Vergleich besteht aus zwei numerischen Ausdrücken und einem Vergleichsoperator:

*numerischer\_ausdruck vergleichs-operator numerischer\_ausdruck*

Die bekannten Vergleichsoperatoren können in SPSS alternativ durch EDV-Varianten der mathematischen Symbole oder durch Schlüsselwörter dargestellt werden:

Symbol	Schlüsselwort	Bedeutung
=	EQ	gleich
<>	NE	ungleich
<	LT	kleiner als
<=	LE	kleiner oder gleich
>	GT	größer als
>=	GE	größer oder gleich

Beispiel:                   beruf > 4

### 6.5.2.2 Logischer Ausdruck

Aus dem einfachen Begriff *Vergleich* wird nun durch eine rekursive Definition der komplexere Begriff *logischer Ausdruck* konstruiert:

- i) Jeder Vergleich ist ein logischer Ausdruck.
- ii) Durch Anwendung des logischen Operators **NOT** auf einen logischen Ausdruck oder durch Anwendung der logischen Operatoren **AND** bzw. **OR** auf zwei logische Ausdrücke entsteht ein neuer logischer Ausdruck:

NOT *logischer\_ausdruck*

*logischer\_ausdruck\_1* AND *logischer\_ausdruck\_2*

*logischer\_ausdruck\_1* OR *logischer\_ausdruck\_2*

Den Wahrheitswert eines zusammengesetzten logischen Ausdrucks erhält man aus den Wahrheitswerten der Argumente nach den Regeln für logische Operatoren, die in den so genannten Wahrheitstafeln festgelegt sind (siehe unten).

Es lassen sich sukzessiv beliebig komplexe logische Ausdrücke aufbauen, die für einen konkreten Fall immer die Wahrheitswerte *wahr*, *falsch* oder *unbestimmt* haben können.

Beispiel:                   (lie1 = 0) and (lie2 = 0)

Das Problem unbestimmter Wahrheitswerte in logischen Ausdrücken löst SPSS analog zu den Regeln für das Rechnen mit fehlenden Werten in numerischen Ausdrücken (siehe Abschnitt 6.4.2.2). Die folgenden Wahrheitstafeln sind gegenüber der klassischen Aussagenlogik um den Wahrheitswert *unbestimmt* erweitert (*la1* und *la2* seien logische Ausdrücke):



<i>la1</i>	NOT <i>la1</i>
wahr	falsch
falsch	wahr
unbestimmt	unbestimmt

<i>la1</i>	<i>la2</i>	<i>la1</i> AND <i>la2</i>	<i>la1</i> OR <i>la2</i>
wahr	wahr	wahr	wahr
wahr	falsch	falsch	wahr
wahr	unbestimmt	unbestimmt	wahr
falsch	wahr	falsch	wahr
falsch	falsch	falsch	falsch
falsch	unbestimmt	falsch	unbestimmt
unbestimmt	wahr	unbestimmt	wahr
unbestimmt	falsch	falsch	unbestimmt
unbestimmt	unbestimmt	unbestimmt	unbestimmt

### 6.5.2.3 Regeln für die Auswertung logischer Ausdrücke

Bei der Auswertung von logischen Ausdrücken gelten in SPSS folgende Regeln:

- Die Abarbeitungsreihenfolge hängt von der Priorität der Operatoren ab. Es gilt folgende Rangordnung:
  - Priorität 1: Funktionen
  - Priorität 2: Potenzfunktion (\*\*)
  - Priorität 3: Multiplikation (\*), Division (/),  
Vorzeichen-Minus (z.B. -a)
  - Priorität 4: Addition (+), Subtraktion (-)
  - Priorität 5: Vergleichsoperatoren
  - Priorität 6: NOT
  - Priorität 7: AND
  - Priorität 8: OR
- Bei gleicher Priorität: Abarbeitung von links nach rechts.
- Eine andere Auswertungsreihenfolge kann durch Klammern erzwungen werden.

Beispiel: Das obige Beispiel für einen zusammengesetzten logischen Ausdruck wegen der voreingestellten Abarbeitungsreihenfolge auch einfacher geschrieben werden:

$$lie1 = 0 \text{ and } lie2 = 0$$

### 6.5.3 Übung

Bei den Variablen METH1 bis METH3 haben wir zur Vereinfachung der Erfassung im Kodierplan festgelegt, dass „unbenutzte“ Variablen einfach leer bleiben sollen. Nun wollen wir aber bei Fällen mit regulärem Antwortmuster die SYSMIS - Werte durch Nullen ersetzen. Die Null soll z.B. bei der Variablen METH2 bedeuten: Die Option, einen zweiten Methodenwunsch zu äußern, wurde nicht genutzt.

Die folgende Tabelle, die wir in Abschnitt 1.4.3.2.3 vereinbart haben, legt im Einzelnen fest, was unter den möglichen Bedingungskonstellationen geschehen soll:

		Mindestens eine speziell interessierende Methode angegeben?	
		Ja	Nein
SMG	1	METH1 ... METH3: SYSMIS → 0 Bem.: Korrektes Antwortverhalten. Variablen zu nicht benutzten Optionen (gem. Kodierplan bisher auf SYSMIS) werden auf 0 gesetzt.	SMG: 1 → SYSMIS Bem.: Irreguläres Antwortverhalten. METH1 bis METH3 behalten SYMIS. SMG wird ebenfalls auf SYMIS gesetzt.
	0	SMG: 0 → 1 METH1 ... METH3: SYSMIS → 0 Bem.: Leicht irreguläres Antwortverhalten. Wir sind großzügig und setzen SMG auf 1.	METH1 ... METH3: SYSMIS → 0 Bem.: Korrektes Antwortverhalten. Die Variablen zu allen Optionen (gem. Kodierplan bisher auf SYSMIS) werden auf 0 gesetzt.
	SYSMIS	SMG: SYSMIS → 1 METH1 ... METH3: SYSMIS → 0 Bem.: Leicht irreguläres Antwortverhalten. Wir sind großzügig und setzen SMG auf 1 sowie die Variablen zu nicht benutzten Optionen auf 0.	Bem.: Irreguläres Antwortverhalten. Alle Variablen behalten den Wert SYSMIS.

In den beiden obersten Zeilen jeder Zelle sind die erforderlichen Korrekturen bei SMG bzw. METH1 bis METH3 angegeben.

Erweitern Sie Ihr Programm **kfat.sps** um passende Transformationsanweisungen.

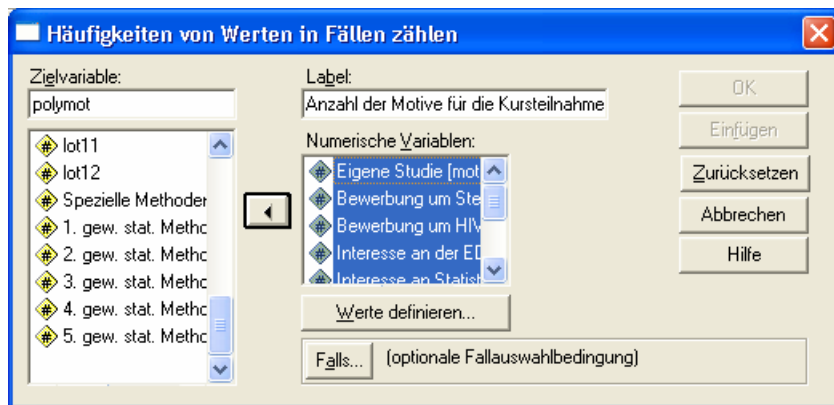
### 6.6 Häufigkeit bestimmter Werte bei einem Fall ermitteln

Mit dem Befehl **Zählen** aus dem Menü **Transformieren** bzw. mit dem zugrunde liegenden COUNT-Kommando kann man eine Variable berechnen lassen, die für jeden Fall festhält, wie oft bestimmte Werte in einer Liste von  $k$  Variablen vorkommen. Das minimale Ergebnis ist Null (keine Variable hat einen der kritischen Werte), und das maximale Ergebnis ist  $k$  (jede Variable hat einen kritischen Wert).

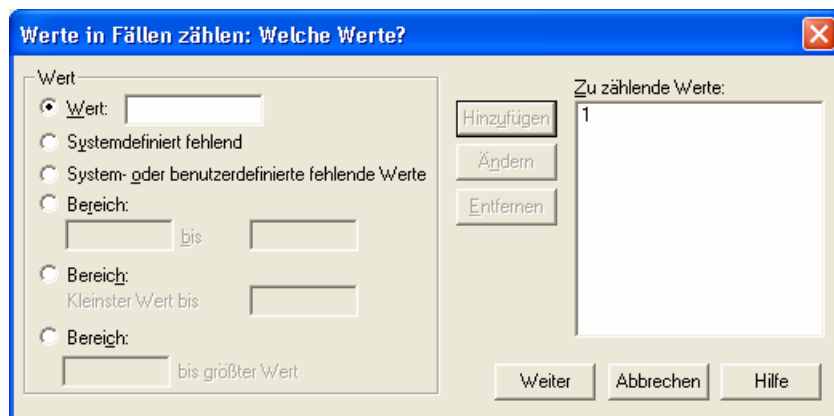
Wir wollen eine neue Variable namens POLYMOT berechnen lassen, die für jede Person festhält, wie viele Motive zur Kursteilnahme sie im Fragebogenteil 4a angegeben hat. Aktivieren Sie die Dialogbox **Häufigkeiten von Werten in Fällen zählen** mit

#### Transformieren > Zählen...

Vergeben Sie für die Zielvariable den Namen POLYMOT sowie das Label *Anzahl der Motive für die Kursteilnahme*, und transportieren Sie die Variablen MOTIV1 bis ANDERE in die Teilnehmerliste. Danach müsste Ihre Dialogbox ungefähr so aussehen:



Wechseln Sie jetzt mit dem Schalter **Werte definieren** in die Subdialogbox **Werte in Fällen zählen: Welche Werte?**, tragen Sie dort den kritischen **Wert** Eins ein, und klicken Sie auf **Hinzufügen**:



Die in dieser Subdialogbox angebotenen sonstigen Möglichkeiten zur Festlegung der Trefferwerte kennen wir übrigens schon aus der Subdialogbox **Umkodieren: Alte und neue Werte** (siehe Abschnitt 6.2).

Da SPSS eine Folge von mehreren Kommandos stets in der natürlichen Reihenfolge abarbeitet, wird beim späteren Ablauf unseres Transformationsprogramms die MD-Problematik bei den Variablen MOTIV1 bis ANDERE bereits gelöst sein, wenn die **Zählen**-Anweisung an die Reihe kommt. Bei Personen, die den Fragebogenteil 4a nicht korrekt bearbeitet haben, wird also gelten:

$$\text{MOTIV1} = \text{MOTIV2} = \dots = \text{ANDERE} = \text{SYSMIS}$$

Wir müssen die folgende wichtige Eigenschaft der **Zählen**-Anweisung beachten: Ihre Ergebnisvariable hat *stets* einen validen Wert größer oder gleich Null. Wenn ein Fall z.B. bei allen kritischen Variablen den - nicht zu zählenden - Wert SYSMIS hat, resultiert das valide Ergebnis Null! In dieser Situation wissen wir aber *nichts* von den Motiven der Person und dürfen ihr keine Motivationslosigkeit (POLYMOT = 0) unterstellen.

Weil im konkreten Beispiel das Zählergebnis Null generell als irregulär einzustufen ist, könnten wir durch ein gewöhnliches (unbedingtes) Umkodieren

$$0 \rightarrow \text{SYSMIS}$$

dafür sorgen, dass ein Fall bei POLYMOT den Wert SYSMIS erhält, wenn er den Fragebogenteil 4a nicht korrekt bearbeitet hat. Im Allgemeinen ist jedoch eine bedingte Datentransformation erforderlich, um MD-belastete Zählergebnisse zu verhindern. Wir wollen das generelle Verfahren der Übung halber auch im konkreten Fall einsetzen und formulieren mit Hilfe der in Abschnitt 6.4.2.1.1 beschriebenen Funktion NMISS die folgende Bedingung

$$\text{NMISS}(\text{MOTIV1 TO ANDERE}) = 0$$

Klicken Sie bitte in der Dialogbox **Häufigkeiten von Werten in Fällen zählen** auf den **Falls**-Schalter, und tragen Sie die vorgeschlagene Bedingung ein. Wenn Sie dann **Weiter** machen und die Hauptdialogbox mit **Einfügen** quittieren, erhalten Sie im Syntaxfenster die folgenden Kommandos:

```
DO IF (nmiss(motiv1 to andere) = 0) .
COUNT
  polymot = motiv1 motiv2 motiv3 motiv4 motiv5 andere (1) .
VARIABLE LABELS polymot 'Anzahl der Motive für die Kursteilnahme' .
END IF .
EXECUTE .
```

Was hier *zählt*, ist offenbar das COUNT-Kommando. Es enthält im Wesentlichen eine Liste der zu untersuchenden Variablen, gefolgt von einer eingeklammerten Liste der kritischen Werte. Das VARIABLE LABELS - Kommando hat SPSS aufgrund unserer Eintragung im **Label**-Textfeld erstellt.

Das Zählergebnis wird nur dann ermittelt und der neuen Variablen POLYMOT als Wert zugewiesen, wenn die Bedingung im DO IF - Kommando erfüllt ist. Anderenfalls behält POLYMOT den Initialisierungswert SYSMIS.

## 6.7 Erstellung der Fertigdatendatei mit dem Transformationsprogramm

Aufgrund der KFA-bezogenen Übungsaufgaben in den Abschnitten 6.2 (Erstellung von DEKA-DE durch Rekodierung von GEBJ, Umkodieren der negativ formulierten LOT-Fragen), 6.4 (Berechnung von IDGEW, LOT, AERGAM und AERGZ), 6.5 (MD-Behandlung für die Motiv- und für die Methoden-Variablen) und 6.6 (Auszählen der Kursmotive) sollten jetzt alle vorläufig im KFA-Projekt benötigten Transformationskommandos in einem Syntaxfenster stehen.

### 6.7.1 Transformationsprogramm vervollständigen

Um daraus ein komfortables SPSS-Programm zu machen, das die Rohdatendatei **kfar.sav** selbstständig einliest, dann die so entstandene Arbeitsdatei transformiert und schließlich als Fertigdatendatei **kfa.sav** auf die Festplatte schreibt, müssen wir an den Anfang des Syntaxfensters noch ein GET-Kommando zum Öffnen von **kfar.sav** und ans Ende noch ein SAVE-Kommando zum Sichern in **kfa.sav** setzen. Wie Sie das GET-Kommando produzieren können, haben Sie schon in Abschnitt 5.2 erfahren. Wenn Sie das Kommando jetzt erzeugen lassen, erscheint es am Ende des Syntaxfensters, und Sie müssen es an den Anfang verschieben. Um das SAVE-Kommando zu generieren, wechseln wir ins Datenfenster und aktivieren mit **Datei > Speichern unter...** die zugehörige Dialogbox. Dann tragen wir den gewünschten Dateinamen **kfa.sav** ein und erzeugen mit **Einfügen** das benötigte SAVE-Kommando.

Zwei Hinweise zur Ausgabedatei eines Transformationsprogramms:

- Verwenden Sie niemals dieselbe Datei als Quelle und Ziel des Transformationsprogramms. Schreiben Sie also keinesfalls mit Ihrem Transformationsprogramm in die Rohdatendatei. Wenn Sie der Empfehlung in Abschnitt 6.1.2 folgend für die Rohdatendatei das Schreibschutzattribut gesetzt haben, kann dieses Desaster auch nicht versehentlich passieren.
- Bei der Ausführung des Transformationsprogramms darf für seine Ausgabedatei, also für die Fertigdatendatei, das Schreibschutzattribut natürlich nicht gesetzt sein.

Schließlich sollte Ihr Syntaxfenster ungefähr so aussehen:

```
GET
  FILE='U:\Eigene Dateien\SPSS\KFAR.SAV'.
EXECUTE .

* DEKADE.
RECODE
  gebj
  (60 thru 69=1) (70 thru 79=2) INTO dekade .
VARIABLE LABELS dekade "Dekade".
EXECUTE .

* LOT-Fragen umkodieren.
RECODE
  lot3 lot4 lot5 lot12 (5=1) (4=2) (2=4) (1=5) .
EXECUTE .

* IDGEW.
COMPUTE idgew = groesse - 100 .
VARIABLE LABELS idgew 'Idealgewicht nach der Formel: Größe - 100' .
EXECUTE .

* LOT berechnen.
COMPUTE lot = MEAN.6(lot1,lot3,lot4,lot5,lot8,lot9,lot11,lot12) .
VARIABLE LABELS lot 'LOT-Optimismus' .
EXECUTE .

* AERGAM berechnen.
COMPUTE aergam = (aergo + aergm)/2 .
VARIABLE LABELS aergam 'Mittel der Ärger-Variablen' .
EXECUTE .

* AERGFZ berechnen.
COMPUTE aergz = aergm - aergo .
VARIABLE LABELS aergz 'Ärger-Zuwachs durch die KFA' .
EXECUTE .

* MD-Behandlung für die Motiv-Variablen.
DO IF (SUM(motiv1 to andere) = 0) .
RECODE
  motiv1 motiv2 motiv3 motiv4 motiv5 andere (0=SYSMIS) .
END IF .
EXECUTE .

* MD-Behandlung für die Methoden-Variablen, Zelle (1,1) der Tabelle.
DO IF (smg=1 and nmiss(meth1 to meth3) < 3) .
RECODE
  meth1 meth2 meth3 (SYSMIS=0) .
END IF .
EXECUTE .

* MD-Behandlung für die Methoden-Variablen, Zelle (1,2) der Tabelle.
DO IF (smg=1 and nmiss(meth1 to meth3) = 3) .
RECODE
  smg (1=SYSMIS) .
END IF .
EXECUTE .

* MD-Behandlung für die Methoden-Variablen, Zelle (2,1) der Tabelle.
DO IF ((smg = 0) and (nmiss(meth1 to meth3) < 3)) .
RECODE
  smg (0=1) .
END IF .
EXECUTE .
DO IF ((smg = 0) and (nmiss(meth1 to meth3) < 3)) .
RECODE
  meth1 meth2 meth3 (SYSMIS=0) .
END IF .
EXECUTE .

* MD-Behandlung für die Methoden-Variablen, Zelle (2,2) der Tabelle.
DO IF (smg=0 and nmiss(meth1 to meth3) = 3) .
RECODE
  meth1 meth2 meth3 (SYSMIS=0) .
END IF .
EXECUTE .
```

```

* MD-Behandlung für die Methoden-Variablen, Zelle (3,1) der Tabelle.
DO IF ((nmiss(smj) = 1) and (nmiss(meth1 to meth3) < 3)) .
RECODE
  smj (SYSMIS=1) .
END IF .
EXECUTE .
DO IF ((nmiss(smj) = 1) and (nmiss(meth1 to meth3) < 3)) .
RECODE
  meth1 meth2 meth3 (SYSMIS=0) .
END IF .
EXECUTE .

* POLYMOT berechnen.
DO IF (nmiss(motiv1 to andere) = 0) .
COUNT
  polymot = motiv1 motiv2 motiv3 motiv4 motiv5 andere (1) .
VARIABLE LABELS polymot 'Anzahl der Motive für die Kursteilnahme' .
END IF .
EXECUTE .

SAVE OUTFILE='U:\Eigene Dateien\SPSS\KFA.SAV'
  /COMPRESSED.

```

Hierzu sind einige Anmerkungen erforderlich:

- Zwischen manchen Kommandos sind der Übersichtlichkeit halber Leerzeilen eingefügt worden. Man darf aber auf keinen Fall *innerhalb* eines Kommandos (d.h. zwischen dem Kommandonamen und dem abschließenden Punkt) eine Leerzeile einfügen (vgl. Abschnitt 5.4).
- Die mit einem Sternchen (\*) eingeleiteten Zeilen beinhalten *Kommentare*, die nachträglich eingefügt wurden, um die spätere Orientierung im Programm zu erleichtern.  
**Wichtig:** Ein Kommentar hat ebenfalls Kommandostatus und muss daher unbedingt mit einem Punkt abgeschlossen werden. Sonst erstreckt sich der Kommentar bis zur nächsten Zeile, die entweder komplett leer ist oder mit einem Punkt endet.
- Das GET-Kommando am Anfang des Programms *überschreibt die aktuelle Arbeitsdatei ohne Nachfrage!* Wenn Sie im Datenfenster manuelle Korrekturen vornehmen, diese nicht sichern, sondern anschließend ein GET-Kommando (via Syntaxfenster) ausführen lassen, dann sind die manuellen Korrekturen verloren.
- Das SAVE-Kommando überschreibt eine eventuell vorhandene Datei **kfa.sav** ohne Nachfrage, was jedoch bei der in diesem Manuskript vorgeschlagenen Arbeitsweise (vgl. Abschnitt 6.1.1) unproblematisch ist.

Eventuell legen Sie Wert darauf, dass auch die neu berechneten Variablen mit einer optimalen Anzahl von Dezimalstellen angezeigt werden. Eine manuelle Einstellung (vgl. Abschnitt 3.2.2) ist wenig attraktiv, weil unser Transformationsprogramm ja mit einiger Wahrscheinlichkeit mehrfach ausgeführt werden muss. Die bessere Alternative besteht darin, unser Programm um ein FORMATS-Kommando zu erweitern, das die Attribute automatisch setzt:

```

formats dekade (f1.0) idgew (f2.0) aergam (f3.1)
  aergz (f3.0) polymot (f1.0).

```

Im Ausdruck „(fb.d)“ legt man mit *b* die Gesamtbreite der Wertausgabe (Attribut **Spaltenformat**) und mit *d* die Anzahl der Dezimalstellen fest.

Mit den folgenden Kommandos wird die Breite der Datenfensterspalte (Attribut **Spalten**) und das Messniveau für die neuen Variablen eingestellt, wobei SCALE für Intervallskalenqualität steht:

```

variable width dekade to polymot (7).
variable level dekade (ordinal) / idgew to polymot (scale).

```

Fügen Sie die Kommandos zur Deklaration von Variablenattributen am Ende des Transformationsprogramms ein (unmittelbar vor dem SAVE-Kommando).

Damit ist das Transformationsprogramm zum KFA-Projekt fertig. Falls noch nicht geschehen, müssen Sie es unbedingt sichern, z.B. in das Verzeichnis **U:\Eigene Dateien\SPSS** unter dem oben vorgeschlagenen Dateinamen **kfat.sps**.

### 6.7.2 Transformationsprogramm ausführen

Lassen Sie das Transformationsprogramm ausführen, z.B. mit

#### **Ausführen > Alles**

Wenn Sie anschließend im Hauptausgabefenster keine Spur des Programmlaufs finden, ist alles glatt gegangen. Anderenfalls erscheinen dort Fehlermeldungen und/oder Warnungen in einem mit **Log** betitelten Ausgabeblock. Da alle Kommandos Ihres Programms von SPSS erstellt wurden, sollte dies eigentlich nicht passieren.

Ältere Warnungen bzw. Fehlermeldungen sollten vor einem Lauf des Transformationsprogramms aus dem Ausgabefenster gelöscht werden, um Unklarheiten zu vermeiden.

Ein gelungener Lauf des Transformationsprogramms hinterlässt zwar im Ausgabefenster keine Spuren, wirkt sich aber nachhaltig auf das Datenfenster aus. Dort erscheinen z.B. am rechten Rand der Datenmatrix die neuen Variablen.

Sie dürfen aber Ihre Erfolgskontrolle keinesfalls auf das Datenfenster beschränken, sondern müssen unbedingt das Ausgabefenster auf Fehlermeldungen und Warnungen überprüfen. SPSS stoppt nämlich die Programmausführung **nicht** beim Auftreten des ersten fehlerhaften Kommandos, sondern ignoriert das fehlerhafte Kommando und macht unverdrossen mit den nächsten Kommandos weiter. Diese arbeiten aber möglicherweise aufgrund des vorangegangenen Fehlers mit falschen Zwischenergebnissen und produzieren Unsinn. Es kann also leicht passieren, dass nach einem fehlerbehafteten Lauf des Transformationsprogramms alle erwarteten neuen Variablen vorhanden sind, jedoch unsinnige Werte enthalten.

---

## 7 Prüfung der zentralen Projekt-Hypothesen

### 7.1 Entscheidungsregeln beim Hypothesentesten

In diesem Abschnitt werden einige Grundprinzipien der Inferenzstatistik am Beispiel unserer allgemeinspsychologischen Hypothese demonstriert. Dabei handelt es sich nicht um eine systematische Behandlung des Themas, die erheblich mehr Platz beanspruchen würde. Im Wesentlichen sollen die statistischen Entscheidungsregeln so präsentiert werden, dass sie mit Hilfe der SPSS-Ausgaben unmittelbar umgesetzt werden können. Zumindest in älteren Statistikbüchern findet man nämlich Formulierungen mit wenig Bezug zu den heute üblichen Ausgaben von Statistikprogrammen.

Wenn mit  $\mu_O$  der Erwartungswert (Populationsmittelwert) des Merkmals AERGO und mit  $\mu_M$  der Erwartungswert des Merkmals AERGM bezeichnet wird, dann lautet unser zentrales, allgemeinspsychologisches KFA-Testproblem:

$$H_0 : \mu_M \leq \mu_O \quad \text{vs.} \quad H_1 : \mu_M > \mu_O$$

Mit Hilfe der Differenzvariablen  $AERZ := AERGM - AERGO$ , deren Erwartungswert mit  $\mu_Z$  bezeichnet werden soll, lässt sich das Testproblem äquivalent noch kompakter formulieren:

$$H_0 : \mu_Z \leq 0 \quad \text{vs.} \quad H_1 : \mu_Z > 0$$

Bei der Reformulierung wird die folgende, generell gültige, Identität ausgenutzt:

$$\mu_Z = \mu_M - \mu_O$$

Wir wollen noch voraussetzen, dass die Differenzvariable AERZ normalverteilt sei mit dem Erwartungswert  $\mu_Z$  und der Varianz  $\sigma_Z^2$ :

$$AERZ \sim N(\mu_Z, \sigma_Z^2)$$

Für die  $n$  AERZ-Beobachtungen in der Stichprobe nehmen wir an, dass sie durch **unabhängiges** „Ziehen“ aus der eben beschriebenen Population entstanden sind. Das schon in Abschnitt 1 betonte Unabhängigkeitsprinzip ist die zentrale Forderung in unserem **Stichprobenmodell** über die Gewinnung der empirischen Daten.

Bei der inferenzstatistischen Lösung des beschriebenen Testproblems benötigen wir eine so genannte **Prüfstatistik**  $T$ , die aus den Stichprobendaten berechnet werden kann und folgende Eigenschaften besitzt:

- Sie ist indikativ für Abweichungen der wahren Populationsverteilung von der Nullhypothesenbehauptung, wird also tendenziell umso größer, je weiter der Verteilungsparameter  $\mu_Z$  in positiver Richtung vom Wert Null entfernt ist. Sie quantifiziert also, wie gut bzw. schlecht die Nullhypothese mit den Daten vereinbar ist.
- Es ist bekannt, welcher Verteilung die Prüfstatistik  $T$  bei gültiger Nullhypothese folgt, also bei  $\mu_Z \leq 0$ . Damit lässt sich für den konkreten Wert  $T_{\text{emp}}$  der Prüfstatistik in einer bestimmten Stichprobe beurteilen, mit welcher Wahrscheinlichkeit eine Nullhypothesen-Population derartige Werte generiert. Ist diese Wahrscheinlichkeit sehr klein, liegt der Schluss nahe, dass die Stichprobe *nicht* aus einer Nullhypothesen-Population stammt.



In der oben beschriebenen Situation hat sich die folgende Prüfstatistik  $T_Z$  bewährt (mit  $Z$  als Abkürzung für AERGGZ):

$$T_Z := \frac{\bar{Z}}{S_Z} \sqrt{n} \quad \text{mit} \quad \bar{Z} := \frac{1}{n} \sum_{i=1}^n Z_i \quad \text{und} \quad S_Z := \sqrt{\frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z})^2}$$

Weil  $\frac{S_Z}{\sqrt{n}}$  gerade der Standardfehler des Stichprobenmittelwerts ist, handelt es sich bei  $T_Z$  um den Quotienten aus dem Stichprobenmittelwert und seinem Standardfehler. Prüfgrößen von analoger Bauart sind uns schon bei den „quick-and-dirty“-Tests zur Schiefe bzw. Wölbung einer Verteilung in Abschnitt 4.6 begegnet.

$T_Z$  erfüllt die obigen Anforderungen:

- Der Stichprobenmittelwert  $\bar{Z}$  wächst stochastisch mit seinem Erwartungswert  $\mu_Z$ . Bleibt gleichzeitig der „Nebenparameter“  $\sigma_Z^2$  konstant, so hat auch dessen erwartungstreuer Schätzer  $S_Z$  keine Wachstumstendenz. Folglich sind von  $T_Z$  umso größere Werte zu erwarten, je weiter der Erwartungswert  $\mu_Z$  bei konstanter Varianz  $\sigma_Z^2$  über den Wert Null hinauswächst.
- Für  $\mu_Z = 0$  besitzt  $T_Z$  (bei beliebigem Nebenparameter  $\sigma_Z^2$ ) eine t-Verteilung mit  $n - 1$  Freiheitsgraden. Damit kennen wir das Verhalten der Prüfstatistik am Rand der Nullhypothese. Dieses Wissen genügt, weil die bei einer Testentscheidung relevante Überschreitungswahrscheinlichkeit unter der  $H_0$  (siehe unten) am Rand der Nullhypothese (also bei  $\mu_Z = 0$ ) maximal wird. Ist sie am Rand klein genug, dann gilt dies auch für alle anderen Verteilungen in der Nullhypothese.

Aufgrund dieser Voraussetzungen kann man zum Wert  $T_{\text{emp}}$  der Prüfstatistik  $T_Z$  für eine konkrete Stichprobe die folgende **Überschreitungswahrscheinlichkeit** bestimmen:

Mit welcher Wahrscheinlichkeit nimmt die Prüfstatistik  $T_Z$  **bei Gültigkeit der Nullhypothese** (genauer: bei  $\mu_Z = 0$ ) einen Wert größer oder gleich  $T_{\text{emp}}$  an?

Diese Wahrscheinlichkeit wollen wir mit  $\mathbf{P}_{H_0}(T_Z \geq T_{\text{emp}})$  bezeichnen. Sie wird von SPSS berechnet und in der Ausgabe zum t-Test für gepaarte Stichproben mit **Sig.** überschrieben.<sup>1</sup>

Bei einem akzeptierten **Fehlerrisiko erster Art** von  $\alpha = 5\%$  verwendet man die folgende **Entscheidungsregel**:

$$\mathbf{P}_{H_0}(T_Z \geq T_{\text{emp}}) \begin{cases} \geq 0,05 & \Rightarrow H_0 \text{ beibehalten} \\ < 0,05 & \Rightarrow H_0 \text{ verwerfen} \end{cases} \quad (7-1)$$

Die Nullhypothese wird also abgelehnt, wenn die Prüfstatistik einen Wert annimmt, der bei Gültigkeit der  $H_0$  nur relativ selten (mit einer Wahrscheinlichkeit von  $< 5\%$ ) erreicht oder gar übertroffen wird.

In Statistiklehrbüchern wird oft ein **kritischer Wert**  $T_{\text{krit}}$  aufgrund der Kenntnis über die Verteilung von  $T_Z$  unter der  $H_0$  (genauer: bei  $\mu_Z = 0$ ) so bestimmt, dass gilt:

$$\mathbf{P}_{H_0}(T_Z \geq T_{\text{krit}}) = 0,05$$

<sup>1</sup> Leider gibt SPSS beim t-Test für gepaarte Stichproben ausschließlich die *zweiseitige* Überschreitungswahrscheinlichkeit aus (siehe unten), während wir unsere alltagspsychologische KFA-Hypothese mit gutem Grund einseitig formuliert haben und daher auch die einseitige Überschreitungswahrscheinlichkeit betrachten.

$T_{\text{krit}}$  ist in unserer Situation gerade das 95%-Quantil der t-Verteilung mit  $n - 1$  Freiheitsgraden.

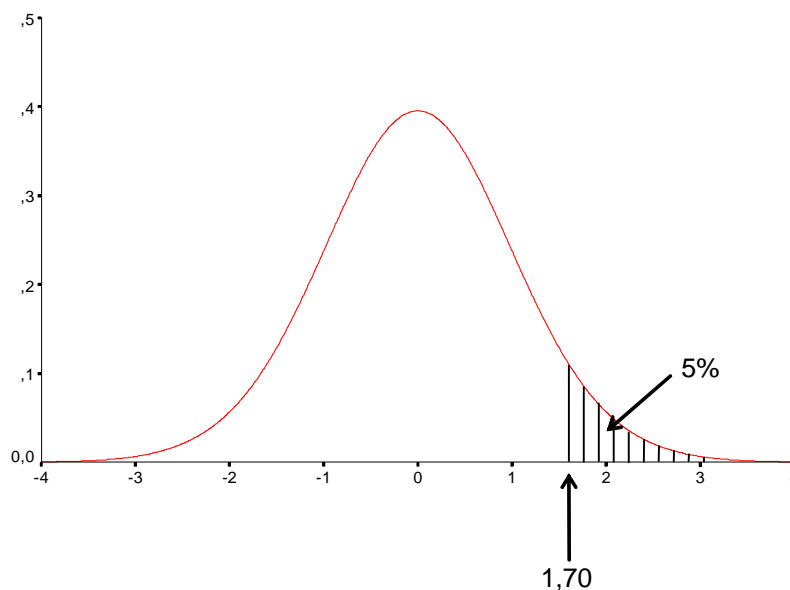
Bei unserer Stichprobengröße  $n = 31$  erhalten wir z.B.  $T_{\text{krit}} = 1,70$ .

Damit kann obige Entscheidungsregel äquivalent folgendermaßen formuliert werden:

$$T_{\text{emp}} \begin{cases} \leq T_{\text{krit}} \Rightarrow H_0 \text{ beibehalten} \\ > T_{\text{krit}} \Rightarrow H_0 \text{ verwerfen} \end{cases} \quad (7-2)$$

Wir haben übrigens bei den „quick-and-dirty“-Tests in Abschnitt 4.6 die Testentscheidung anhand von kritischen Werten kennen gelernt. Dort waren wir ausnahmsweise in der Lage, keine Überschreitungswahrscheinlichkeiten zu kennen, aber die kritischen Werte (als Quantile der Standardnormalverteilung) besonders leicht ermitteln zu können.

Die folgende Abbildung zeigt die Dichte unserer Prüfverteilung ( $t_{30}$ ) sowie den  $H_0$ -**Ablehnungsbereich** bei einseitiger Fragestellung im Sinne unserer KFA-Hypothese (|||):



Die rot eingezeichnete Dichte beschreibt das Verteilungsverhalten der Zufallsgröße  $T_Z$ , für die eine einzelne Realisation folgendermaßen zu ermitteln ist: Ziehe aus einer Population mit

$$\text{AERGZ} \sim N(0, \sigma_Z^2)$$

eine Zufallsstichprobe der Größe  $n = 31$ , ermittle die AERGZ-Werte und berechne  $\frac{\bar{Z}}{S_Z} \sqrt{n}$ .

Wir kommen zu einer Testentscheidung, indem wir unser Stichprobenergebnis  $T_{\text{emp}}$  vor dem Hintergrund dieses Erwartungshorizonts beurteilen. Wir lehnen die Nullhypothese ab, wenn sie als Generator unserer Daten unplausibel ist.

Wenn wir aus einer Nullhypothesen-Population (genauer: bei  $\mu_Z = 0$ ) eine Zufallsstichprobe der Größe  $n = 31$  ziehen und  $T_{\text{emp}}$  ermitteln, werden wir mit der Wahrscheinlichkeit  $\alpha = 0,05$  einen Wert größer oder gleich  $T_{\text{krit}} = 1,70$  erhalten und falsch gegen die  $H_0$  entscheiden, also einen Fehler erster Art begehen. Der  $\alpha$ -Wert sollte umso niedriger angesetzt werden, je gravierender (schädlicher, teurer) das irrtümliche Ablehnen einer gültigen Nullhypothese ist.

Das Risiko, bei Gültigkeit der *Alternativhypothese* falsch zu entscheiden (**Fehler zweiter Art**,  $\beta$ -Fehler), ist umso kleiner, ...

- je stärker der wahre Lageparameter  $\mu_Z$  von der Nullhypothese  $\{\mu_Z \leq 0\}$  entfernt ist,

- je kleiner die Streuung  $\sigma_Z^2$  ist,
- je größer die Teststärke (Power) des verwendeten Signifikanztests ist, d.h. je wahrscheinlicher unter der Alternativhypothese ein signifikantes Ergebnis erzielt wird.

Wie Sie aus der Stichprobenumfangsplanung in Abschnitt 1.3.2 wissen, kann man für jede konkret vorgegebene Alternativhypothesenverteilung, also für bestimmte Werte der unbekanntenen Verteilungsparameter  $\mu_Z$  und  $\sigma_Z^2$  ...

- für eine erwünschte Teststärke (z.B.  $1 - \beta = 95\%$ ) die erforderliche Stichprobengröße ermitteln,
- das  $\beta$ -Fehler-Risiko zu einer festen Stichprobengröße ausrechnen.

Passend zu unserer allgemeinspsychologischen KFA-Hypothese haben wir bislang das einseitige Testproblem behandelt. Wir wollen noch das folgende **zweiseitige Testproblem** betrachten:

$$H_0 : \mu_M = \mu_O \quad \text{vs.} \quad H_1 : \mu_M \neq \mu_O$$

bzw.

$$H_0 : \mu_Z = 0 \quad \text{vs.} \quad H_1 : \mu_Z \neq 0$$

Die  $H_0$  des zweiseitigen Tests ist gerade identisch mit dem *Rand* der  $H_0$  zum einseitigen Test.

Wir verwenden beim zweiseitigen Test dieselbe Prüfstatistik  $T_Z$  wie beim einseitigen Test. Nun sind aber *betragsmäßig* große  $T_{\text{emp}}$ -Werte mit *positivem oder negativem* Vorzeichen indikativ für eine Abweichung von der Nullhypothese. Nach einem generellen Prinzip der Testkonstruktion müssen *alle* Elemente der Alternativhypothese (mit  $\mu_Z < 0$  oder  $\mu_Z > 0$ ) eine faire Chance haben, sich in einem signifikanten Ergebnis zu artikulieren. Anderenfalls resultiert ein so genannter *verfälschter Test*. Daher ist die Überschreitungswahrscheinlichkeit

$$P_{H_0}(|T_Z| \geq |T_{\text{emp}}|)$$

zu ermitteln und folgende Entscheidungsregel zu verwenden:

$$P_{H_0}(|T_Z| \geq |T_{\text{emp}}|) \begin{cases} \geq 0,05 & \Rightarrow H_0 \text{ beibehalten} \\ < 0,05 & \Rightarrow H_0 \text{ verwerfen} \end{cases} \quad (7-3)$$

Der kritische Werte  $T_{\text{krit},2}$  zum zweiseitigen Test ist so zu bestimmen, dass gilt:

$$P_{H_0}(|T_Z| \geq |T_{\text{krit},2}|) = 0,05$$

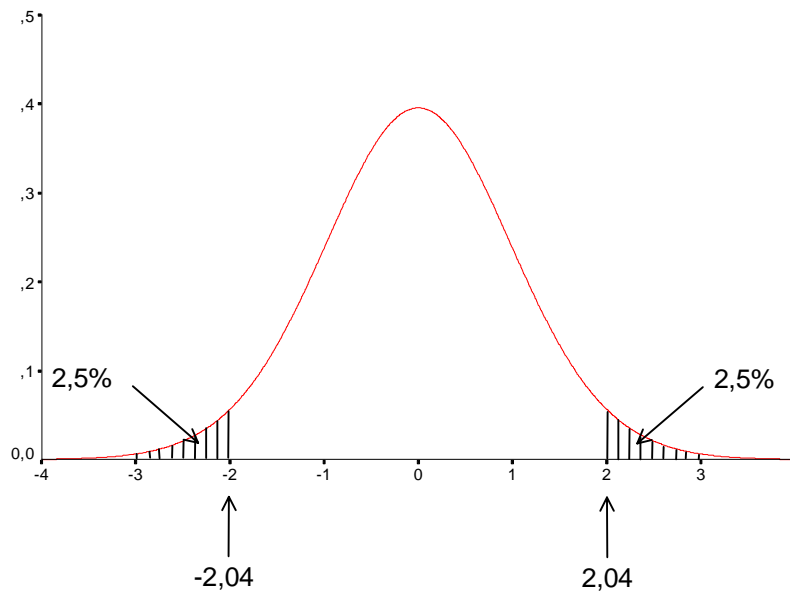
Bei unserer Stichprobengröße  $n = 31$  erhalten wir  $T_{\text{krit},2} = \pm 2,04$ .

Aufgrund der Symmetrie der Prüfverteilung gilt für  $T_{\text{emp}} \geq 0$ :

$$P_{H_0}(T_Z \geq T_{\text{emp}}) = \frac{1}{2} \cdot P_{H_0}(|T_Z| \geq |T_{\text{emp}}|) \quad (7-4)$$

Die Überschreitungswahrscheinlichkeit des einseitigen t-Tests ergibt sich also durch Halbieren aus der Überschreitungswahrscheinlichkeit des zweiseitigen t-Tests, sofern die Prüfgröße das behauptete Vorzeichen besitzt. Dieser Zusammenhang ist wichtig in der statistischen Praxis mit SPSS, weil dieses Programm bei t-Tests häufig nur die zweiseitige Überschreitungswahrscheinlichkeit ausgibt. Sie dürfen aber den Zusammenhang in Gleichung (7-4) keinesfalls auf beliebige Tests generalisieren. Wir werden z.B. im Zusammenhang mit der Kreuztabellenanalyse den exakten Test von Fisher kennen lernen, bei dem eine analoge Gleichung *nicht* gilt.

Bei zweiseitiger Fragestellung haben wir zwei symmetrisch angeordnete Ablehnungsbereiche:



## 7.2 Zu den Voraussetzungen unserer Hypothesentests

Der t-Test für gepaarte Stichproben, mit dem wir unsere allgemeinspsychologische Hypothese prüfen wollen, setzt voraus, dass die Differenzvariable AERGZ normalverteilt ist (vgl. Abschnitt 7.1). Diese Normalverteilungsannahme soll anschließend mit der SPSS-Prozedur zur explorativen Datenanalyse geprüft werden.

Unsere differentialpsychologische Hypothese bezieht sich auf den Steigungskoeffizienten  $\beta_1$  in der linearen Regression von AERGAM auf LOT:

$$\text{AERGAM} = \beta_0 + \beta_1 \text{LOT} + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

Die Hypothesen des Testproblems lauten:

$$H_0 : \beta_1 \geq 0 \quad \text{vs.} \quad H_1 : \beta_1 < 0$$

Es kommt eine Prüfstatistik zum Einsatz, die sich im vorliegenden Fall der bivariaten Regression besonders bequem mit Hilfe der Stichprobenkorrelation  $r$  zwischen Kriterium und Regressor notieren lässt:

$$T_r := \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

Sie ist bei gültiger Nullhypothese (genauer: bei  $\beta_1 = 0$ ) t-verteilt mit  $n - 2$  Freiheitsgraden, sofern die Voraussetzungen des Regressionsmodells erfüllt sind, die anschließend der bequemeren Schreibweise halber für ein Kriterium  $Y$  und einen Regressor  $X$  angegeben sind:

### 1) Linearität

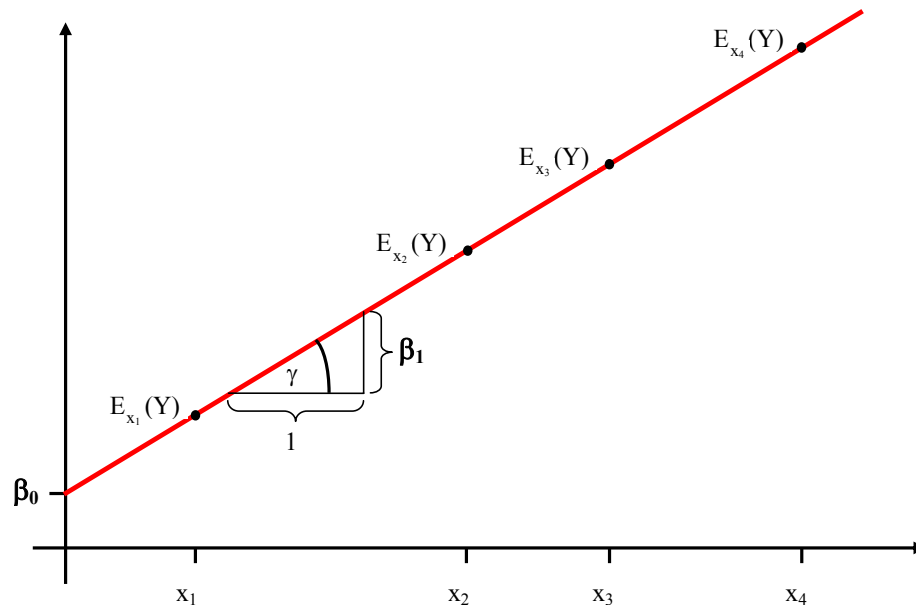
Der Erwartungswert (Mittelwert)  $E_X(Y)$  von  $Y$  für einen bestimmten  $X$ -Wert hängt **linear** von  $X$  ab:

$$E_X(Y) = \beta_0 + \beta_1 X$$

Für beliebige  $X$ -Ausprägungen liegen die zugehörigen Erwartungswerte  $E_X(Y)$  auf der **Regressionsgeraden** durch die Punktepaare

$$(X, \beta_0 + \beta_1 X)$$

Dabei ist  $\beta_0$  der Schnittpunkt der Regressionsgeraden mit der Y-Achse (Ordinatenabschnitt) und  $\beta_1$  die Steigung der Regressionsgeraden (der Tangens des Winkels  $\gamma$  der Regressionsgeraden mit der X-Achse).



Zur Interpretation des Koeffizienten  $\beta_1$ : Erhöht man  $X$  um eine Einheit, so steigt modellgemäß der Erwartungswert  $E_X(Y)$  um  $\beta_1$  Einheiten an.

## 2) Normalität der Residuen

Für die (nicht direkt beobachtbare) Fehler- bzw. Residualvariable  $\varepsilon$  wird angenommen, dass sie **normalverteilt** ist mit Erwartungswert 0 und Varianz  $\sigma^2$ . Sie dürfen sich vorstellen, dass es für jede  $X$ -Ausprägung eine **Normalverteilung** potentieller  $\varepsilon$ -Werte gibt, aus der **zufällige** Realisationen gezogen werden, die zusammen mit dem konstanten Anteil  $\beta_0 + \beta_1 X$  die Realisationen der abhängigen Variablen  $Y$  ergeben.

## 3) Varianzhomogenität der Residuen (Homoskedastizität)

Die Normalverteilungen der  $\varepsilon$ -Variablen zu den verschiedenen  $X$ -Ausprägungen haben alle **die-selbe Varianz**  $\sigma^2$ .

## 4) Unabhängigkeit der Residuen

Die Residuen zu den einzelnen Beobachtungen (Fällen) in der Stichprobe sind unkorreliert. Wegen ihrer Normalverteilung sind sie damit auch stochastisch unabhängig.

Hinsichtlich der Verteilungsvoraussetzungen ist zu betonen:

- Es wird keine Annahme über die Verteilung des Regressors gemacht.
- Es wird keine Annahme über die univariate Verteilung des Kriteriums gemacht.
- Es sind die **Residuen des Modells**, die bestimmte Verteilungsvoraussetzungen erfüllen müssen (Unabhängigkeit, Normalität, Homoskedastizität)

Für methodisch besonders Interessierte soll noch eine alternative Darstellung für  $T_r$  vorgeführt werden, die von eher anwendungsorientierten Lesern gefahrlos übersprungen werden kann.

Weil der Stichprobenschätzer  $b_1$  des Steigungskoeffizienten in folgender Beziehung zur Stichprobenkorrelation  $r$  und den Schätzern  $s_Y$  und  $s_X$  für die Standardabweichungen des Kriteriums  $Y$  und des Regressors  $X$  steht

$$b_1 = r \frac{s_Y}{s_X}$$

und der Standardfehler zu  $b_1$  gleich

$$sf_{b_1} = \frac{s_Y}{s_X} \frac{\sqrt{1-r^2}}{\sqrt{n-2}}$$

ist (siehe z.B. Cohen et al. 2003, S. 42), kann auch die Prüfgröße  $T_r$  als Quotient aus einem Stichprobenschätzer und seinem Standardfehler geschrieben werden:

$$T_r = b_1 \frac{s_X}{s_Y} \frac{\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{b_1}{sf_{b_1}}$$

### 7.3 Verteilungsanalyse zu AERGF, AERGF und LOT

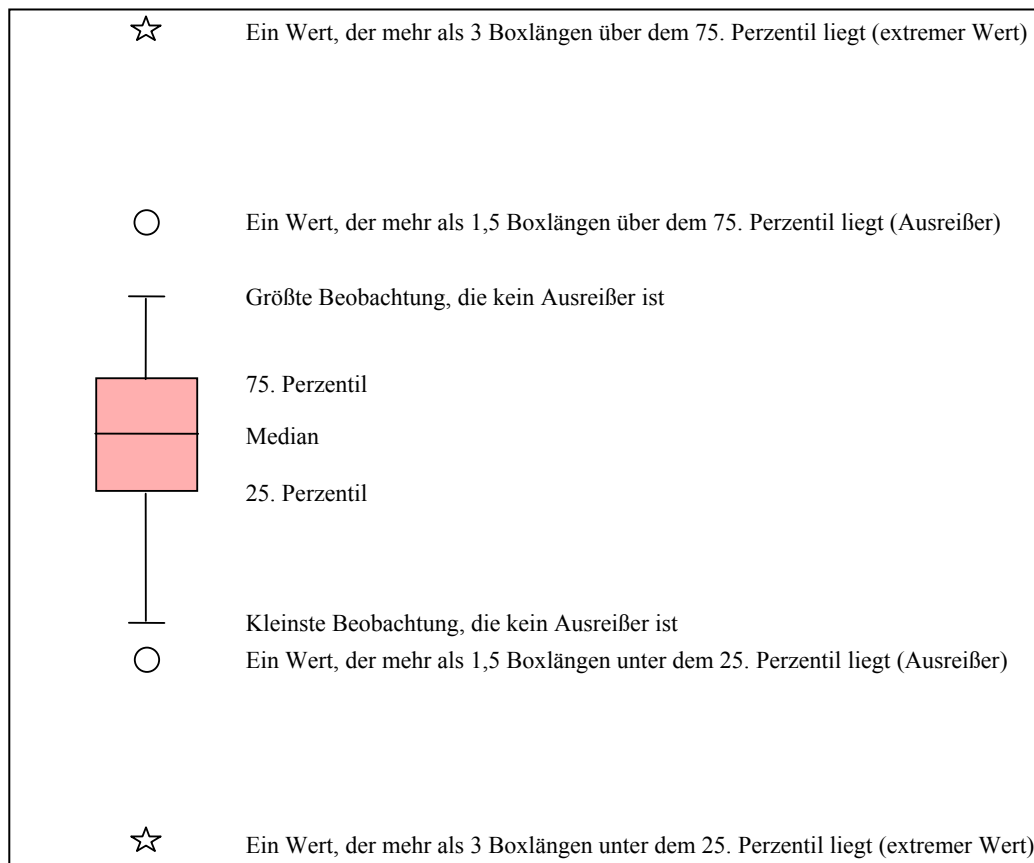
Für die folgenden Schritte wird eine aktive SPSS-Sitzung mit geöffneter Projekt-Fertigdaten-datei **kfa.sav** vorausgesetzt. Ob Sie die SPSS-Kommandos zu den anstehenden Analysen für spätere Wiederverwendung konservieren wollen, bleibt Ihnen überlassen.

Wir wollen zunächst die univariaten Verteilungen der abgeleiteten Variablen AERGF, AERGF und LOT untersuchen. Analog zu den Verteilungsanalysen in Abschnitt 4, die auch zur Datenprüfung dienten, wollen wir bei den Verteilungen der abgeleiteten Variablen auch auf Anomalien infolge fehlerhafter oder schlecht durchdachter Berechnungsvorschriften achten. Außerdem wollen wir noch eine weitere Gefahrenquelle für unser Forschungsprojekt ins Visier nehmen:

#### 7.3.1 Diagnose von Ausreißern

Als **Ausreißer** bezeichnet man extreme Werte, die zwar innerhalb des logisch möglichen Wertebereichs liegen, aber doch mit großer Wahrscheinlichkeit nicht aus der interessierenden Verteilung bzw. Population stammen. Diese Werte haben insbesondere auf parametrische Auswertungsverfahren einen starken, verzerrenden Einfluss. Daher wollen wir ab jetzt auch auf Ausreißer achten.

Dazu lassen wir uns für jede Variable einen **Boxplot** erstellen. Dieses beliebte Instrument der explorativen Datenanalyse zeigt auf prägnante Weise wesentliche Verteilungsinformationen, und ist zur Identifikation von Ausreißern sehr gut geeignet. Die Bestandteile eines Boxplots haben folgende Bedeutung:



Als Ursachen für Ausreißer kommen in Frage:

- Erhebungs- bzw. Erfassungsfehler  
Messwerte können falsch ermittelt oder fehlerhaft in die EDV übernommen worden sein.
- Besondere Umstände beim Merkmalsträger  
Bei einer Agrarstudie zum Ertrag verschiedene Getreidesorten kann z.B. der Boden in einer bestimmten Versuchsparzelle durch einen Ölunfall verseucht worden sein.

Eindeutig irreguläre Daten müssen natürlich entfernt werden. Sie können z.B. mit dem Dateneditor in der Rohdatendatei:

- einen Wert löschen
- einen Wert als MD-Indikator deklarieren
- einen kompletten Fall löschen

Natürlich dürfen Sie keine Daten eliminieren, weil sie Ihren Hypothesen widersprechen.

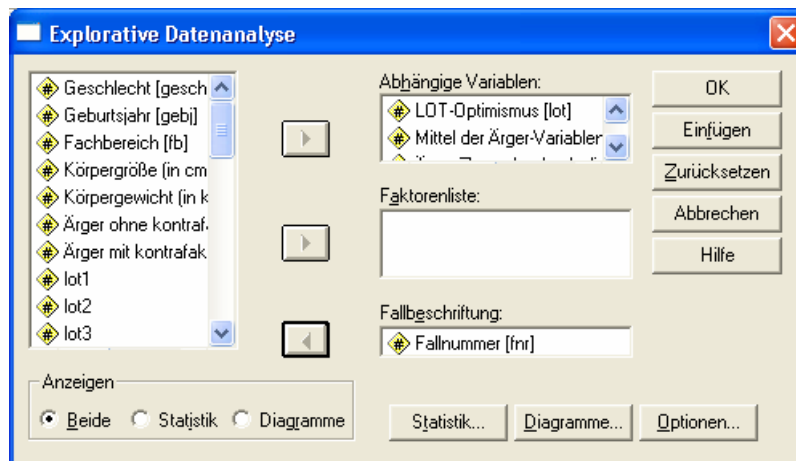
### 7.3.2 Die SPSS-Prozedur zur explorativen Datenanalyse

Für die eben geplanten Aufgaben (Ausreißerdiagnose und Verteilungsprüfung) eignet sich die SPSS-Prozedur zur explorativen Datenanalyse besser als die in Abschnitt 4 der Einfachheit halber bevorzugte Häufigkeitsanalyse. Natürlich können Sie in Zukunft auch die Verteilungen von Rohvariablen mit der leistungsfähigeren explorativen Datenanalyse untersuchen.

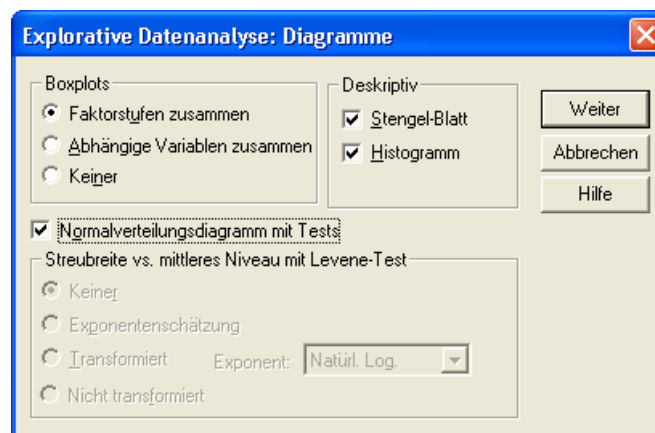
Starten Sie deren Dialogbox mit:

#### **Analysieren > Deskriptive Statistiken > Explorative Datenanalyse**

Transportieren Sie die Namen der drei zu untersuchenden Variablen in die Liste der **abhängigen Variablen**, und wählen Sie die Variable FNR zur Fallbeschriftung aus, damit mögliche Ausreißer durch ihre Fallnummer identifiziert werden können:



Fordern Sie in der **Diagramme**-Subdialogbox zusätzlich **Histogramme** sowie **Normalverteilungdiagramme mit Tests** an:



Das Kontrollkästchen zum Anfordern von Normalverteilungsanpassungstests (Kolmogorov-Smirnov und Shapiro-Wilk) hat SPSS wirklich sehr gut in der **Diagramme**-Subdialogbox der explorativen Datenanalyse versteckt.

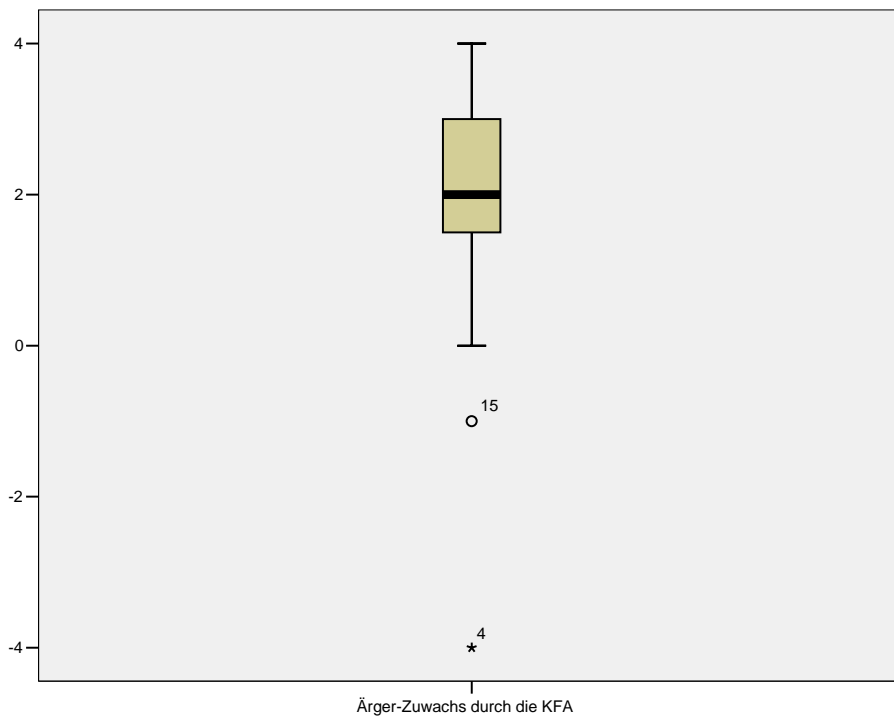
Der Klarheit halber soll nochmals betont werden, dass wir nur für die Variable AERGZ einen Normalverteilungsanpassungstest benötigen (vgl. Abschnitt 7.2). Allerdings sind die teilweise irrelevanten Ausgaben für AERGAM und LOT kein starker Grund dafür, zwei verschiedene Analysen anzufordern.

Wir erhalten im Viewer-Fenster u.a. für jede abhängige Variable einen **Boxplot**.

### 7.3.3 Ergebnisse für AERGZ


Bei der Ausreißer-Analyse gibt es nur einen Problemfall und zwar ausgerechnet bei der Variablen AERGZ, über die unsere zentrale KFA-Hypothese geprüft werden soll. Hier tanzt Fall Nr. 4 aus der Reihe:

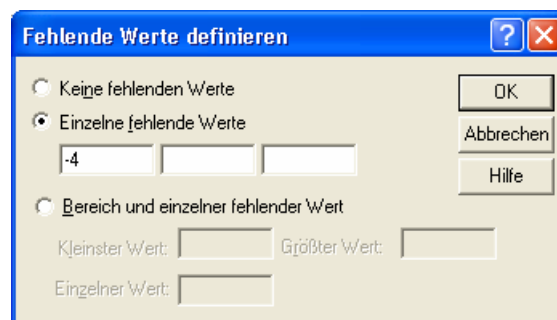




Diese Person hatte ohne KFA eine Ärgertemperatur von  $60^\circ$  gemeldet, die sich dann durch die KFA-Komponente auf  $20^\circ$  abkühlte. Zwar darf dieses Muster nicht a-priori als verdächtig gelten, weil es unserer Hypothese widerspricht, doch der Boxplot gibt eine klare Empfehlung, den Fall bei dieser Analyse auszuschließen. Allerdings scheut sich ein redlicher Forscher, Daten zu neutralisieren, die der eigenen Hypothese widersprechen.

Vor einer endgültigen Entscheidung wollen wir die Verteilung von AERGSZ noch weiter analysieren, da beim geplanten t-Test zur allgemeinspsychologischen KFA-Hypothese vorausgesetzt werden muss, dass AERGSZ (in der Population) normalverteilt ist.

Damit der extreme AERGSZ-Wert von Fall Nr. 4 die weitere Verteilungsanalyse nicht beeinflusst, soll er vorübergehend neutralisiert werden. Weil wir noch keine Methode kennen, komplette Fälle von einer Analyse fern zu halten (siehe Abschnitt 9), deklarieren wir den betroffenen Wert (= -4) als MD-Indikator. Auf diese Weise findet sich doch noch eine Gelegenheit, die Deklaration von benutzerdefinierten MD-Indikatoren zu üben. Markieren Sie in der Variablenansicht des Datenfensters die Zelle mit den **Fehlenden Werten** der Variablen AERGSZ, klicken Sie auf den Erweiterungsschalter , und tragen Sie den Wert -4 als einzelnen MD-Indikator ein:



Das folgende Histogramm zeigt, dass die AERGSZ-Verteilung *auch nach* Elimination von Fall Nr. 4 noch relativ deutlich von der Normalität abweicht:



Tatsächlich lehnen auch *nach* der Elimination des Ausreißers die beiden von SPSS angebotenen Normalverteilungstests (Kolmogorov-Smirnov und Shapiro-Wilk) die im t-Test benötigte Normalverteilungsannahme ab:

**Tests auf Normalverteilung**

	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistik	df	Signifikanz	Statistik	df	Signifikanz
Ärger-Zuwachs durch die KFA	,207	30	,002	,913	30	,018

a. Signifikanzkorrektur nach Lilliefors

Auch diese Testentscheidung folgt der in Abschnitt 7.1 beschriebenen Logik, wobei folgende Hypothesen zur Konkurrenz stehen:

$H_0$ : AERGFZ ist normalverteilt versus  $H_1$ : AERGFZ ist nicht normalverteilt

Die von SPSS berechnete Überschreitungswahrscheinlichkeit (**Signifikanz**) ist bei beiden Prüfstatistiken kleiner als 5%, so dass beide Tests übereinstimmend die Nullhypothese verwerfen. Dies ist vor allem deshalb ein ernst zu nehmender Befund, weil unsere Stichprobe relativ klein und damit die Power der Tests eher gering ist.

Bei einer *großen* Stichprobe besitzen die Normalitätstests eine hohe Power und decken auch kleinste (für die Validität des geplanten t-Tests irrelevante) Abweichungen von der Nullhypothese auf. Folglich ist dann ein signifikantes Testergebnis „nicht tragisch“. Wenn bei einer *kleinen* Stichprobe ein Normalitätstest „anschlägt“, ist jedoch von einer relevanten Verletzung der Normalitätsannahme auszugehen.

Aufgrund der problematischen Verteilungsverhältnisse entscheiden wir uns, statt des geplanten parametrischen t-Tests für gepaarte Stichproben einen verteilungsfreien Lagevergleich mit dem **Vorzeichentest** durchzuführen (siehe z.B. Hartung 1989, S. 242f).

Dieser Test entscheidet sich zwischen folgenden Hypothesen:

$H_0$ : Der Median der Differenzvariablen AERGFZ ist kleiner oder gleich Null.

versus

$H_1$ : Die Differenzvariable AERGFZ hat einen positiven Median. (Mehr als 50% der Fälle haben einen positiven AERGFZ-Wert)

Statt der in Abschnitt 7.1 ausführlich vorgestellten Prüfstatistik  $T_Z$  verwendet der Vorzeichentest eine Prüfgröße, die im Wesentlichen auf der Anzahl der positiven AERGZ-Ausprägungen in der Stichprobe basiert. Sie wird üblicherweise mit  $Z$  bezeichnet, weil sie unter der  $H_0$  (genauer: bei einem Median von Null) approximativ  $z$ -verteilt (d.h. standardnormalverteilt) ist. Die Übereinstimmung der Bezeichnung mit der oben eingeführten Abkürzung für unsere Ärgerzuwachs-Variable ist also rein zufällig.

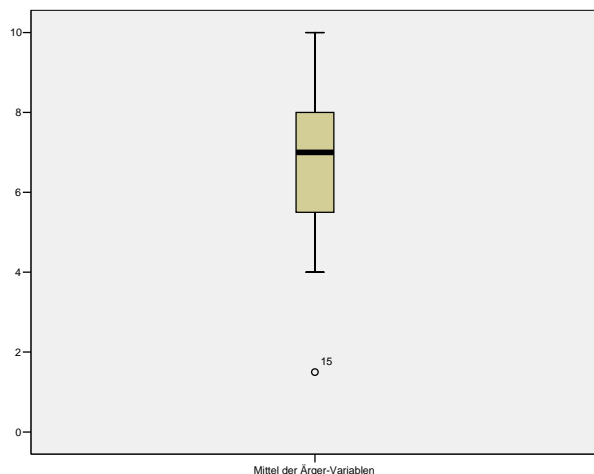
Man geht davon aus, dass die Verteilungs-Approximation ab  $n \geq 20$  hinreichend genau ist, so dass wir den Test bei unserer Stichprobe ( $n = 31$ ) in der üblichen approximativen Form anwenden dürfen. Bei kleineren Stichproben muss eine exakte Variante des Tests eingesetzt werden, die von SPSS ebenfalls unterstützt wird.

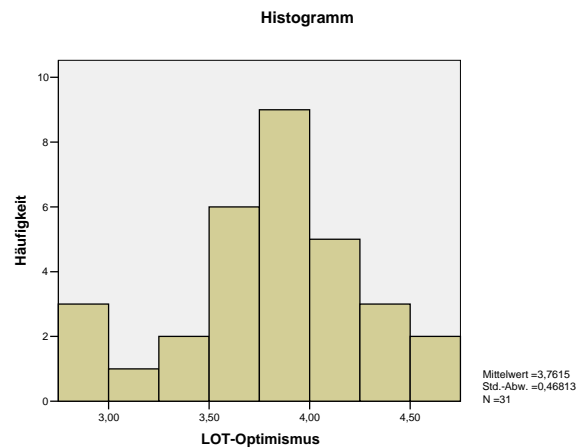
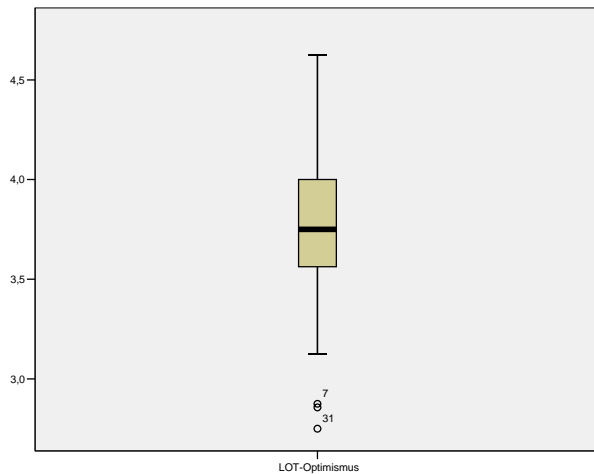
Weil der Vorzeichentest weit weniger empfindlich auf Ausreißer reagiert als der parametrische  $t$ -Test, können wir den kritischen Fall Nr. 4 in der Auswertung belassen. Damit vermeiden wir den Verdacht, die Daten zu unseren Gunsten bereinigt zu haben. Heben Sie also bitte die MD-Deklaration für den Wert  $-4$  bei der Variablen AERGZ wieder auf.

Die bisherige Diskussion der AERGZ-Verteilung hat sich auf Gefahrenquellen für die Interpretierbarkeit des geplanten zentralen Hypothesentests konzentriert. Es ist jedoch keinesfalls verboten, sondern sogar dringend empfohlen, sich anhand obiger Verteilungsdiagramme und sonstiger deskriptiver Informationen einen Eindruck von der empirischen Bewährung der KFA-Hypothese zu verschaffen. Das Histogramm spricht für einen starken KFA-Effekt in der erwarteten Richtung. Eine genaue Kenntnis des deskriptiven Ergebnisbilds kann verhindern, dass wir von einem durch technische Defekte verfälschten Testergebnis in die Irre geführt werden.

### 7.3.4 Ergebnisse für AERGAM und LOT

Bei den Variablen AERGAM und LOT finden sich keine Hinweise auf Fehler in den Berechnungsanweisungen oder auf extreme Ausreißer:





Die in den Boxplots auftauchenden Ausreißer sind nicht extrem (Abstand vom 25. Perzentil kleiner als 3 Boxlängen), und sollten aufgrund einer relativ kleinen Stichprobe, welche die Populationsverteilungen nur grob charakterisiert, *nicht* ausgeschlossen werden.

Bei der mit diesen Variablen geplanten Regressionsanalyse hat zudem die Ausreißeranalyse auf der Basis der Modellresiduen das weit größere Gewicht.

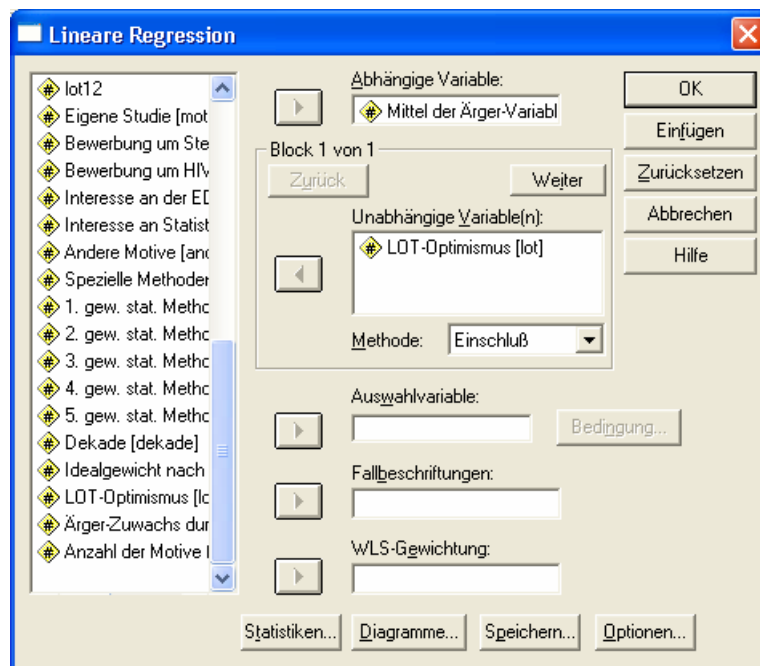
## 7.4 Prüfung der differentialpsychologischen Hypothese

### 7.4.1 Regression von AERGAM auf LOT

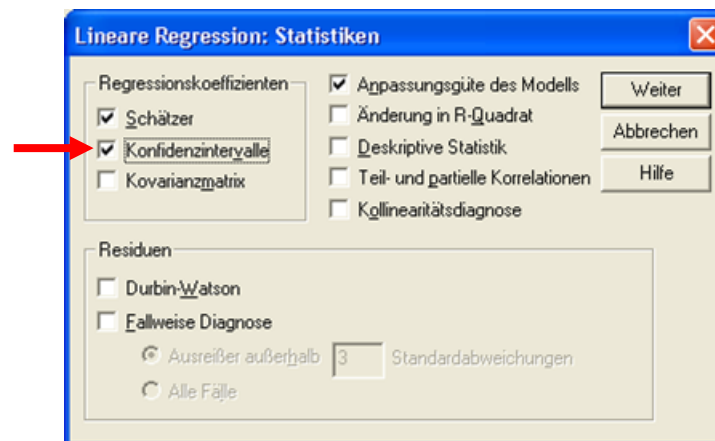
Nun wollen wir die lineare Regression von AERGAM auf LOT untersuchen, die wir nach dem Menübefehl

#### Analysieren > Regression > Linear

in der folgenden Dialogbox anfordern können:



Verlangen Sie in der Statistiken - Subdialogbox über die Voreinstellung hinausgehend noch die Berechnung von Konfidenzintervallen:



Wir erhalten zwar, wie erwartet, einen negativen Regressionskoeffizienten, doch ist dieser bei weitem nicht signifikant:

**Koeffizienten<sup>a</sup>**

Modell	Nicht standard. Koeffizienten		Standard. Koeff.	T	Signifikanz	95%-Konfid.intervall für B	
	B	Standardfehler	Beta			Untergrenze	Obergrenze
1 (Konstante)	7,669	2,947		2,602	,014	1,641	13,697
LOT-Optimismus	-,264	,778	-,063	-,339	,737	-1,854	1,327

a. Abhängige Variable: Mittel der Ärger-Variablen

SPSS ermittelt eine zweiseitige Überschreitungswahrscheinlichkeit von 0,737, die auch nach der zulässigen Halbierung aufgrund unserer einseitigen Fragestellung von der Signifikanzgrenze weit entfernt ist.

Der LOT-Optimismus zeigt entgegen unserer Annahme fast keinen *linearen* Zusammenhang mit dem summativen Ärger in unserer fiktiven Situation.

Von den bislang ungeprüften Voraussetzungen des Regressionsmodells ist in dieser Lage vor allem die Linearitätsannahme relevant. Vielleicht ist ein Zusammenhang vorhanden, aber nicht von linearer Form. Wir gehen gleich auf Versuche zur Rettung der differentialpsychologischen Hypothese ein.

Wer sich ausführlich über die Regressionsanalyse mit SPSS informieren möchte, kann eine elektronische Publikation des Rechenzentrums zu diesem Thema auf dem WWW-Server der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend folgendermaßen erreichen:

[Weitere Serviceangebote](#) > [EDV-Dokumentationen](#) > [Elektronische Publikationen](#) > [Statistische Spezialthemen](#) > [Lineare Regressionsanalyse mit SPSS](#)

Hier wird u.a. diskutiert, wie man die statistischen Voraussetzungen einer linearen Regressionsanalyse (Linearität, Unabhängigkeit, Normalität und Homoskedastizität der Residuen) überprüfen kann.

## 7.4.2 Methodologische Anmerkungen

### 7.4.2.1 Explorative Analysen im Anschluss an einen „gescheiterten“ Hypothesentest

Nach dem „Scheitern“ einer konfirmatorischen Forschungsbemühung wird sich in der Regel eine exploratorische Phase anschließen. Im Fall unserer differentialpsychologischen Hypothese sollten wir uns spätestens jetzt mit Hilfe eines Streudiagramms (siehe unten) einen Eindruck von der bivariaten Verteilung der beiden Variablen AERGAM und LOT verschaffen. Oben wurde schon zu Recht festgestellt, dass man (wegen potentieller technischer Probleme) einem statis-

tischen Test nur dann glauben sollte, wenn seine Entscheidung mit den deskriptiven Befunden harmoniert. Wir mussten bislang auf das Streuungsdiagramm verzichten, weil uns die dazu nötigen SPSS-Kenntnisse noch fehlen.

Außer dem Streuungsdiagramm kommen in unserem Beispiel auch noch andere statistische und graphische Analysen in Frage, um neue Information über empirische Gesetzmäßigkeiten zu gewinnen. Bei der explorativen Analyse der Stichprobendaten können Hypothesen generiert oder verbessert werden. Wir werden uns in Abschnitt 8 z.B. dafür interessieren, ob eventuell das Geschlecht den Zusammenhang zwischen Optimismus und Ärger moderiert. Allerdings ist es *nicht* möglich, die neuen oder revidierten Hypothesen anhand *derselben* Stichprobe zu testen. Also: Sie dürfen und sollen aus Ihren Daten etwas lernen, aber ein Test der dabei generierten Hypothese erfordert eine neue, unabhängige Stichprobe.

Außerdem sollten Sie es nicht unterlassen, das „Scheitern“ einer Hypothese zu veröffentlichen. Ansonsten tragen Sie dazu bei, in der Fachliteratur ein systematisch verzerrtes Bild der Wirklichkeit aufzubauen.

#### 7.4.2.2 *Post hoc* - Poweranalyse

Bei der Interpretation des obigen Resultates ist außerdem zu beachten, dass die Power des t-Tests zum Regressionskoeffizienten in unserer relativ kleinen Stichprobe recht bescheiden ist, so dass kleine Effekte leicht übersehen werden können. Unser Testergebnis kann nicht als Beleg für die Nullhypothese interpretiert werden, doch spricht es wohl gegen die Existenz eines *starken* Effektes. Um zu genaueren Aussagen zu kommen, betrachten wir die Power unseres t-Tests bei unterschiedlichen Effektstärken in der Population, quantifiziert durch die Korrelation  $\rho$ .

Dabei verwenden wir erneut das Programm **GPowEr 3** (Faul et al., im Druck), das schon bei der Stichprobenumfangsplanung in Abschnitt 1.3.2 zum Einsatz kam.<sup>1</sup>

Auf den Pool-PCs der Universität Trier unter dem Betriebssystem MS-Windows ist GPowEr 3 folgendermaßen zu starten

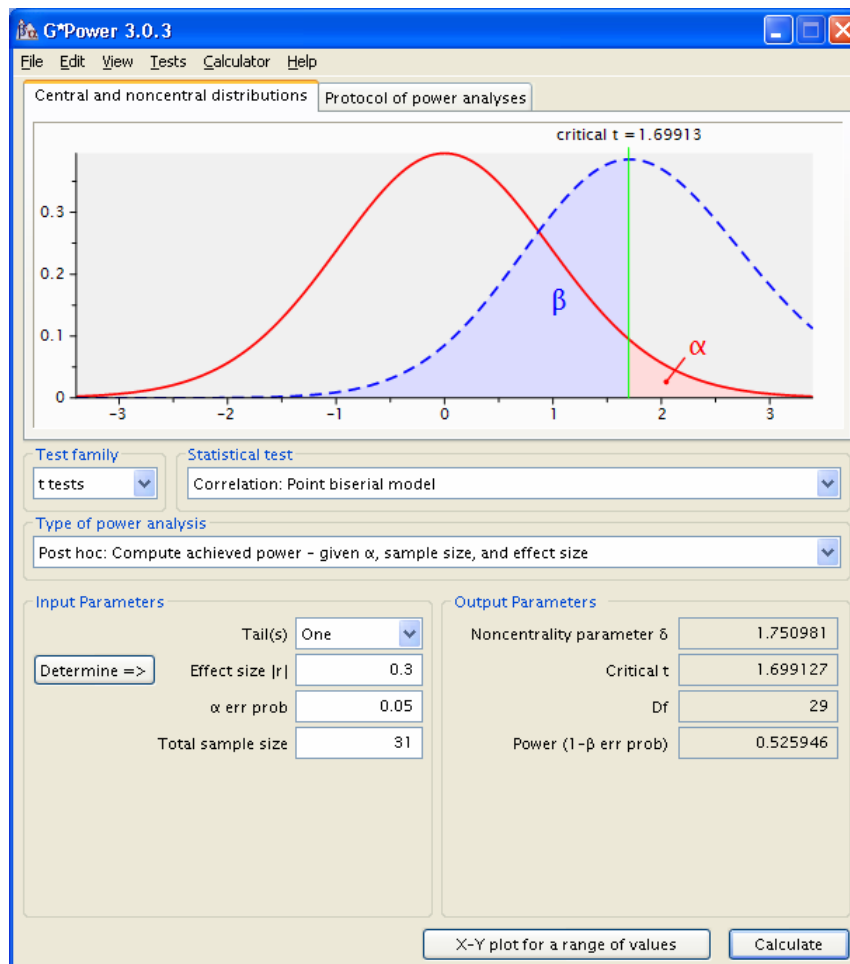
#### **Start > Programme > Wissenschaftliche Programme > GPowEr**

Wir wählen

- |                                       |  |
|---------------------------------------|--|
| • <b>Test family</b>                  | <b>t-Tests</b>                           |
| • <b>Statistical test</b>             | <b>Correlation: Point biserial model</b> |
| • <b>Type of power analysis</b>       | <b>Post hoc</b>                          |
| • <b>Effect size  r </b>              | 0.3                                      |
| • <b>Tail(s)</b>                      | <b>One</b>                               |
| • <b><math>\alpha</math> err prob</b> | 0.05                                     |
| • <b>total sample size</b>            | 31                                       |

<sup>1</sup> GPowEr 3 kann für MS-Windows und MacOS kostenlos über folgende Webseite bezogen werden:

<http://www.psych.uni-duesseldorf.de/abteilungen/aap/gpower3/>

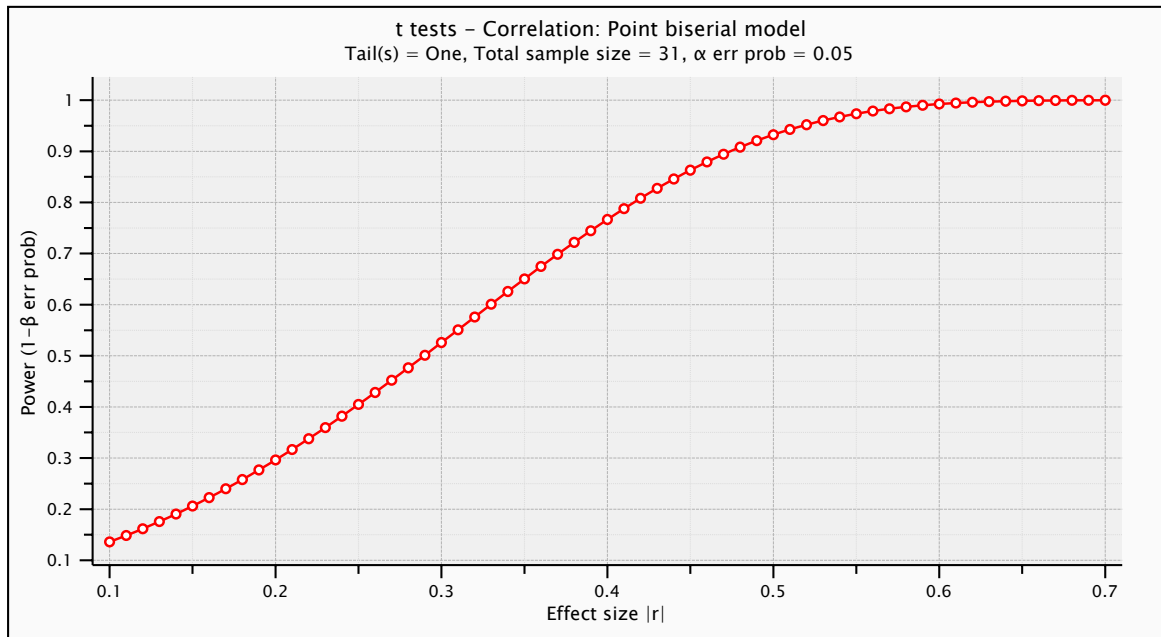


Nach einem Klick auf den Schalter **Calculate** wird eine Power von lediglich 0,53 für den Test der differentialpsychologischen Hypothese berechnet.

Um zur Darstellung der Power als Funktion der Effektstärke zu gelangen, klicken wir auf den Schalter **X-Y plot for a range of values** und wählen

- **Plot (on y axis)**                      **Power (1 - β err prob)**
- **as a function of**                      **Effect size |r|**
- **from**                                      0,1
- **in steps of**                              0.01
- **through to**                              0,7
- **Plot**                                      1

Nach einem Klick auf den Schalter **Draw Plot** zeigt die folgende Abbildung zeigt, wie bei fester Stichprobengröße ( $n = 31$ ) die Power des einseitigen t-Tests von der Effektstärke, d.h. vom Betrag der Korrelation abhängt:



Erst ab einer Effektstärke von ca.  $|r| = 0,5$  ist die Power so groß (ca. 0,95), dass man die ausgebliebene Signifikanz als Beleg gegen einen Effekt dieser Stärke werten kann. Unserer Studie hat also keinesfalls die differentialpsychologische Nullhypothese bewiesen, aber doch ein Argument gegen die Existenz eines starken Effektes ( $|r| \geq 0,5$ ) geliefert.

#### 7.4.2.3 Paarweiser oder fallweiser Ausschluss fehlender Werte

Wir müssen uns leider wieder einmal mit dem Problem fehlender Werte befassen, wenn auch ohne direkten Bezug zum Demoprojekt: Wenn Sie die Korrelationsmatrix zu gewissen Variablen A, B, C und D anfordern, dann kann SPSS fehlende Werte auf zweierlei Weise berücksichtigen:

- **Paarweiser Ausschluss fehlender Werte**

Zur Berechnung der Korrelation zwischen den Variablen A und B werden alle Fälle herangezogen, die *bei diesen beiden* Variablen einen validen Wert haben.

Vorteil: Alle verfügbaren validen Beobachtungen werden ausgenutzt.

Nachteil: In der entstehenden Korrelationsmatrix beruhen die einzelnen Koeffizienten im Allgemeinen auf unterschiedlichen Teilstichproben. Daher fehlt dieser Matrix eine gewisse mathematische Eigenschaft (die positive Semidefinitheit), die bei normalen Korrelationsmatrizen vorhanden ist und die in vielen Statistikprozeduren vorausgesetzt wird. Es kann dadurch (z.B. in einer multiplen Regressionsanalyse) zu artifiziellen Ergebnissen kommen.

- **Fallweiser Ausschluss fehlender Werte**

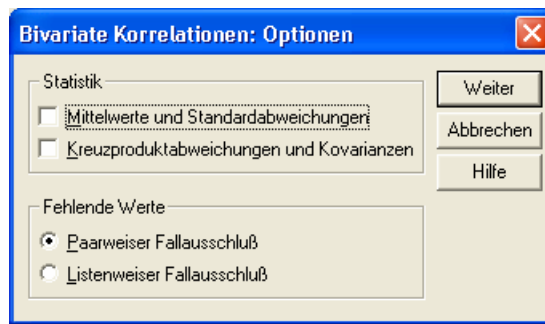
Zur Berechnung der Korrelation zwischen den Variablen A und B werden nur Fälle herangezogen, die *bei allen Variablen*, also bei A, B, C und D, einen validen Wert haben.

Vorteil: Die entstehende Korrelationsmatrix ist intakt (positiv semidefinit).

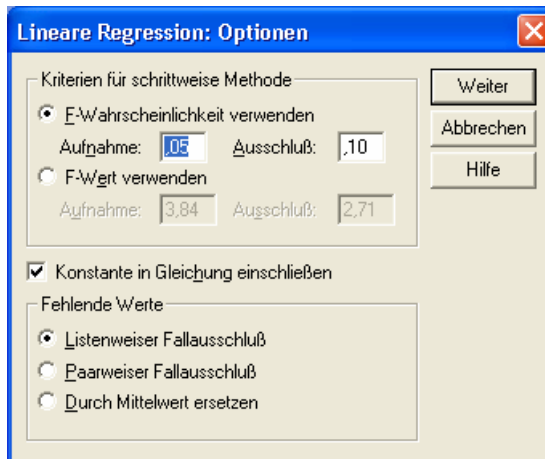
Nachteil: Ist die Gesamtmenge der beteiligten Variablen groß, gehen eventuell sehr viele Fälle verloren.

Per Voreinstellung benutzt SPSS bei der Korrelationsberechnung die *paarweise* Methode. Mit dem Schalter **Optionen** in der Dialogbox **Bivariate Korrelationen** erhalten Sie folgende Subdialogbox, die ein Umschalten auf die fallweise Methode erlaubt:





Bei der linearen Regressionsanalyse benutzt SPSS die Voreinstellung *fallweise* und bietet außerdem an, fehlende Werte durch den Mittelwert der jeweiligen Variablen zu ersetzen:



Über den Menübefehl

### **Analysieren > Analyse fehlender Werte**

kann man die permanente Ersetzung fehlender Werte per EM- oder Regressionsalgorithmus veranlassen.

## **7.5 Prüfung der KFA-Hypothese**

Nun wollen wir die allgemeinspsychologische Kernhypothese unserer Studie prüfen, dass die Verfügbarkeit kontrafaktischer (also positiver) Alternativen den Ärger über ein ungünstiges Ereignis steigert. Aufgrund der Ausreißer- und Verteilungsanalyse in Abschnitt 7.3.3 haben wir uns entschieden, statt des ursprünglich geplanten (parametrischen) t-Tests für abhängige Stichproben den verteilungsfreien Vorzeichenstest zu verwenden.

Suchen Sie die zuständige Dialogbox zunächst über das **Analysieren**-Menü. Bei Misserfolg können Sie auch den Index des Hilfesystems benutzen. Steigen Sie ein mit:

### **Hilfe > Themen > Index**

und beginnen Sie dann, in das aktive Textfeld *Vorzeichenstest* zu schreiben. Schon nach dem vierten Buchstaben wird der gesuchte Beitrag aufgelistet und ist per Doppelklick auf seinen Titel zu öffnen. Hier ist u.a. der Weg zur benötigten Dialogbox erklärt:

### **Analysieren > Nichtparametrische Tests > Zwei verbundene Stichproben**

In der Dialogbox müssen Sie die beiden Variablen angeben und den gewünschten Test markieren:



Wir erhalten folgendes Ergebnis:

#### Häufigkeiten

		N
Ärger mit kontrafaktischer Alternative - Ärger ohne kontrafaktische Alternative	Negative Differenzen <sup>a</sup>	2
	Positive Differenzen <sup>b</sup>	26
	Bindungen <sup>c</sup>	3
	Gesamt	31

- Ärger mit kontrafaktischer Alternative < Ärger ohne kontrafaktische Alternative
- Ärger mit kontrafaktischer Alternative > Ärger ohne kontrafaktische Alternative
- Ärger ohne kontrafaktische Alternative = Ärger mit kontrafaktischer Alternative

#### Statistik für Test<sup>a</sup>

	Ärger mit kontrafaktischer Alternative - Ärger ohne kontrafaktische Alternative
Z	-4,347
Asymptotische Signifikanz (2-seitig)	,000

a. Vorzeichentest

Selbst die von SPSS ausgegebene zweiseitige Überschreitungswahrscheinlichkeit (Bezeichnung: **Asymptotische Signifikanz (2-seitig)**) ist deutlich kleiner als unser vorgegebenes  $\alpha$ -Niveau (0,05). Das unserer einseitigen Fragestellung entsprechende *einseitige* p-level ergibt sich (wegen der Symmetrie der zugrunde liegenden Prüfverteilung) durch Halbierung des zweiseitigen p-levels, ist also erst recht kleiner als die kritische Grenze 0,05.

Damit kann die KFA-Nullhypothese (*Kein Ärgerzuwachs durch eine kontrafaktische Alternative*) deutlich zurückgewiesen werden.

Nach Klärung der zentralen Hypothesen ist unser Projekt nun eigentlich abgeschlossen, aber es gibt noch viele SPSS-Optionen kennen zu lernen, und unsere Daten enthalten sicher auch noch einige interessante Details.

## 7.6 Übung

Für die Differenzvariable (GEWICHT - IDGEW) akzeptieren beide Normalverteilungstests die Nullhypothese:

### Tests auf Normalverteilung

	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistik	df	Signifikanz	Statistik	df	Signifikanz
GEWICHT - IDGEW	,092	31	,200*	,984	31	,905

\*. Dies ist eine untere Grenze der echten Signifikanz.

a. Signifikanzkorrektur nach Lilliefors

Führen Sie mit den Variablen GEWICHT und IDGEW einen t-Test für gepaarte Stichproben zu folgendem Testproblem durch:

$H_0$ : Das Realgewicht der Trierer Studierenden liegt im Mittel nicht unter dem Idealgewicht nach der Formel „Größe - 100“.

versus

$H_1$ : Die Trierer Studierenden sind in Relation zur Idealgewichtsformel „Größe - 100“ im Mittel zu leicht.

Die Ergebnisse werden im nächsten Abschnitt wiedergegeben.

## 7.7 Arbeiten mit dem Ausgabefenster (Teil III)

Oben wurde gelegentlich in didaktischer Nachlässigkeit ohne Erläuterung der Begriff *Pivot-Tabelle* verwendet. Unter dem *Pivotieren* einer Tabelle versteht SPSS u.a. die folgenden Operationen:

- Austauschen ihrer Zeilen-, Spalten- und Schichtendimensionen
- Änderung der Schachtelungsordnung
- Kategorien ausblenden

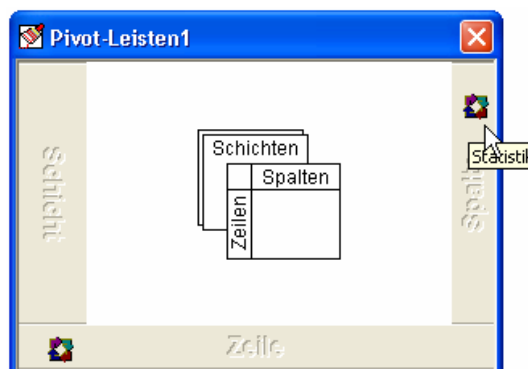
Nachdem wir den Pivot-Editor im zweiten Teil der Serie *Arbeiten mit dem Ausgabefenster* bereits für konventionelle Tabellengestaltungen benutzt haben, beschäftigen wir uns nun mit den Leistungen, die seinen Namen begründen.


### 7.7.1 Pivot-Editor starten

Man startet den Pivot-Editor zum **Bearbeiten** einer Tabelle per Doppelklick oder über das Kontextmenü. Es empfiehlt sich, anschließend nötigenfalls mit dem Menübefehl

#### Pivot > Pivot-Leisten

das folgende Fenster einzuschalten:



Es enthält je eine Leiste für die Zeilen, Spalten und Schichten der Tabelle und je ein Pivotsymbol  für die dargestellten Tabellendimensionen. Die Zeilenleiste enthält z.B. die Pivotsymbole zu allen in den Zeilen dargestellten Tabellendimensionen. Welche Dimension ein Symbol repräsentiert, erfährt man per Quickinfo-Text, wenn man den Mauszeiger einige Zeit darauf ruhen lässt.

Wir wollen als Beispiel die in obiger Übung von Ihnen erstellte Tabelle mit dem t-Test zum Vergleich von Real- und Idealgewicht betrachten:

		Gepaarte Differenzen					T	df	Sig. (2-seitig)
		Mittelwert	Standardabweichung	Standardfehler des Mittelwertes	95% Konfidenzintervall der Differenz				
					Untere	Obere			
Paaren 1	Körpergewicht (in kg) - Idealgewicht nach der Formel: Größe - 100	-9,3226	6,1881	1,1114	-11,5924	-7,0528	-8,388	30	,000

Diese Tabelle enthält leider nur eine Schicht, so dass wir den Umgang mit Mehrschichttabellen nicht üben können.

In den Zeilen der Tabelle wird die Dimension **Paare** dargestellt. Da wir nur ein einziges Variablenpaar untersucht haben, hat diese Dimension nur eine Kategorie, deren Beschriftung aus den Labels der beiden Variablen abgeleitet wurde.

Die Spaltendimension **Statistik** sorgt mit ihren zahlreichen Kategorien für eine überbreite Tabelle, die schlecht auf ein DIN-A4-Blatt im Hochformat passt.

### 7.7.2 Dimensionen verschieben

Durch das Verschieben ihres Pivotsymbols kann man für eine Dimension neu festlegen, ob ihre Kategorien durch Spalten, Zeilen oder Schichten dargestellt werden sollen. Wenn in unserem Beispiel die beiden Pivotsymbole ihre Plätze tauschen, benötigt die Tabelle in horizontaler Richtung deutlich weniger Platz:

		Paaren 1
		Körpergewicht (in kg) - Idealgewicht nach der Formel: Größe - 100
Gepaarte Differenzen	Mittelwert	-9,3226
	Standardabweichung	6,1881
	Standardfehler des Mittelwertes	1,1114
95% Konfidenzintervall der Differenz	Untere	-11,5924
	Obere	-7,0528
T		-8,388
df		30
Sig. (2-seitig)		,000

### 7.7.3 Gruppierungen

Man kann mehrere Kategorien einer Dimension zusammenfassen und mit einem Gruppentickett kennzeichnen. In der aktuellen Version unserer Beispieltabelle sind z.B. die **Untere** und die

**Obere** Konfidenzschranke gruppiert mit dem Etikett **95% Konfidenzintervall der Differenz**. Beseitigen Sie bitte diese Gruppierung folgendermaßen:

- Rechtsklick auf das Kategorienetikett
- Aus dem Kontextmenü wählen: **Gruppierung aufheben**

Welche Gruppierungen in einer Tabelle vorhanden sind, erkennt man am besten nach dem Einschalten der Gitterlinien mit

### Ansicht > Gitterlinien

In unserem Beispiel zeigt sich bei der Statistikdimension eine weitere Gruppe mit dem Etikett **Gepaarte Differenzen**:<sup>1</sup>

Test bei gepaarten Stichproben

		Paaren 1
		Körpergewicht (in kg) - Idealgewicht nach der Formel: Größe - 100
Gepaarte Differenzen	Mittelwert	-9,3226
	Standardabweichung	6,1881
	Standardfehler des	1,1114
	Untere	-11,5924
	Obere	-7,0528
T		-8,388
df		30
Sig. (2-seitig)		,000

Beseitigen Sie bitte der Übersichtlichkeit halber auch diese Gruppierung.

Wenn Sie schließlich noch bei der **Paare**-Dimension das Gruppenetikett **Paaren 1** entfernen, erhalten Sie folgendes Zwischenergebnis:

Test bei gepaarten Stichproben

	Körpergewicht (in kg) - Idealgewicht nach der Formel: Größe - 100
Mittelwert	-9,3226
Standardabweichung	6,1881
Standardfehler des	1,1114
Untere	-11,5924
Obere	-7,0528
T	-8,388
df	30
Sig. (2-seitig)	,000

Wenn Sie mehrere Kategorien einer Dimension zu einer Gruppe zusammenfassen wollen, können Sie folgendermaßen vorgehen:

<sup>1</sup> Eingblendete Gitterlinien sind nur bei aktivem Pivot-Editor sichtbar. Um diese Hilfslinien im Manuskript darzustellen, wurden über **Format > Tabelleneigenschaften > Rahmen** zusätzliche Trennlinien aktiviert (und später wieder abgeschaltet).

- Alle Kategorien markieren
- Kontextmenü zu einer markierten Kategorie öffnen und Option **Gruppieren** wählen
- Gruppenbeschriftung anpassen

In der folgenden Version unserer Tabelle wurde eine Gruppe mit den drei Kategorien zum t-Test gebildet:

**Test bei gepaarten Stichproben**

		Körpergewicht (in kg) - Idealgewicht nach der Formel: Größe - 100
Mittelwert		-9,32
Standardabweichung		6,19
Standardfehler des Mittelwertes		1,11
Untere		-11,59
Obere		-7,05
	T	-8,39
Signifikanztest	df	30
	Sig. (2-seitig)	,00

Außerdem wurde bei einigen Zellen die Anzahl der Dezimalstellen reduziert (über **Format > Zelleigenschaften**).

#### 7.7.4 Kategorien aus- und einblenden

Wenn eine SPSS-Tabelle zu ausführlich erscheint, können Kategorien einer Dimension ausgeblendet werden. In unserem Beispiel wollen wir bei der Statistikdimension auf den Standardfehler des Mittelwertes und die Konfidenzintervalle verzichten:

**Test bei gepaarten Stichproben**

		Körpergewicht (in kg) - Idealgewicht nach der Formel: Größe - 100
Mittelwert		-9,32
Standardabweichung		6,19
	T	-8,39
Signifikanztest	df	30
	Sig. (2-seitig)	,000

Gehen Sie beim Ausblenden einer Kategorie folgendermaßen vor:

- Bei gedrückter Tastenkombination **Strg+Alt** einen (linken) Mausklick auf das Kategorienetikett setzen
- Rechtsklick auf das Kategorienetikett
- Aus dem Kontextmenü wählen: **Kategorie ausblenden**

In *Spalten* untergebrachte Kategorien kann man auch auf intuitive Weise eliminieren:

- linker Mausklick auf den rechten Spaltenrand, Maustaste gedrückt halten
- Spaltenbreite durch Verschieben der Maus reduzieren, bis die Quick-Info **Ausblenden** erscheint:

Gepaarte Differenzen			
Standardabweichung	Standardfehler des Mittelwerts	95% Konfidenzintervall der Differenz	
		Ausblenden	Obere
6,188	1,111	-11,592	-7,053

- Maustaste loslassen

Zum *Einblenden* von vorher abgeschalteten Kategorien kenne ich nur die global wirksame Methode:

**Ansicht > Alles einblenden**

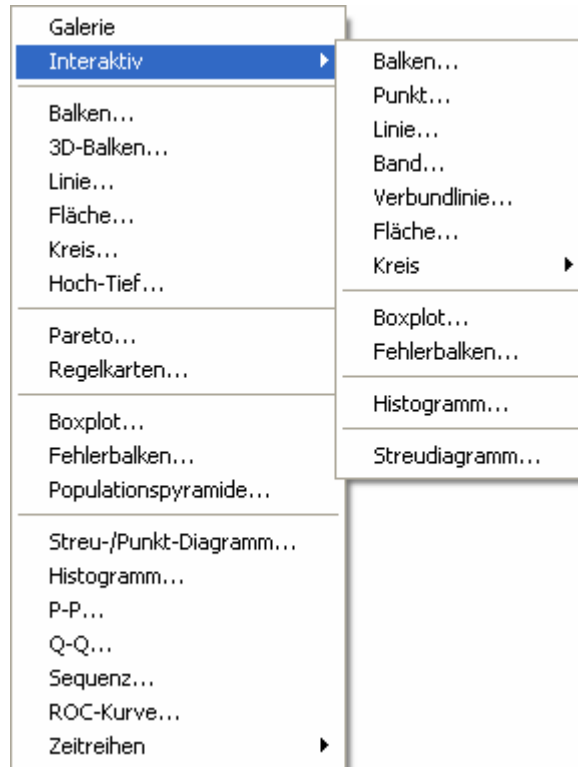
Nach diesem Befehl können Tabellenbestandteile auftauchen (z.B. Dimensionsbeschriftungen), die (je nach verwendeter Vorlage) bei *neuen* Tabellen nicht eingeschaltet sind.

---

## 8 Grafische Datenanalyse

Wir haben schon einige grafische Darstellungsmöglichkeiten kennen gelernt, die im Rahmen von Statistikprozeduren angeboten werden (z.B. Histogramm, Boxplot). In diesem Abschnitt arbeiten wir erstmals mit dem **Grafiken**-Menü und vor allem mit dem Editor zur individuellen Nachbearbeitung von Diagrammen.

SPSS-Einsteiger werden vermutlich durch das **Grafiken**-Menü leicht irritiert, weil viele Grafiktypen sowohl auf der Hauptebene als auch im Untermenü **Interaktiv** auftauchen:



Ursache ist die Koexistenz der Standardgrafik (verknüpft mit dem SPSS-Kommando GRAPH) mit der so genannten interaktiven Grafik (verknüpft mit dem Kommando IGRAPH). War über einige SPSS-Versionen hinweg die interaktive Grafik moderner und leistungsstärker, ist seit der SPSS-Version 12 die Standardgrafik deutlich variabler und attraktiver. Wenn sich eine spezielle Darstellung mit der Standardgrafik nicht zufrieden stellend realisieren lässt, ist das alternative Grafiksystem aber auf jeden Fall einen Versuch wert.

Von den zahlreich angebotenen Grafiktypen können aus Zeitgründen nur wenige Beispiele behandelt werden. Im aktuellen Abschnitt 8 wird die einfache Variante des Streudiagramms vorgestellt, in Abschnitt 10.1 kommt ein Balkendiagramm zum Einsatz.

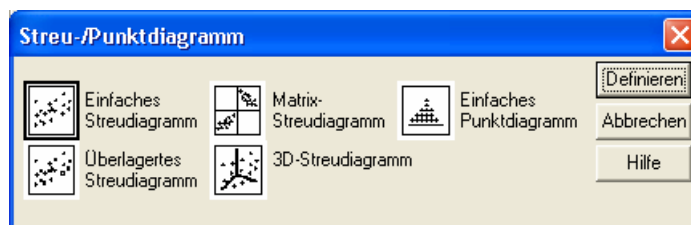
### 8.1 Streudiagramm anfordern

Um die empirische Regression von Gewicht auf Größe visuell beurteilen zu können, fordern wir über den folgenden Menübefehl ein Streudiagramm mit den beiden Variablen an:

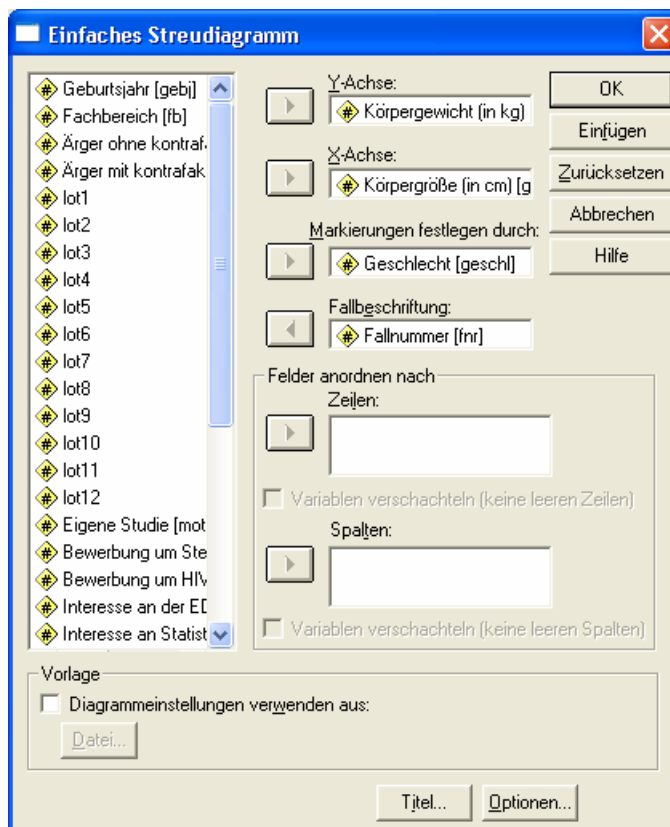
#### **Grafiken > Streu-/Punkt-Diagramm**

In der nun erscheinenden Palette akzeptieren wir die voreingestellte **einfache** Diagrammvariante



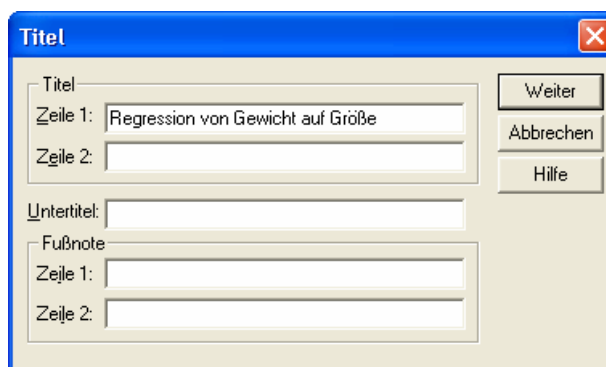


und wechseln per Mausklick auf den Schalter **Definieren** zur Dialogbox **Einfaches Streudiagramm**, wo die beteiligten Variablen per Transportschalter  ihre Rollen erhalten:



Durch die Verwendung von GESCHL als **Markierungsvariable** werden weibliche und männliche Datenpunkte verschieden dargestellt, so dass geschlechtsbedingte Unterschiede bei der Regression von Gewicht auf Größe ggf. sichtbar werden.

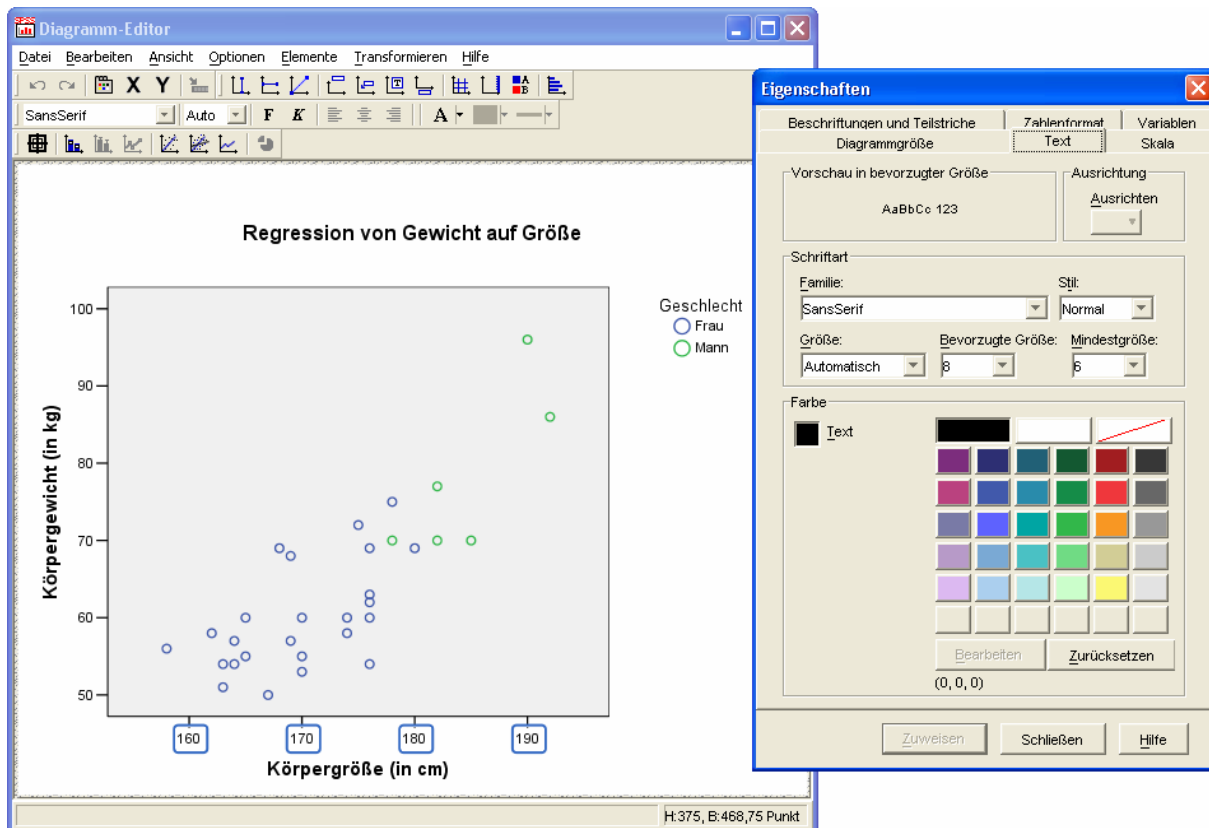
Nach einem Mausklick auf den Schalter **Titel** tragen wir eine Titelzeile ein:



Quittieren Sie die Subdialogbox mit **Weiter** und die Hauptdialogbox mit **OK**.



## 8.2 Streudiagramm modifizieren

Wenn Sie im Viewer einen Doppelklick auf die fertige Grafik setzen, wird sie im Diagramm-Editor geöffnet:




Im Unterschied zum Pivot-Editor für Tabellen, der seine Tätigkeit per Voreinstellung an Ort und Stelle entfaltet, öffnet der Grafikeditor stets ein eigenes Fenster.

Anschließend werden am Beispiel des Streudiagramms einige allgemeine Bedienungsmöglichkeiten des Diagramm-Editors vorgestellt.

Deren Effekte lassen über die Schalter   (mehrstufig) rückgängig machen bzw. wiederherstellen.

### Eigenschaftsfenster

Zum aktuell im Grafikeditor markierten Objekt bzw. zur markierten Objektgruppe (erkennbar an einer blauen Umrahmung) bietet das Eigenschaftsfenster (siehe oben) auf jeweils dynamisch erstellten Registerkarten alle modifizierbaren Attribute. Bei Bedarf kann es über den Schalter , die Tastenkombination **Strg+T** oder den Menübefehl

#### Bearbeiten > Eigenschaften

aktiviert werden.

Wer im Beispiel X-Achsenteilstrichwerte im Abstand von 5 cm wünscht, kann so vorgehen:

- X-Achsenteilstrichwerte per Mausklick auf einen Wert markieren
- im Eigenschaftsfenster die Registerkarte **Skala** wählen
- bei der **ersten Unterteilung** die **Auto**-Markierung aufheben
- den **benutzerdefinierten** Wert 5 eintragen
- **Zuweisen**, um das Ergebnis sofort inspizieren zu können

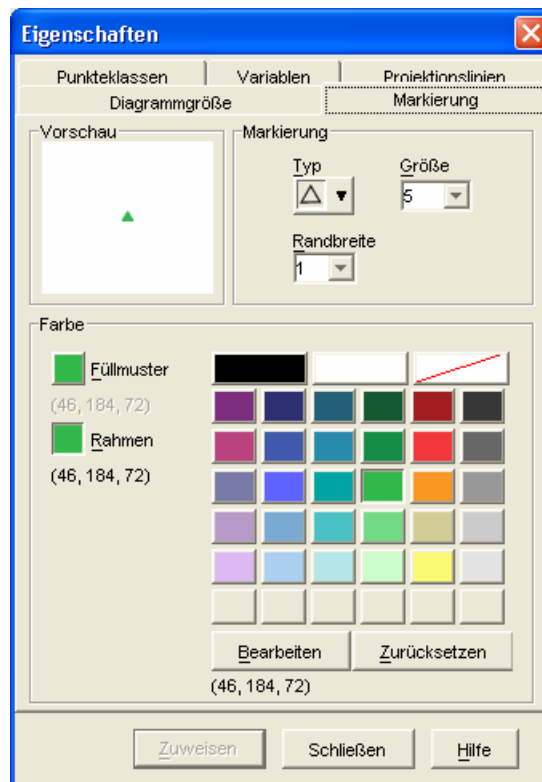
## Markieren von gruppierten Objekten

Sind gruppierte Objekte vorhanden (z.B. die Datenpunkte für Frauen bzw. Männer in unserem Streudiagramm), dann wendet SPSS beim Markieren folgende Logik an:

- Ist gerade *kein* Objekt markiert, bewirkt ein Mausklick auf ein beliebiges Objekt aus einer beliebigen Gruppe die Markierung aller Objekte (aus sämtlichen Gruppen).
- Ein weiterer Mausklick schränkt die Markierung auf die getroffene Gruppe ein. Durch einen Mausklick auf ein Objekt einer anderen Gruppe wird diese komplett markiert.
- Ein weiterer Mausklick in derselben Gruppe schränkt die Markierung auf das getroffene Objekt ein.
- Sobald ein einzelnes Objekt markiert ist, wandert bei weiteren Mausklicks die Einzelmarkierung über Gruppengrenzen hinweg zum getroffenen Objekt.
- Bei gedrückter **Strg**-Taste ist ein gruppenunabhängiges kumulierendes Markieren möglich.


Eine alternative Möglichkeit zum Markieren aller Elemente einer Gruppe ist der Mausklick auf das zugehörige Symbol in der Legende.

Im Beispiel könnte man nach dem Markieren der Datenpunkte zu jeweils einer Teilstichprobe die Form, Größe, Randfarbe und Füllfarbe der Symbole ändern:



## Menüs

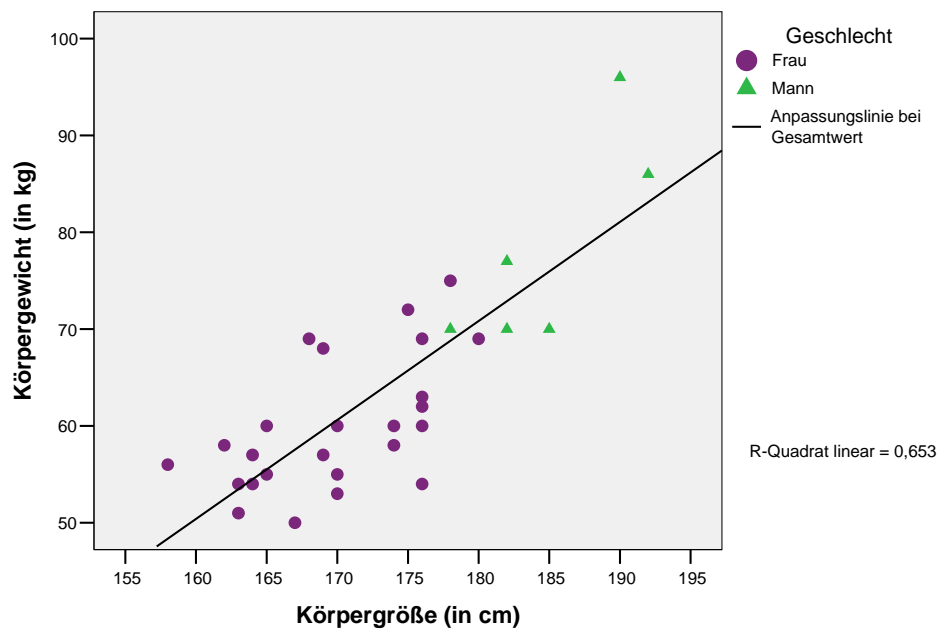
Viele Angebote sind über die Untermenüs zu den Items **Optionen** und **Elemente** im Grafiker-Hauptmenü sowie über äquivalente Symbolleisten verfügbar (z.B. Anpassungs- oder Interpunktionslinien, Datenbeschriftungen, Legende, Anmerkungen). Zumindest bei Streudiagrammen sind die Kontextmenüs zu den meisten Objekten sehr ähnlich aufgebaut und fassen die Angebote der Hauptmenüitems **Optionen** und **Elemente** zusammen.


Im Beispiel bietet es sich an, über das Symbol  oder den Menübefehl

### Elemente > Anpassungslinie bei Gesamtwert

die empirische Regressionsgerade einzeichnen zu lassen:

## Regression von Gewicht auf Größe

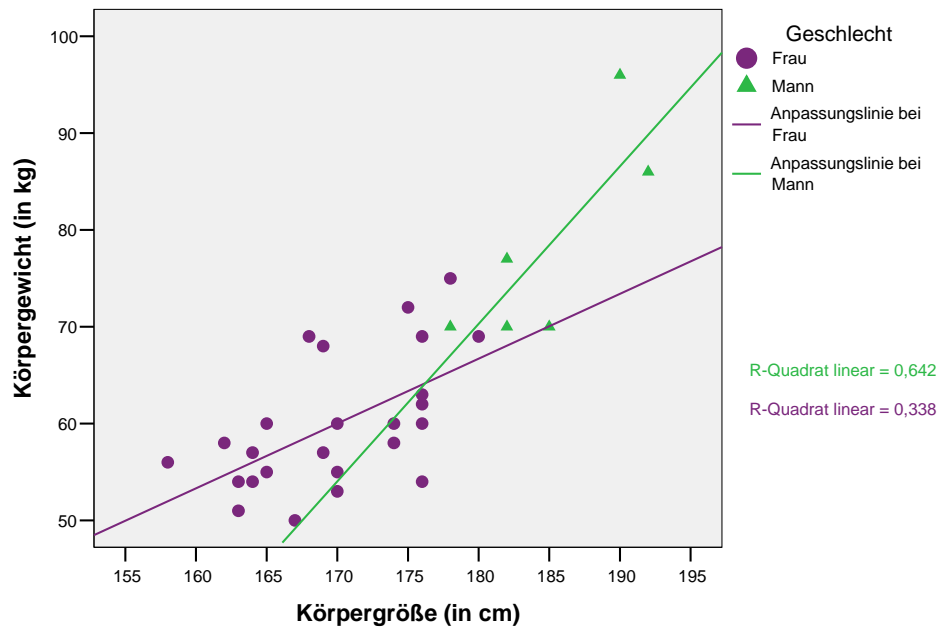


Überflüssige Objekte lassen sich über ihr Kontextmenü oder (im markierten Zustand) per **Entf**-Taste löschen. Im Beispiel könnte man so die Regressionsgerade wieder verschwinden lassen, um anschließend über das Symbol  oder den Menübefehl

## Elemente &gt; Anpassungslinie bei Untergruppen

gruppenspezifische (geschlechtsbedingte) Regressionsgeraden einzufügen:

## Regression von Gewicht auf Größe



Man erkennt in der Grafik einen Geschlechtsunterschied hinsichtlich der Regressionssteigung, der durch Unterschiede im Körperbau zu erklären ist:




Bei zwei Männern mit 10 cm Größenunterschied ist ein stärkerer Gewichtsunterschied zu erwarten als bei zwei Frauen mit derselben Größendifferenz. Es ist also zu vermuten, dass Geschlecht den Effekt der Größe auf das Gewicht *moderiert*.

### Beschriftungen

Viele **Beschriftungen** (z.B. Überschriften, Legenden, Erläuterungen) besitzen nach dem Markieren einen Textrahmen mit acht Anfassern zur Größenänderung:





Solche Rahmen lassen sich auch verschieben, wobei die Transportfunktionalität des Mauszeigers am Rand aktiv wird, signalisiert durch die Zeigergestalt .

Um einen Text zu ändern, markiert man ihn und setzt nach Erscheinen des Rahmens einen weiteren Mausklick darauf. Zum Beenden der Texteingabe drückt man die **Enter**-Taste oder setzt einen Mausklick außerhalb des Textrahmens.

Bei der Textformatierung kann alternativ zum Eigenschaftsfenster auch die folgende Symbolleiste verwendet werden:



Über die Schaltfläche  (de)aktiviert man das Werkzeug  zur Datenbeschriftung, das zu angeklickten Datenpunkten den Wert der vereinbarten Fallbeschriftungsvariablen oder aber die laufende Nummer in die Grafik einfügt bzw. wieder entfernt, z.B.:



Nach einem rechten Mausklick auf einen Datenpunkt mit dem Fallbeschriftungswerkzeug kann man per Kontextmenü veranlassen, dass die zugehörige Zeile im Datenfenster markiert wird.

### 8.3 Grafiken verwenden

Wie Tabellen lassen sich auch die Grafiken aus dem Viewer-Fenster über die Windows-Zwischenablage in andere Anwendungen übertragen:

- Mit **Bearbeiten > Kopieren** oder **Strg+C** überträgt man eine markierte Grafik vom Viewer in die Zwischenablage.
- Mit **Bearbeiten > Einfügen** oder **Strg+V** übernimmt man sie in ein Dokument der Zielanwendung.

Als Viewer-Bestandteile lassen sich Grafiken sichern, drucken oder exportieren.

Zur Verwendung als **Vorlage** kann man eine Grafik aus dem Diagramm-Editor mit dem Menübefehl

#### Datei > Diagrammvorlage speichern

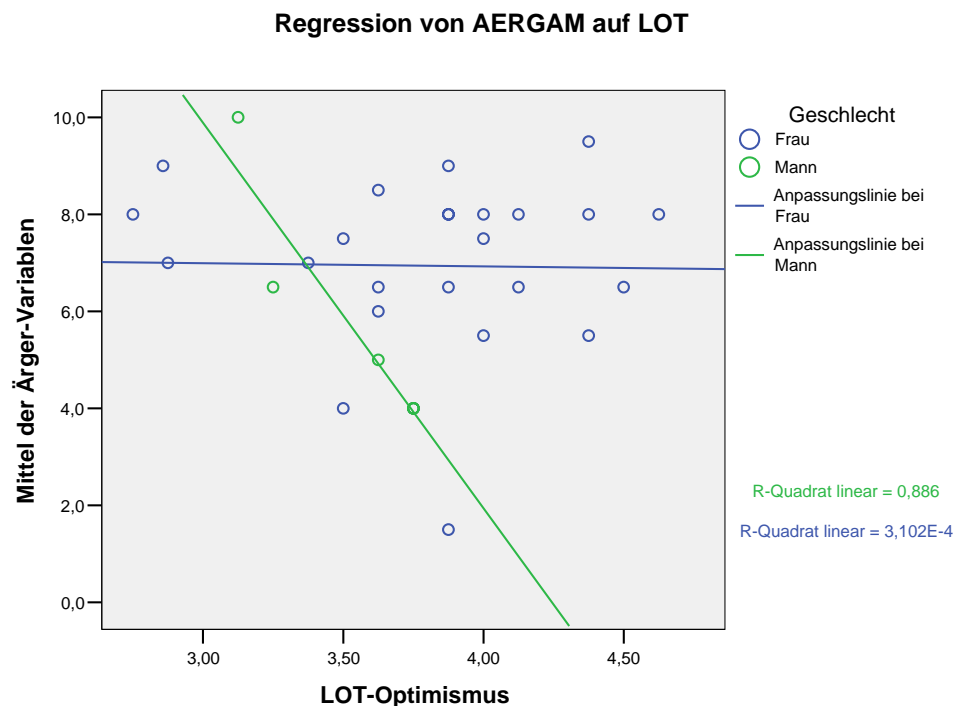
in eine Datei mit der Namenserweiterung **sgt** sichern.

Auf andere Grafiken kann man eine Vorlage bereits beim Erstellen (siehe Dialogbox **Einfaches Streudiagramm** in Abschnitt 8.1) oder im Diagrammeditor anwenden:

#### Datei > Diagrammvorlage zuweisen

### 8.4 Übung

Um Fehlentscheidungen aufgrund von technischen Fehlern zu vermeiden, sollten wir uns zu jedem statistischen Test die zugrunde liegenden deskriptiven Datenverhältnisse möglichst genau ansehen. Dies muss für die „gescheiterte“ differentialpsychologische Hypothese (siehe Abschnitt 7.4) noch nachgeholt werden. Erzeugen Sie bitte dazu ein Streudiagramm mit den Variablen AERGAM und LOT, und verwenden Sie wie in obigem Beispiel GESCHL als Markierungsvariable. Mit eingezeichneten Regressionsgeraden für die Untergruppen sollten Sie ungefähr folgendes Ergebnis erhalten:



Während bei den Frauen offenbar *kein* Zusammenhang zwischen LOT und AERGAM besteht, zeigt sich bei den Männern ein Effekt im Sinne unserer differentialpsychologischen Hypothese. Allerdings sollten wir die Beobachtung sehr zurückhaltend interpretieren, weil unsere Stichprobe lediglich sechs Männer enthält.

Immerhin resultiert bei einer regressionsanalytischen Auswertung für den Moderatoreffekt<sup>1</sup> eine relativ kleine Überschreitungswahrscheinlichkeit (0,01):

Koeffizienten<sup>a</sup>

Modell		Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten	T	Signifikanz
		B	Standardfehler	Beta		
1	(Konstante)	-19,356	11,285		-1,715	,098
	GESCHL * LOT	-7,883	2,860	-5,633	-2,756	,010
	Geschlecht	26,543	10,211	5,426	2,600	,015
	LOT-Optimismus	7,818	3,121	1,863	2,505	,019

a. Abhängige Variable: Mittel der Ärger-Variablen

Hier haben wir es aber **nicht** mit dem signifikanten Ergebnis eines statistischen Tests zu tun, sondern mit einem deskriptiven Maß zu einer interessanten Vermutung, die sich bei der explorativen Datenanalyse ergeben hat. Eine Testentscheidung über die Moderatorhypothese ist nur in einer unabhängigen Stichprobe möglich.

<sup>1</sup> Über die Analyse von Moderatoreffekten mit Hilfe der SPSS-Regressions-Prozedur informiert eine elektronische Publikation des Rechenzentrums, die Sie auf dem WWW-Server der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend folgendermaßen erreichen:

[Weitere Serviceangebote](#) > [EDV-Dokumentationen](#) > [Elektronische Publikationen](#) > [Statistische Spezialthemen](#) > [Moderatoranalyse per multipler Regression mit SPSS](#)

## 9 Fälle auswählen

Es kommt durchaus vor, dass man sich bei einer Analyse auf eine Teilstichprobe beschränken möchte.

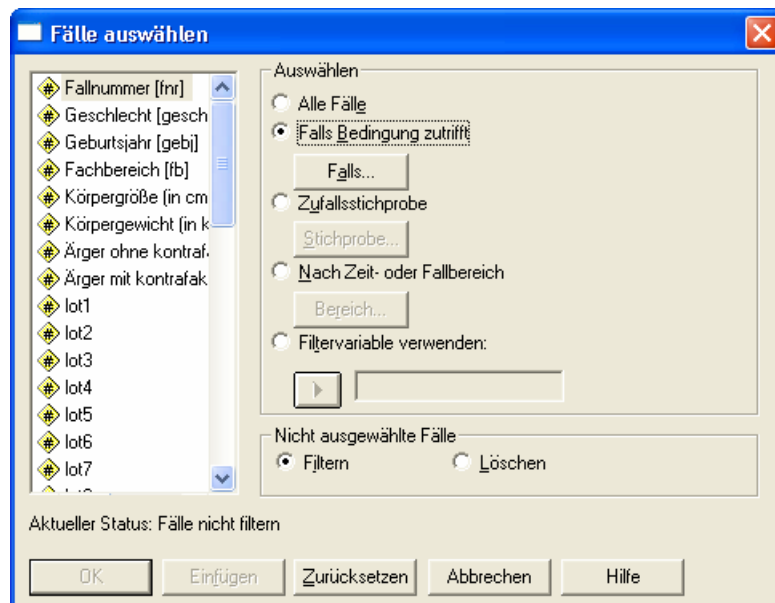
Bei unserer KFA-Studie ist es von Interesse, die Personen mit einem *negativen* KFA-Effekt ( $AERGZ < 0$ ) näher kennen zu lernen. Wir können dazu nach geeigneter Fallauswahl einen Bericht mit interessanten Variablenausprägungen anfordern.

### 9.1 Auswahl über eine Bedingung

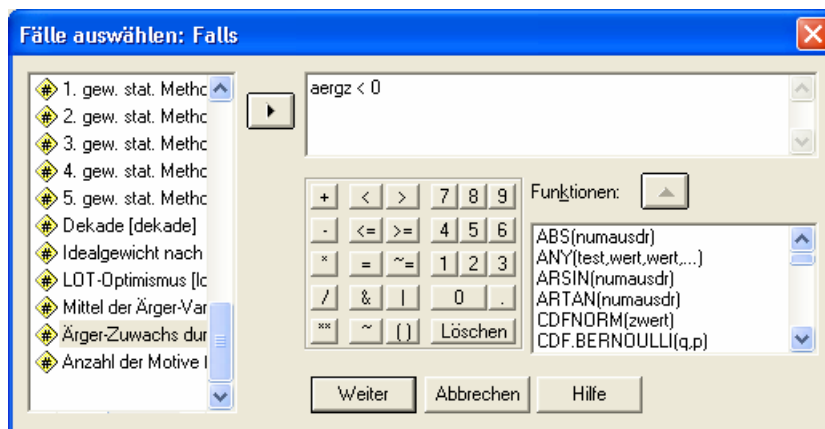
SPSS erlaubt es, Fälle in Abhängigkeit von einer Bedingung temporär oder permanent aus der Arbeitsdatei auszuschließen. Die zuständige Dialogbox erreichen Sie über den Menübefehl:

#### Daten > Fälle auswählen

Um eine Bedingung für die Teilnahme an den weiteren Auswertungen zu setzen, müssen Sie im Optionenfeld **Auswählen** die Alternative **Falls Bedingung zutrifft** markieren und anschließend die zugehörige Subdialogbox mit dem **Falls**-Schalter aktivieren:



Im **Falls**-Dialogfenster haben Sie die Möglichkeit, einen beliebigen logischen Ausdruck (vgl. Abschnitt 6.5.2) als Teilnahme Kriterium zu definieren, z.B.:



Wenn Sie nach erfolgreicher Definition des Teilnahme Kriteriums **Weiter** machen, können Sie im Optionenfeld **Nicht ausgewählte Fälle** der Hauptdialogbox (siehe oben) entscheiden, was mit den Negativ-Fällen geschehen soll:



- **Filtern** SPSS erzeugt aufgrund Ihres logischen Ausdrucks eine Hilfsvariable namens FILTER\_\$ mit folgenden Werten:
  - 1 falls bei einem Fall der logische Ausdruck wahr ist,
  - 0 sonst (also auch bei unbestimmtem Ausdruck).

Diese Variable wird als **Filter** aktiviert, d.h. bis zu einer Deaktivierung des Filters werden bei allen Analysen nur noch Fälle mit Wert 1 bei FILTER\_\$ einbezogen. Die in den einstweiligen Ruhezustand versetzten Null-Fälle sind im Datenfenster an der durchgestrichenen Zeilennummer zu erkennen:

	meth1	meth2	meth3	meth4	meth5	dekade	idgew	lot	aergam	aergz	polymot	filter \$
1	1	2	3	0	0	1	63	4,13	6,5	3	1	0
2	1	2	0	0	0	2	58	3,88	6,5	3	1	0
3	4	0	0	0	0	1	74	3,63	6,0	4	1	0
4	1	2	5	0	0	1	82	3,75	4,0	-4	1	1
5	3	2	4	0	0	1	80	3,88	8,0	0	1	0

**Wichtig:** Filter wirken sich nur bei statistischen und graphischen Analysen aus. Bei Datentransformationen werden auch die ausgefilterten Fälle einbezogen. Wer eine bedingte Datentransformation benötigt, muss die Methoden aus Abschnitt 6.5 verwenden.

Wenn ein Filter aktiv ist, wird dies in der Statuszeile angezeigt (siehe Abbildung). Um den Filter später zu deaktivieren, müssen Sie die Dialogbox **Fälle auswählen** erneut aufrufen und dann im **Auswählen**-Optionenfeld wieder den Ausgangszustand **Alle Fälle** aktivieren.

- **Löschen** Die Negativ-Fälle werden aus der (temporären) Arbeitsdatei entfernt. Aus der *externen* Datei (z.B. auf der Festplatte) verschwinden die Fälle dabei *nicht*. Wenn Sie allerdings das teilentleerte Datenfenster „sichern“, haben Sie eventuell anschließend ein kleines Problem.

#### Hinweise:

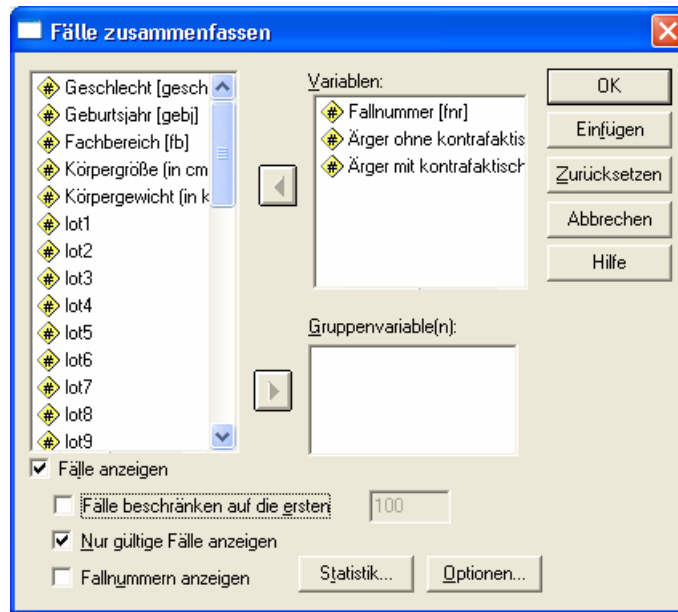
- Ist beim Sichern der Arbeitsdatei ein Filter aktiv, wird die zugrunde liegende Variable FILTER\_\$ mit abgespeichert. Beim nächsten Öffnen der Datei ist der Filter allerdings *nicht* aktiv, sondern muss nötigenfalls erneut vereinbart werden. Dazu muss aber lediglich die Variable FILTER\_\$ in der Dialogbox **Fälle auswählen** als **Filtervariable verwendet** werden. Weil Filtervariablen mit beliebigem Namen akzeptiert werden, kann man in einer SPSS-Datendatei mehrere Filtervariablen bereithalten. Außerdem kann man die einem Filter zugrunde liegende Syntax abspeichern und später wieder verwenden.
- Mit der oben dargestellten Dialogbox **Fälle auswählen** kann man auch eine *zufällige Teilstichprobe* ziehen oder eine Analyse auf die ersten *n* Fälle beschränken.

## 9.2 Bericht anfordern

Gelegentlich benötigt man für eine bestimmte Teilmenge von Fällen eine übersichtliche Liste mit den Ausprägungen bestimmter Variablen. Um z.B. für Personen mit negativem Ärgerzuwachs eine Liste mit den Variablen FNR, AERGO und AERGM zu erhalten, vereinbart man zunächst die Filterbedingung „AERGZ < 0“ und fordert dann über

## Analysieren > Berichte > Fälle zusammenfassen

die gewünschte Auflistung an:



Wir erhalten folgende Liste:

**Zusammenfassung von Fällen**

	Fallnummer	Ärger ohne kontrafaktische Alternative	Ärger mit kontrafaktischer Alternative
1	4	6	2
2	15	2	1
Insgesamt N	2	2	2

---

## 10 Analyse von Kreuztabellen

Wir wollen die Hypothese prüfen, dass Frauen und Männer unterschiedliche Präferenzen bei der Wahl des Studienfachs haben.

Unsere Fachbereichs-Variable (FB) enthält Information über die Studienfächer der Untersuchungsteilnehmer(innen) auf einem angemessenen Aggregationsniveau. Ihre Werte stehen für die folgenden Fachbereiche der Universität Trier:

Fachbereich	Fächer
I	Pädagogik, Philosophie, Psychologie
II	Sprachorientierte Fächer
III	Historische und politische Wissenschaften
IV	BWL, Ethnologie, Informatik, Mathematik, Soziologie, VWL, Wirtsch.-Informatik
V	Jura
VI	Geowissenschaften

Nachdem die Begriffe aus der eingangs formulierten inhaltlichen Hypothese hinreichend präzisiert sind, können wir die empirisch zu prüfenden *Nullhypothese* formulieren:

**Die Merkmale Geschlecht und Fachbereich sind unabhängig voneinander.**

Die Unabhängigkeitsbehauptung der Nullhypothese bedeutet, dass sich aus dem Wissen über das Geschlecht eines Untersuchungsteilnehmers keinerlei Information über seine Fachbereichszugehörigkeit ableiten lässt, dass also die bedingten Fachbereichsverteilungen bei Frauen und Männern identisch sind.

Zur Illustration des *Unabhängigkeitsbegriffs* wurde hier auf eine Verteilungshomogenität verwiesen. Später folgen noch einige Erläuterungen zu den beiden Begriffen und zu ihrer Beziehung.

Unsere Nullhypotheseformulierung ist „zweiseitig“, wozu es auch gar keine Alternative gibt, weil die Fachbereichsvariable mehr als zwei Stufen hat. Bei  $(2 \times 2)$ -Kreuztabellen sind aber auch einseitige Hypothesen möglich (siehe Abschnitt 10.3.3.2).

Weil der Zusammenhang zwischen den beiden *nominalskalierten* Merkmalen Fachbereich und Geschlecht zu untersuchen ist, wählen wir als Auswertungsmethode die Kreuztabellenanalyse mit  $\chi^2$ -Test.

Weil Kreuztabellenanalysen recht häufig benötigt werden, erläutert der vorliegende Abschnitt die wichtigsten statistischen Grundlagen und die Regeln für eine korrekte Interpretation der SPSS-Ergebnisse.

Leider erweist sich unsere KFA-Stichprobe bei näherer Betrachtung als ungeeignet zur Prüfung der Präferenz-Divergenz-Hypothese, denn

- Sie ist recht klein (geringe Teststärke).
- Die Stichprobe ist wenig repräsentativ, weil nur SPSS-Interessierte enthalten sind. Folglich sind manche Fachbereiche (z.B. III, V) fast nicht vertreten.

Daher wurde eine Zufallsstichprobe der Größe  $n = 283$  aus der Datenbank mit allen Studierenden der Universität Trier im WS 1993/94 gezogen<sup>1</sup>. Bei jedem Fall wurden die Variablen Geschlecht (GESCHL) und Fachbereich (FB) festgestellt.

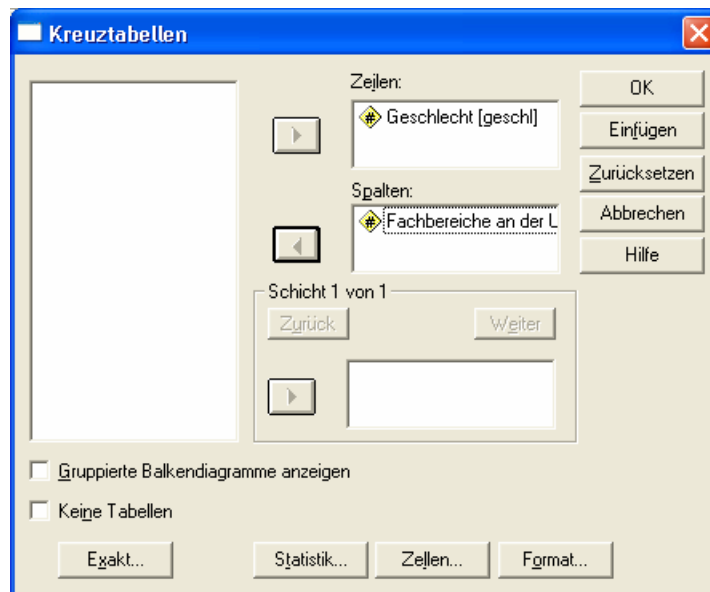
Die SPSS-Datendatei **fbgeschl.sav** mit den beiden Variablen finden Sie an der im Vorwort für Kursdateien vereinbarten Stelle.

### 10.1 Beschreibung der bivariaten Häufigkeitsverteilung

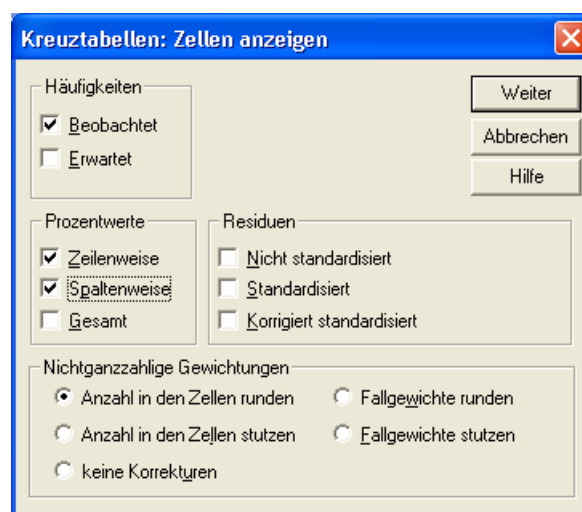
Die SPSS-Dialogbox zur Analyse zweidimensionaler Kontingenztabelle erscheint nach dem Menübefehl:

#### Analysieren > Deskriptive Statistiken > Kreuztabellen

Wir wählen GESCHL als Zeilen- und FB als Spaltenvariable:



In der **Zellen**-Subdialogbox kann man u.a. zeilen- und spaltenbezogene Prozentangaben für die Zellen der Kontingenztabelle anfordern:



Aufgrund dieser Spezifikationen erhalten wir für unsere Stichprobe die folgende Kreuztabelle<sup>1</sup>:

<sup>1</sup> Aufmerksame Leser(innen) werden zu Recht fragen, warum nicht *alle* Trierer Studierenden einbezogen wurden. Eine größere Stichprobe bringt stabilere Ergebnisse und hätte in dieser speziellen Situation kaum mehr „gekostet“. Allerdings habe ich aus didaktischen Gründen eine Stichprobe mit „typischem“ Umfang vorgezogen.

Geschlecht \* Fachbereiche an der Universität Trier Kreuztabelle

		Fachbereiche an der Universität Trier						Gesamt
		I	II	III	IV	V	VI	
Frauen	Anzahl	29	26	18	22	26	23	144
	% von Geschlecht	20,1%	18,1%	12,5%	15,3%	18,1%	16,0%	100,0%
	% von FB	63,0%	66,7%	50,0%	31,0%	54,2%	53,5%	50,9%
Männer	Anzahl	17	13	18	49	22	20	139
	% von Geschlecht	12,2%	9,4%	12,9%	35,3%	15,8%	14,4%	100,0%
	% von FB	37,0%	33,3%	50,0%	69,0%	45,8%	46,5%	49,1%
Gesamt	Anzahl	46	39	36	71	48	43	283
	% von Geschlecht	16,3%	13,8%	12,7%	25,1%	17,0%	15,2%	100,0%
	% von FB	100,0%	100,0%	100,0%	100,0%	100,0%	100,0%	100,0%

Durch die Einträge in den Zellen wird die gemeinsame Verteilung der beiden Variablen GESCHL und FB beschrieben:

- Oben ... steht die absolute Häufigkeit der Zelle  
Z.B. befanden sich in der Stichprobe 17 Studenten aus dem Fachbereich I.
- In der Mitte ... steht der prozentuale Anteil der Zelle an allen Fällen in der zugehörigen Zeile.  
Z.B. gehörten von den 139 *männlichen* Untersuchungsteilnehmern 12,2% zum Fachbereich I.  
Diese auf die Zeile bezogenen relativen Häufigkeiten beschreiben also die bedingte Verteilung der Spaltenvariablen (FB) für einen festen Wert der Zeilenvariablen (GESCHL). Wir erhalten z.B. für die Männer die folgende bedingte Verteilung der Fachbereichs-Variablen:

I	II	III	IV	V	VI
12,2%	9,4%	12,9%	35,3%	15,8%	14,4%

- Unten ... steht der prozentuale Anteil der Zelle an allen Fällen in der zugehörigen Spalte  
Z.B. waren von den 46 Personen aus dem Fachbereich I 37% Männer.  
Diese auf die Spalte bezogenen relativen Häufigkeiten beschreiben also die bedingte Verteilung der Zeilenvariablen (GESCHL) für einen festen Wert der Spaltenvariablen (FB). Wir erhalten z.B. für den Fachbereich I die folgende bedingte Geschlechtsverteilung:

Frauen	63%
Männer	37%

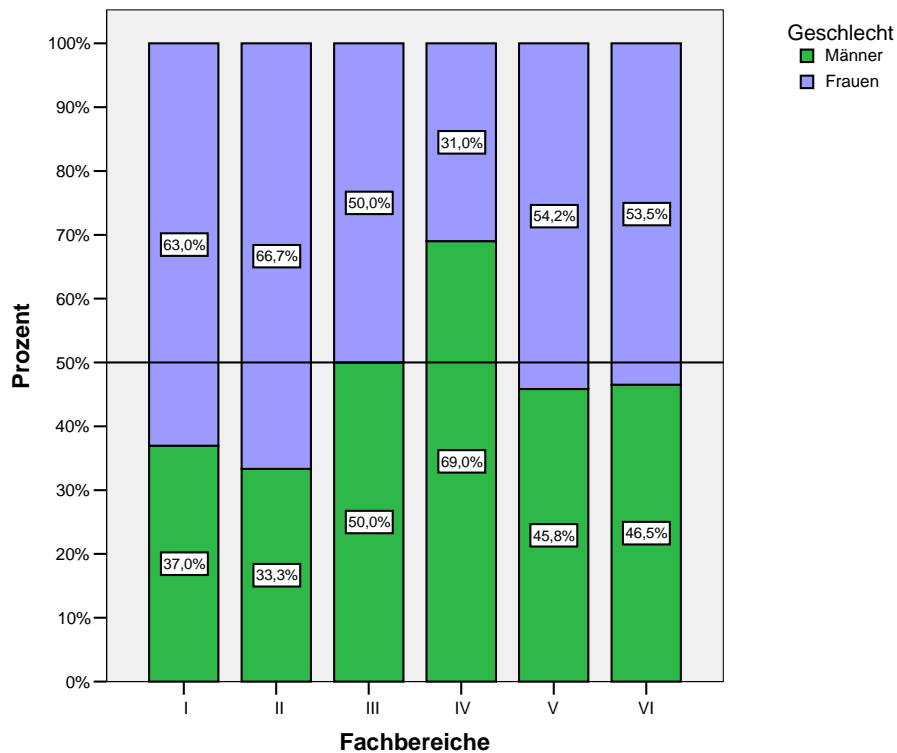
In der **Zellen**-Subdialogbox können auch noch weitere Informationen zu den Zellen angefordert werden (z.B. der prozentuale Anteil der Zelle an der Gesamtstichprobe).

Beim Vergleich der fachbereichsbedingten Geschlechtsverteilungen zeigen sich erhebliche Unterschiede:

- In den Fachbereichen I und II dominieren die Frauen mit einem Anteil von 63 bzw. 66,7%.
- Im Fachbereich IV sind die Frauen mit einem Anteil von nur 31% in der Minderheit.
- In den übrigen Fachbereichen III, V und VI zeigt sich ein relativ ausgeglichenes Geschlechtsverhältnis.

<sup>1</sup> Die Tabelle wurde mit dem Pivot-Editor durch Aufheben der Gruppierung **Geschlecht** etwas schlanker gemacht.

In diesem gestapelten Balkendiagramm werden die bedingten Verteilungen veranschaulicht:



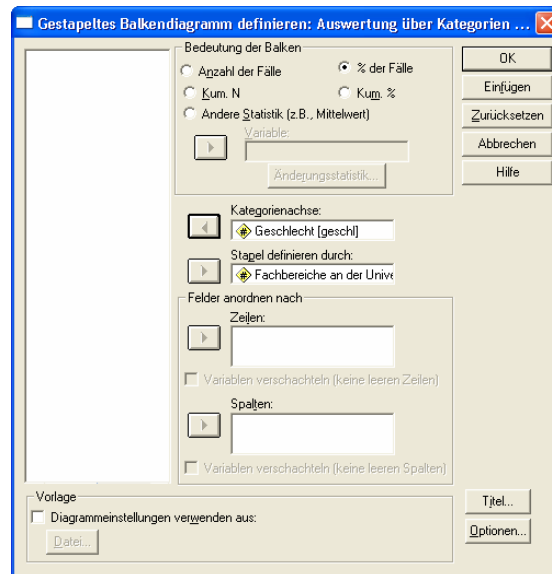
Sie können es nach dem Menübefehl

### Grafiken > Balken

und der Entscheidung für ein **gestapeltes** Balkendiagramm mit den **Kategorien einer Variablen** als **Daten im Diagramm**



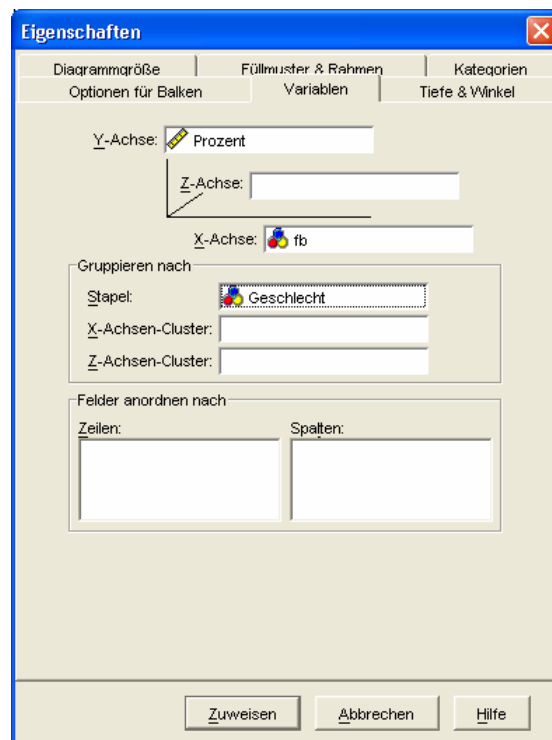
mit folgender Dialogbox anfordern:



Machen Sie **% der Fälle** zur **Bedeutung der Balken**. Indem zunächst GESCHL als Kategorienvariable fungiert, erzielt man den gewünschten Bezug für die Prozentangaben auf den Balken.

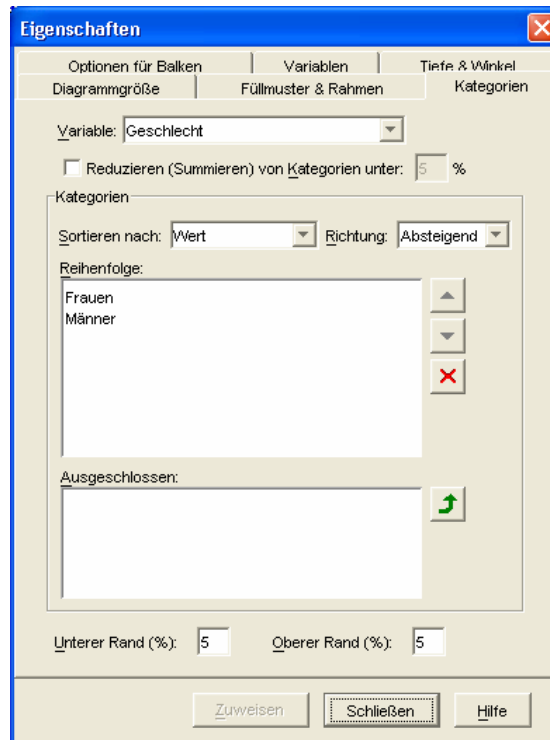
Nehmen Sie im Grafikeditor folgende Anpassungen vor:

- Bei markierten Balken tauschen GESCHL und FB ihre Rollen:

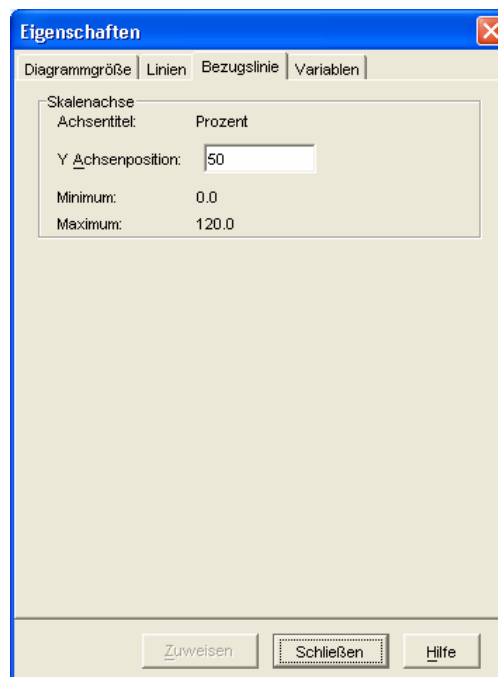


Befördern Sie z.B. die Variable FB per „Mauskralle“ (Ziehen und Ablegen) an ihren neuen Einsatzort.

- Die Reihenfolge der Kategorien wird geändert, z.B. durch die Wahl der **absteigenden Richtung**:

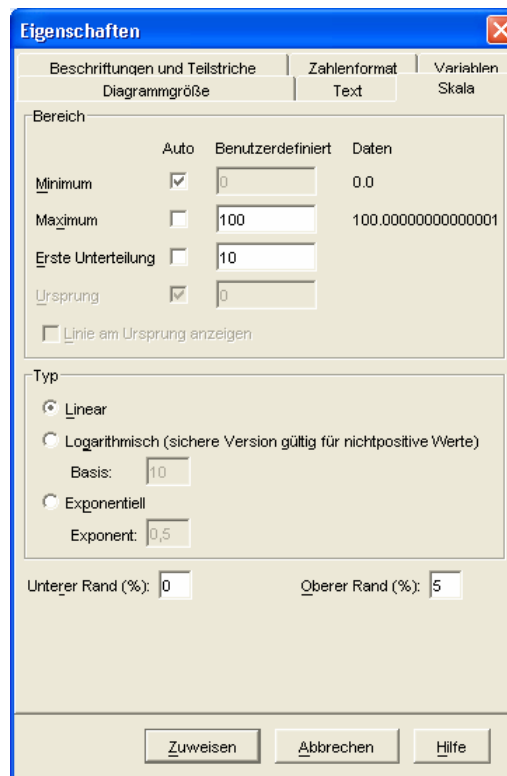


- Über **Optionen > Bezugslinie für Y-Achse** wird die 50% - Marke hervorgehoben:



- Nach dem Markieren der Y-Achse wird auf dem Eigenschaftsfenster-Registerblatt **Skala** das **Maximum** 100 und die **erste Unterteilung** 10 eingestellt:





- Über **Elemente > Datenbeschriftungen einblenden** sorgen wir für eine Anzeige der Prozentwerte.

## 10.2 Die Unabhängigkeits- bzw. Homogenitätshypothese

Hypothesen zum Zusammenhang zwischen zwei kategorialen Merkmalen lassen sich auf letztlich äquivalente Weise durch Verwendung verschiedener wahrscheinlichkeitstheoretischer Begriffen formulieren. Dies soll an unserem Beispiel demonstriert werden, damit Sie die Äquivalenz verstehen und ausnutzen lernen. Es ist ja generell sinnvoll, einen Sachverhalt aus verschiedenen Blickrichtungen zu betrachten.

### 1. Formulierung: Unabhängigkeitshypothese

- $H_0$ : Die Merkmale Geschlecht und Fachbereich sind unabhängig, d.h. die Wahrscheinlichkeit für jedes Verbundereignis (z.B. Mann im Fachbereich V) ist gleich dem Produkt aus den Wahrscheinlichkeiten der Randereignisse (im Beispiel: Mann, Fachbereich V).
- $H_1$ : Die Merkmale Geschlecht und Fachbereich sind abhängig, d.h. die Wahrscheinlichkeit für mindestens ein Verbundereignis ist ungleich dem Produkt aus den Wahrscheinlichkeiten der Randereignisse.

### 2. Formulierung: Homogenitätshypothese

- $H_0$ : Der Frauenanteil ist in allen Fachbereichen gleich.
- $H_1$ : Die Frauenanteile in den Fachbereichen sind verschieden.

Man kann leicht zeigen (vgl. Hartung 1989, S. 412):

*Perfekte Homogenität liegt genau dann vor, wenn die Merkmale Geschlecht und Fachbereich unabhängig sind.*

### 10.3 Testverfahren

#### 10.3.1 Asymptotische $\chi^2$ - Tests

Die bekannteste Prüfgröße zur Testung der Unabhängigkeits- bzw. Homogenitätshypothese ist die folgende  $\chi^2_{\text{P}}$  - Statistik nach Pearson:

$$\chi^2_{\text{P}} := \sum_{i=1}^z \sum_{j=1}^s \frac{(n_{ij} - m_{ij})^2}{m_{ij}}, \quad \text{mit } m_{ij} = \frac{n_{i.} \cdot n_{.j}}{n}$$

Darin bedeuten:

$z, s$	=	Anzahl der Zeilen bzw. Spalten
$n_{ij}$	=	beobachtete Häufigkeit in Zelle $ij$
$m_{ij}$	=	geschätzte erwartete Häufigkeit in Zelle $ij$ unter der $H_0$
$n_{i.}$	=	beobachtete Häufigkeit in Zeile $i$
$n_{.j}$	=	beobachtete Häufigkeit in Spalte $j$
$n$	=	Umfang der Gesamtstichprobe

Wir wollen kurz überlegen, wie die erwarteten Häufigkeiten  $m_{ij}$  unter der Nullhypothese geschätzt werden. Zunächst soll die Wahrscheinlichkeit  $p_{ij}$  der Zelle  $ij$  unter der  $H_0$  bestimmt werden. Da es sich hier um ein Verbundereignis aus zwei *unabhängigen* ( $H_0$ !) Einzelereignissen handelt (Zeile  $i$  und Spalte  $j$ ), ergibt sich  $p_{ij}$  als Produkt der Wahrscheinlichkeiten  $p_{i.}$  bzw.  $p_{.j}$  für die beiden verknüpften Einzelereignisse.

$$p_{ij} = p_{i.} \cdot p_{.j}$$

Die Einzelwahrscheinlichkeiten  $p_{i.}$  und  $p_{.j}$  sind allerdings nicht bekannt, sondern müssen durch die entsprechenden relativen Häufigkeiten in der Stichprobe geschätzt werden<sup>1</sup>. Z.B. wird die Wahrscheinlichkeit  $p_{i.}$  zur Zeile  $i$  geschätzt durch die relative Häufigkeit der Zeile  $i$  in der Stichprobe:

$$\hat{p}_{i.} := \frac{n_{i.}}{n}$$

Analog ergibt sich die geschätzte Wahrscheinlichkeit  $p_{.j}$  der Spalte  $j$ :

$$\hat{p}_{.j} := \frac{n_{.j}}{n}$$

Damit gilt für die geschätzte Wahrscheinlichkeit der Zelle  $ij$ :

$$\hat{p}_{ij} = \hat{p}_{i.} \cdot \hat{p}_{.j} = \frac{n_{i.}}{n} \cdot \frac{n_{.j}}{n} = \frac{n_{i.} \cdot n_{.j}}{n^2}$$

Die Wahrscheinlichkeit  $p_{ij}$  lässt sich interpretieren als Erwartungswert der Indikator-Zufallsvariablen  $X_{ij}$  zur Zelle  $(i, j)$  beim Ziehen *eines* Falles:

<sup>1</sup> Diese Formulierung geht davon aus, dass man *eine* Stichprobe gezogen und bei jedem Fall die *beiden* Merkmale Geschlecht und Fachbereich beobachtet hat. Ein anderes Stichprobenmodell läge vor, wenn man in jedem Fachbereich eine Stichprobe der festen Größe 50 gezogen und bei jedem Fall die eine Variable Geschlecht beobachtet hätte. Dann wären die Randwahrscheinlichkeiten der FB-Kategorien bekannt. Allerdings bleiben auch unter dem alternativen Stichprobenmodell alle vorgestellten Rechnungen und Entscheidungsregeln korrekt.

Tritt Zelle  $(i, j)$  auf, nimmt  $X_{ij}$  den Wert Eins an,  
bei jedem anderen Ergebnis nimmt  $X_{ij}$  den Wert Null an.

Werden  $n$  Fälle unabhängig gezogen, realisieren sich  $n$  unabhängige Zufallsvariablen  $X_{ij}^k$ ,  $k = 1, \dots, n$ , mit identischem Erwartungswert  $p_{ij}$ , und der Erwartungswert der Summenvariablen

$$E\left(\sum_{k=1}^n X_{ij}^k\right) = \sum_{k=1}^n E(X_{ij}^k) = n \cdot p_{ij}$$

ist die erwartete Häufigkeit der Zelle  $(i, j)$ .

Mit der geschätzten Wahrscheinlichkeit  $\hat{p}_{ij}$  ergibt sich sofort die geschätzte erwartete Häufigkeit  $m_{ij}$  in Pearsons Prüfstatistik:

$$m_{ij} = n \cdot \hat{p}_{ij} = n \cdot \frac{n_{i \cdot} \cdot n_{\cdot j}}{n^2} = \frac{n_{i \cdot} \cdot n_{\cdot j}}{n}$$

In Pearsons  $\chi_p^2$ -Statistik werden die quadrierten Abweichungen der beobachteten Häufigkeiten von den geschätzten Erwartungswerten unter der  $H_0$  aufsummiert. Durch das Quadrieren werden größere Diskrepanzen besonders stark gewichtet. Jede quadrierte Abweichung wird außerdem *normiert*, indem sie durch ihren erwarteten Wert dividiert wird. Steht etwa dem erwarteten Wert 5 die Häufigkeit 15 gegenüber, so resultiert die quadrierte und normierte Diskrepanz 20:

$$\frac{(15 - 5)^2}{5} = 20$$

Dieselbe Abweichung einer beobachteten Häufigkeit 2010 vom erwarteten Wert 2000 erbringt jedoch sinnvollerweise nur eine quadrierte und normierte Diskrepanz von 0,05:

$$\frac{(2010 - 2000)^2}{2000} = 0,05$$

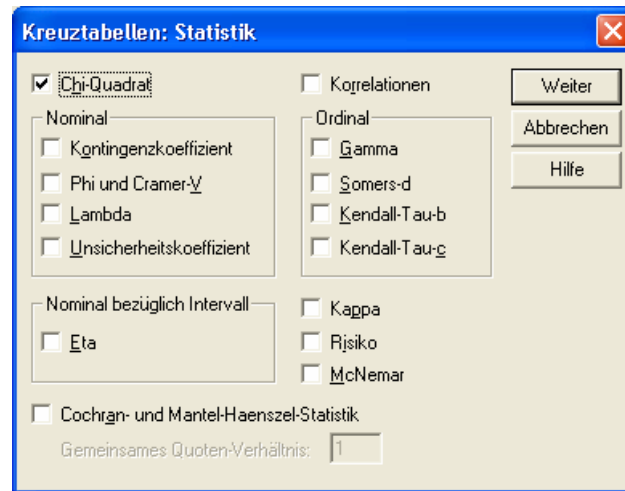
Der  $\chi_p^2$ -Wert ist offenbar, wie es in Abschnitt 7.1 von einer Prüfstatistik gefordert wird, indikativ für Abweichungen von der Nullhypothese.

Außerdem erfüllt die  $\chi_p^2$ -Teststatistik nach Pearson auch die Verteilungsbedingung aus Abschnitt 7.1, wenn auch nur approximativ. Unter der Nullhypothese ist die  $\chi_p^2$ -Statistik asymptotisch, d.h. für  $n \rightarrow \infty$ ,  $\chi^2$ -verteilt mit  $df = (z - 1) \cdot (s - 1)$  Freiheitsgraden.<sup>1</sup> Für unsere Kreuztabelle erhalten wir also:  $df = 1 \cdot 5 = 5$ .

Folglich kann mit Pearsons  $\chi_p^2$ -Statistik nicht nur die Plausibilität der  $H_0$  deskriptiv beurteilt werden, sondern es kann eine empirische Überschreitungswahrscheinlichkeit berechnet und nach den Regeln aus Abschnitt 7.1 ein Signifikanztest durchgeführt werden.

In SPSS wird die  $\chi_p^2$ -Statistik samt Signifikanztest mit dem Kontrollkästchen **Chi-Quadrat** in der **Kreuztabellen**-Subdialogbox **Statistik** angefordert:

<sup>1</sup> In diesem Satz treten zwei Symbole mit ähnlicher Gestalt aber deutlich verschiedener Bedeutung auf:  $\chi_p^2$  steht für eine (letztlich heuristisch definierte) Prüfgröße, mit  $\chi^2$  ist hingegen eine theoretische Verteilung gemeint.



Für unsere Daten erhalten wir folgendes Ergebnis:

#### Chi-Quadrat-Tests

	Wert	df	Asymptotische Signifikanz (2-seitig)
Chi-Quadrat nach Pearson	18,191 <sup>a</sup>	5	,003
Likelihood-Quotient	18,570	5	,002
Zusammenhang linear-mit-linear	3,197	1	,074
Anzahl der gültigen Fälle	283		

a. 0 Zellen (,0%) haben eine erwartete Häufigkeit kleiner 5. Die minimale erwartete Häufigkeit ist 17,68.

Es ergibt sich ein  $\chi_p^2$ -Wert von ca. 18,19, der bei  $df = 5$  unter der  $H_0$  eine Überschreitungswahrscheinlichkeit (**Asymptotische Signifikanz**) von ca. 0,003 hat, d.h. ein  $\chi_p^2$  - Wert  $\geq 18,19$  bei  $df = 5$  ist unter der  $H_0$  extrem unwahrscheinlich. Insbesondere ist die empirisch ermittelte Überschreitungswahrscheinlichkeit deutlich kleiner als die üblicherweise akzeptierte Irrtumswahrscheinlichkeit von  $\alpha = 0,05$ . Folglich entscheidet sich der  $\chi_p^2$  - Test klar für die  $H_1$ . In Abschnitt 7.1 wurde dieses Argumentationsmuster der Inferenzstatistik ausführlich erläutert.

Neben der  $\chi_p^2$ -Statistik nach Pearson, die aus heuristischen Überlegungen hervorgegangen zu sein scheint, berechnet SPSS noch die alternative Prüfgröße  $\chi_{LQ}^2$ , die auf dem **Likelihood-Quotienten - Prinzip** basiert. Letztere ist unter der  $H_0$  ebenfalls asymptotisch, d.h. für  $n \rightarrow \infty$ ,  $\chi^2$  - verteilt mit  $df = (z-1) \cdot (s-1)$  Freiheitsgraden, und trotz unterschiedlicher Herleitung sind beide Statistiken asymptotisch äquivalent, d.h. mit wachsender Stichprobengröße werden sie immer ähnlicher. Während bei größeren Stichproben wegen der asymptotischen Äquivalenz die Entscheidung für eine der beiden Prüfgrößen beliebig ist, sprechen einige Befunde dafür, bei kleineren Stichproben die  $\chi_p^2$ -Statistik nach Pearson wegen der besseren Verteilungsapproximation zu bevorzugen (siehe z.B. Hartung 1989, S. 439). Damit ist es also vertretbar, die  $\chi_p^2$ -Statistik nach Pearson grundsätzlich gegenüber der Likelihood-Quotienten - Prüfgröße zu bevorzugen. SPSS liefert stets beide Prüfgrößen. In unserem Fall sind die Unterschiede geringfügig und für die Testentscheidung irrelevant.

Die Pearson- und die Likelihood-Quotienten-Statistik zur Beurteilung der Unabhängigkeits- bzw. Homogenitätshypothese sind nur **asymptotisch**, d.h. für  $n \rightarrow \infty$ ,  $\chi^2$ -verteilt. Für die Zulässigkeit der zugehörigen Hypothesentests setzt man üblicherweise voraus, dass alle **erwarteten** Häufigkeiten  $m_{ij}$  mindestens gleich 5 sind. SPSS protokolliert daher für jede Kreuztabelle die minimale erwartete Häufigkeit. In unserem Fall beträgt sie 17,682, so dass keine Einwände gegen Tests auf Basis der  $\chi_p^2$ - bzw.  $\chi_{LQ}^2$ -Statistik bestehen.

Manche Autoren formulieren etwas abgeschwächte Voraussetzungen für die erwarteten Häufigkeiten. Siegel (1976, S. 107) verlangt z.B. für  $\chi_p^2$ -Tests mit  $df > 1$ , dass die beiden folgenden Bedingungen erfüllt sind:

- Weniger als 20% der Zellen haben eine erwartete Häufigkeit kleiner als 5.
- Keine Zelle hat eine erwartete Häufigkeit kleiner als 1.

Neben den beiden Statistiken zur Prüfung der Unabhängigkeits- bzw. Homogenitätshypothese liefert SPSS unter der Bezeichnung **Zusammenhang linear-mit-linear** auch noch den  $\chi_{MH}^2$ -Wert nach **Mantel-Haenszel** zur Beurteilung der **linearen** Beziehung zwischen den beiden Variablen. Diese Statistik darf nur interpretiert werden, *wenn beide Variablen Intervallskalqualität besitzen*. Es handelt sich nämlich schlicht um die mit  $(n - 1)$  multiplizierte quadrierte Produkt-Moment-Korrelation zwischen den beiden Variablen:

$$\chi_{MH}^2 := r^2(n - 1)$$

Da wir zwei kategoriale Variablen betrachten, ist diese Statistik in unserem Fall völlig sinnlos.

### 10.3.2 Exakte Tests

Für die  $(2 \times 2)$ -Kreuztabellen gibt es seit Jahrzehnten mit dem **exakten Test von Fisher** eine glänzende Alternative zu den approximativen  $\chi^2$ -Tests. Wie sein Name sagt, kommt Fishers Test ohne Approximationen aus und ist daher bei jeder Stichprobe anwendbar. Erfreulicherweise bietet SPSS mittlerweile exakte Tests auch für beliebige  $(z \times s)$ -Kreuztabellen.

Eine ausführliche Beschreibung der neuen statistischen Verfahren, die durch das SPSS-Zusatzmodul **Exact Tests** implementiert werden, finden Sie auf dem WWW-Server der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend über:

[Weitere Serviceangebote > EDV-Dokumentationen > Elektronische Publikationen > Statistische Spezialthemen > Exakte Tests mit SPSS](#)

Allerdings sind die traditionellen asymptotischen Verfahren nun keinesfalls obsolet, weil der exakte Test für  $(z \times s)$ -Kreuztabellen wegen seines enormen Rechenaufwandes nur für kleine Stichproben durchführbar ist. Insgesamt steht für die meisten Situationen ein angemessenes Verfahren zur Verfügung:

- Wenn die Anwendbarkeitskriterien für die asymptotischen Verfahren erfüllt sind, sollten Sie den Pearson-Test verwenden.
- Anderenfalls sollten Sie einen exakten Test versuchen.

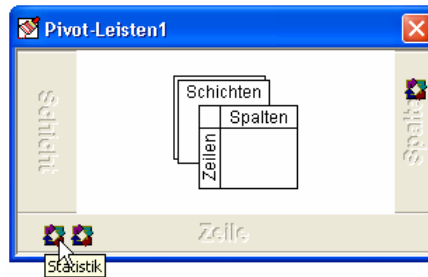
Wenn bei einer Kreuztabelle die Minimalanforderungen an die erwarteten Häufigkeiten *nicht* erfüllt sind, *und* der exakte Test aufgrund des insgesamt zu großen Stichprobenumfangs scheitert, müssen Sie die verantwortlichen schwach besetzten Zeilen bzw. Spalten entweder löschen oder miteinander bzw. mit anderen Zeilen/Spalten zusammenlegen.

In einem Anwendungsbeispiel wollen wir die Daten aus dem ersten Abschnitt des SPSS-Handbuchs zum Modul **Exact Tests** (1996, S. 1) verwenden. Es handelt sich um Prüfungsergebnisse weißer, schwarzer, asiatischer und hispanoider Feuerwehrbewerber in einer amerikanischen Kleinstadt.

		Hautfarbe				Gesamt
		Weiß	Schwarz	Asiatisch	Mittel- und Südamerika	
Anzahl	Bestanden	5	2	2	0	9
	Unklar	0	1	0	1	2
	Durchgefallen	0	2	3	4	9
	Gesamt	5	5	5	5	20
Prozent	Bestanden	100,0%	40,0%	40,0%	,0%	45,0%
	Unklar	,0%	20,0%	,0%	20,0%	10,0%
	Durchgefallen	,0%	40,0%	60,0%	80,0%	45,0%
	Gesamt	100,0%	100,0%	100,0%	100,0%	100,0%

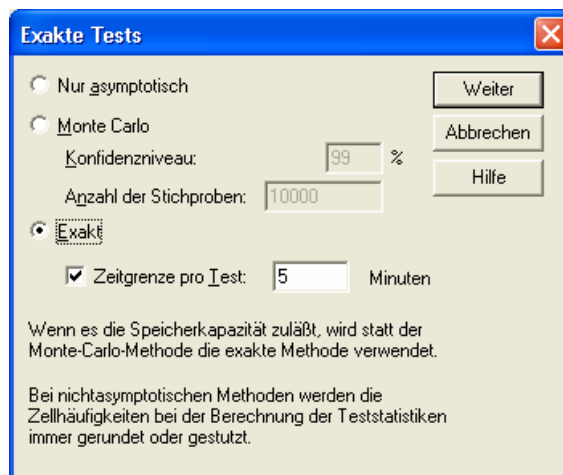
Technische Hinweise:

- Die Tabelle enthält spaltenbezogene relative Häufigkeiten (Subdialogbox **Zellen**).
- Für die beiden Zeilendimensionen wurde per Pivot-Werkzeug die Schachtelungsordnung geändert:



Wir wollen die Nullhypothese testen, dass die Prüfungsergebnisse von der Hautfarbe unabhängig sind.

Nach einem Mausklick auf den **Exakt**-Schalter in der Dialogbox zur Kreuztabellenanalyse können wir in der folgenden Subdialogbox die **exakte** Testmethode wählen:



Daraufhin erhalten wir neben den approximativen Ergebnissen auch exakte Überschreitungswahrscheinlichkeiten für die Pearson- und die Likelihood-Quotienten – Prüfstatistik. Außerdem führt SPSS noch eine Verallgemeinerung des exakten Tests von Fisher durch, der in seiner klassischen Variante auf  $(2 \times 2)$ -Tabellen beschränkt ist:

Chi-Quadrat-Tests

	Wert	df	Asymptotische Signifikanz (2-seitig)	Exakte Signifikanz (2-seitig)	Exakte Signifikanz (1-seitig)	Punkt-Wahrscheinlichkeit
Chi-Quadrat nach Pearson	11,556 <sup>a</sup>	6	,073	,040		
Likelihood-Quotient	15,673	6	,016	,040		
Exakter Test nach Fisher	11,239			,040		
Zusammenhang linear-mit-linear	8,276 <sup>b</sup>	1	,004	,004	,002	,001
Anzahl der gültigen Fälle	20					

a. 12 Zellen (100,0%) haben eine erwartete Häufigkeit kleiner 5. Die minimale erwartete Häufigkeit ist ,50.

b. Die standardisierte Statistik ist 2,877.

Die approximativen  $\chi^2$  - Unabhängigkeitstests (Pearson und Likelihood-Quotient) sind nicht anwendbar, weil in allen 12 Zellen die erwartete Häufigkeit kleiner als 5 ist. Wer dieses Problem ignoriert, andererseits aber weiß, dass der Pearson-Test gegenüber dem Likelihood-Quotienten - Test im Allgemeinen wegen der besseren Verteilungsapproximation zu bevorzugen ist, gelangt zu einer falschen Testentscheidung. Die korrekte Überschreitungswahrscheinlichkeit beträgt 0,04, was zur Ablehnung der Nullhypothese führt. Der asymptotische Pearson- $\chi^2$  - Test empfiehlt durch eine Überschreitungswahrscheinlichkeit von 0,07 hingegen, die Nullhypothese beizubehalten.

### 10.3.3 Besonderheiten bei (2 × 2)-Tabellen

#### 10.3.3.1 Ein klarer Fall für Fishers Test

Im beliebten Spezialfall der (2 × 2)-Tabelle ist Fishers Test nicht nur *exakt* für beliebige Stichproben, sondern er besitzt sogar unter allen „vernünftigen“, nämlich unter den so genannten unverfälschten, Tests die besten Güteeigenschaften. Daher sollten Sie in dieser Situation grundsätzlich Fishers Test verwenden.

Die oben beschriebenen Rechenzeitprobleme bei exakten Tests für allgemeine (z × s)-Kreuztabellen treten bei Fishers Test für die (2 × 2)-Tabelle *nicht* auf.

#### 10.3.3.2 Einseitige Hypothesen

Bei einer (2 × 2)-Tabelle lässt sich im Unterschied zu allen anderen Tabellen die Unabhängigkeits- bzw. Homogenitätshypothese auch *einseitig* formulieren. Wenn wir uns z.B. beim Vergleich der Frauenanteile unter den Studierenden der Universität Trier auf die Fachbereiche III und IV beschränken, können wir die folgende einseitige Homogenitätshypothese aufstellen:

H<sub>0</sub>: Der Frauenanteil ist im FB IV mindestens genauso groß wie im FB III.

H<sub>1</sub>: Der Frauenanteil ist im FB IV kleiner als im FB III.

Aus den (z.B. per Filterbedingung, vgl. Abschnitt 9) eingeschränkten Beispieldaten erhalten wir folgende Ergebnisse:

Kreuztabelle

	Fachbereiche an der Universität Trier		Gesamt
	III	IV	
Frauen	18 45,0% 50,0%	22 55,0% 31,0%	40 100,0% 37,4%
Männer	18 26,9% 50,0%	49 73,1% 69,0%	67 100,0% 62,6%
Gesamt	36 33,6% 100,0%	71 66,4% 100,0%	107 100,0% 100,0%

Chi-Quadrat-Tests

	Wert	df	Asymptotische Signifikanz (2-seitig)	Exakte Signifikanz (2-seitig)	Exakte Signifikanz (1-seitig)
Chi-Quadrat nach Pearson	3,689 <sup>b</sup>	1	,055		
Kontinuitätskorrektur <sup>a</sup>	2,922	1	,087		
Likelihood-Quotient	3,643	1	,056		
Exakter Test nach Fisher				,061	,044
Zusammenhang linear-mit-linear	3,655	1	,056		
Anzahl der gültigen Fälle	107				

a. Wird nur für eine 2x2-Tabelle berechnet

b. 0 Zellen (,0%) haben eine erwartete Häufigkeit kleiner 5. Die minimale erwartete Häufigkeit ist 13,46.

Wie wir bereits wissen, beträgt der Frauenanteil im FB III 50% und im FB IV 31%, die deskriptiven Statistiken fallen also klar im Sinne der Alternativhypothese aus. Der nach den obigen Überlegungen zu verwendende exakte Test von Fisher liefert für die *zweiseitige* Fragestellung eine Überschreitungswahrscheinlichkeit von 0,061, so dass die Nullhypothese beibehalten werden müsste. Bei *einseitiger* Testung erhalten wir jedoch eine Überschreitungswahrscheinlichkeit von 0,04, so dass die Nullhypothese verworfen werden kann.

Beachten Sie abschließend noch, dass sich bei Fishers Test die einseitige Überschreitungswahrscheinlichkeit keinesfalls durch Halbieren der zweiseitigen Überschreitungswahrscheinlichkeit ergibt. Die in Abschnitt 7.1 für den Spezialfall des t-Tests angegebene Regel zur Berechnung der einseitigen Überschreitungswahrscheinlichkeit aus der zweiseitigen darf also nicht generalisiert werden.

### 10.3.3.3 Kontinuitätskorrektur nach Yates

Bei  $(2 \times 2)$ -Tabellen berechnet SPSS traditionell auch eine  $\chi^2_Y$ -Größe mit Kontinuitätskorrektur nach Yates. Sie soll bei kleineren Stichproben der Pearson- $\chi^2_P$ -Statistik überlegen sein. Gemäß Abschnitt 10.3.3.1 ist sie allerdings irrelevant, weil in der  $(2 \times 2)$ -Situation Fishers exakter Tests in jedem Fall vorzuziehen ist.





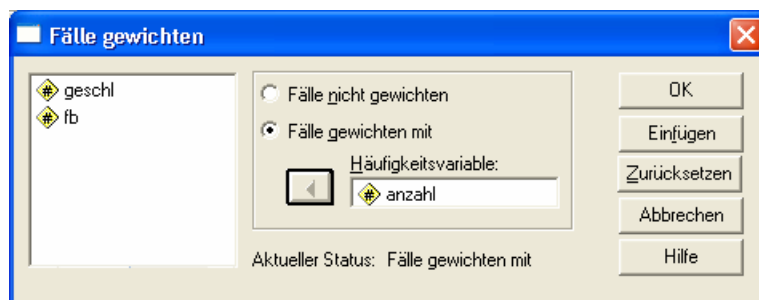
- Die Fälle werden mit der Variablen ANZAHL gewichtet. Damit tun wir z.B. so, als seien 16 Fälle mit dem Geschlecht 1 und dem Fachbereich 1 vorhanden gewesen. Aber das stimmt ja wirklich. Offenbar ist die Fallgewichtung doch nicht so sinnlos.

Um eine Gewichtsvariable zu vereinbaren, rufen wir mit dem Menübefehl

### Daten > Fälle gewichten

eine Dialogbox auf, die folgende Optionen anbietet:

- **Fälle nicht gewichten**  
Damit wird eine bestehende Gewichtung wieder aufgehoben.
- **Fälle gewichten mit**  
Die gewünschte Variable wird mit dem Transportschalter in die Position der **Häufigkeitsvariablen** gebracht, z.B.:



In der Dialogbox wird außerdem angezeigt, ob momentan eine Gewichtungsvariable vereinbart ist. Dieselbe Information erscheint auch in der Statuszeile des Datenfensters (siehe oben).

Beim Einsatz von Gewichtungsvariablen ist noch zu beachten:

- Zur Gewichtung kann natürlich nur eine numerische Variable verwendet werden; diese darf allerdings auch gebrochene Werte enthalten. Negative und fehlende Werte werden auf 0 gesetzt, d.h. die betroffenen Fälle werden nicht berücksichtigt, solange die Gewichtungsvariable aktiv ist.
- Ist beim Speichern der Arbeitsdatei eine Gewichtung aktiv, so wird diese mit abgespeichert und ist bei späterer Verwendung der Datendatei in Kraft.
- Bei der in diesem Abschnitt beschriebenen Anwendung der Gewichtungsoption wird dafür gesorgt, dass alle tatsächlich in der Studie vorhandenen Beobachtungen mit dem Gewicht 1 in die Kreuztabellenanalyse eingehen. Wenn die vorhandenen Beobachtungen individuelle Gewichte ( $\neq 1$ ) erhalten, werden natürlich Signifikanztests erheblich beeinflusst. Auf jeden Fall muss dann die Gewichtungsvariable einen Mittelwert von 1 haben, d.h. die Summe der Gewichte muss gerade den Stichprobenumfang ergeben.

## 11.2 Übung

Prüfen Sie anhand der Daten aus der Tabelle am Anfang von Abschnitt 11.1 die Nullhypothese, dass die Merkmale Geschlecht und Fachbereich unabhängig sind.

## 12 Auswertung von Mehrfachwahlfragen

In Abschnitt 1.4.2.3 wurde betont, dass mit einer Mehrfachwahlfrage nicht etwa *ein* mysteriöses Merkmal mit mehreren Ausprägungen erfasst wird, wie es wohl durch manche Köpfe bzw. Alpträume spukt, sondern *eine Familie* inhaltlich verwandter dichotomer Merkmale. Eine leichte

Komplikation tritt erst auf, wenn zur Vereinfachung der Erfassung ein sparsames Set aus kategorialen Variablen definiert worden ist, das für viele Auswertungen erst „ausgepackt“ werden muss.

Grundsätzlich besteht kein Bedarf für spezielle Auswertungsverfahren für die mit Mehrfachwahlfragen erfassten Variablen. Es ist allerdings gelegentlich sinnvoll, eine Häufigkeits- oder Kreuztabellenanalyse für *alle* Mitglieder einer Familie dichotomer Variablen (ob aus einer Mehrfachwahlfrage entstanden oder wie auch immer) in gleicher Form auszuführen. Für diese Situation bietet SPSS gewisse Rationalisierungsmöglichkeiten, die in diesem Abschnitt vorgestellt werden sollen. Außerdem kann SPSS für die mit einem sparsamen Set aus kategorialen Variablen erfassten dichotomen Merkmale Häufigkeits- und Kreuztabellenanalysen ohne vorheriges Auspacken durchführen.

In den Abschnitten 12.1 bzw. 0 wird die Häufigkeits- bzw. Kreuztabellenanalyse für eine Familie von dichotomen Variablen beschrieben. In Abschnitt 12.3 wird demonstriert, wie man mit Hilfe einiger SPSS-Kommandos aus einem sparsamen Set kategorialer Variablen ein vollständiges Set dichotomer Variablen erzeugen kann.

### 12.1 Häufigkeitstabellen

Im Teil 4a unseres Fragebogens haben die Teilnehmer für fünf konkrete Motive, den SPSS-Kurs zu besuchen, und eine Restkategorie alles zutreffende angekreuzt. Es liegt nahe, eine Übersicht zu erstellen, aus der für die einzelnen Motive hervorgeht, wie häufig sie gewählt worden sind. Natürlich können wir die Zustimmungsfrequenzen bei den Motiv-Variablen z.B. auch mit der längst bekannten Häufigkeitsanalyse (**Analysieren > Deskriptive Statistiken > Häufigkeiten**) bestimmen lassen. SPSS bietet jedoch für solche *Gruppen zusammengehöriger Variablen* eine Prozedur an, welche die Zustimmungshäufigkeiten sowie einige zusätzliche Ergebnisse in besonders kompakter Form ausgibt. Wir erhalten für unsere Daten die folgende Tabelle:

\$Motive Frequencies

		Antworten		Prozent der Fälle
		N	Prozent	
Motive zur Kursteilnahme <sup>a</sup>	Eigene Studie	23	56,1%	76,7%
	Bewerbung um Stelle	1	2,4%	3,3%
	Bewerbung um HIWI-Job	1	2,4%	3,3%
	Interesse an der EDV	5	12,2%	16,7%
	Interesse an Statistik	10	24,4%	33,3%
	Andere Motive	1	2,4%	3,3%
Gesamt		41	100,0%	136,7%

a. Dichotomie-Gruppe tabellarisch dargestellt bei Wert 1.

Es zeigt sich etwa, dass 23 Personen (= 76,7% aller validen Fälle) dem ersten Motiv zugestimmt haben. Diese 23 positiven Antworten machen 56,1% der insgesamt 41 von allen Teilnehmern geäußerten Zustimmungen aus. Ein Fall, auf den wir später noch eingehen müssen, fand keines der fünf konkreten Motive für sich passend und markierte die Restkategorie (*Andere Motive*).

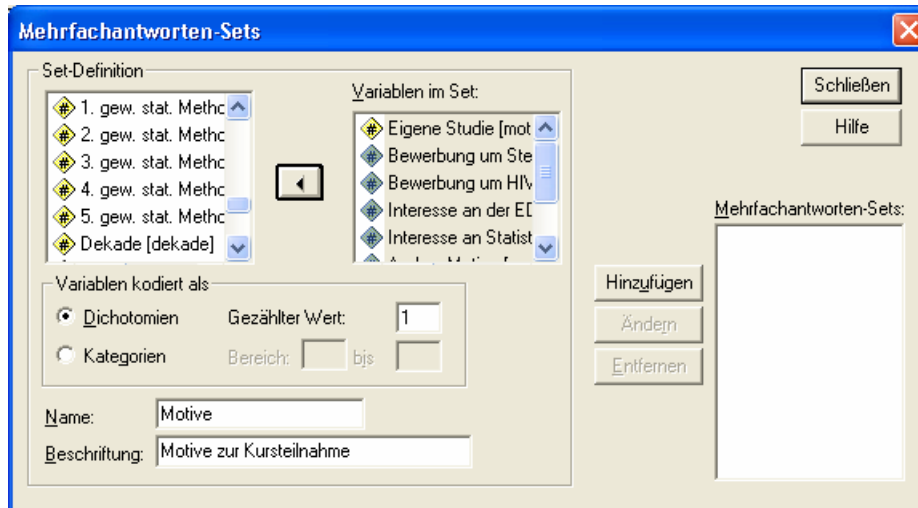
Wie die Titelzeile zu obiger Tabelle zeigt, wurde eine Prozedur für die **Variablengruppe** \$MOTIVE ausgeführt, die natürlich zuvor definiert werden muss. Wählen Sie dazu den Menübefehl:

#### **Analysieren > Mehrfachantwort > Sets definieren**

In der nun erscheinenden Dialogbox sind folgende Aktionen nötig:

- Befördern Sie die Variablen MOTIV1 bis MOTIV5 sowie ANDERE in die Liste **Variablen im Set**.
- Tragen Sie im Rahmen **Variablen kodiert als** für die bei uns zutreffende **dichotome** Option die Eins als **gezählten Wert** ein.
- Vereinbaren Sie für das Set den Namen *Motive* und das Label *Motive zur Kursteilnahme*.

Danach müsste Ihre Dialogbox ungefähr so aussehen:

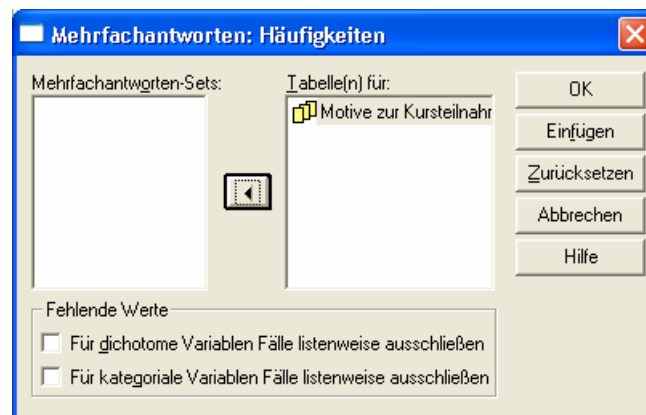


Nehmen Sie abschließend mit **Hinzufügen** die neue Gruppe in die Liste der **Mehrfachantworten-Sets** auf, und **schließen** Sie die Dialogbox.

Nun können Sie obige Tabelle anfordern mit dem Menübefehl

**Analysieren > Mehrfachantwort > Häufigkeiten**

und der zugehörigen Dialogbox:



Entfernt man die Variable ANDERE zur Restkategorie der sonstigen Motive aus dem Set \$MOTIVE, dann resultieren folgende Ergebnisse:

#### Fallzusammenfassung

	Fälle					
	Gültig		Fehlend		Gesamt	
	N	Prozent	N	Prozent	N	Prozent
\$Motive <sup>a</sup>	29	93,5%	2	6,5%	31	100,0%

a. Dichotomie-Gruppe tabellarisch dargestellt bei Wert 1.

\$Motive Frequencies

		Antworten		Prozent der Fälle
		N	Prozent	
Motive zur Kursteilnahme <sup>a</sup>	Eigene Studie	23	57,5%	79,3%
	Bewerbung um Stelle	1	2,5%	3,4%
	Bewerbung um HIWI-Job	1	2,5%	3,4%
	Interesse an der EDV	5	12,5%	17,2%
	Interesse an Statistik	10	25,0%	34,5%
Gesamt		40	100,0%	137,9%

a. Dichotomie-Gruppe tabellarisch dargestellt bei Wert 1.

Bei der ersten Tabelle erstaunt, dass nur **29** gültige Fälle gemeldet werden, obwohl sich in unserer KFA-Datendatei **30** Fälle mit vollständig vorhandenen MOTIV-Werten befinden. Des Rätsels Lösung ist eine SPSS-Eigenart bei der Analyse von Mehrfachwahl-Sets aus dichotomen Variablen: Als gültig werden nur solche Fälle betrachtet, die bei mindestens einer Set-Variablen den zu zählenden Wert besitzen (bei uns also die Eins). Daher wird neben dem Fall 13 mit SYSMIS bei den Variablen MOTIV1 bis MOTIV5 auch der dritte Fall ausgeschlossen, der *alle konkreten Motive verneint*, aber die Restkategorie markiert hat. Wenn SPSS in obiger Ausgabe z.B. zum Motiv 1 meldet, dass 79,3% der Fälle (23 von 29) zugestimmt hätten, ist dies schlicht falsch.

SPSS ignoriert nicht nur Fälle, die bei keiner Set-Variablen den zu zählenden Wert besitzen, sondern auch Variablen, bei denen der zu zählende Wert nicht auftritt. Hätte in unserem Beispiel kein Teilnehmer das Motiv 5 bejaht, würde es in der Häufigkeitstabelle komplett fehlen.

Der Mangel in obiger Ausgabe wurde aufgrund der protokollierten Anzahl fehlender Werte entdeckt. Sie sollten grundsätzlich bei allen SPSS-Ausgaben die protokollierten Fallzahlen überprüfen, weil sich viele technische Fehler durch eine zu niedrige oder zu hohe Anzahl auswertbarer Fälle verraten. Im aktuellen Beispiel ist SPSS für den „Fehler“ verantwortlich; in der Regel werden Sie auf diese Weise Ihre eigenen Fehler entdecken.

Die einzige Möglichkeit, definierte Mehrfachwahl-Sets zu speichern, besteht darin, die zur Häufigkeitsanalyse bzw. zur anschließend beschriebenen Kreuztabellenanalyse gehörige Syntax zu sichern. In den korrespondierenden SPSS-Kommandos sind die Set-Definitionen nämlich enthalten, z.B.:

```
MULT RESPONSE
  GROUPS=$Motive 'Motive zur Kursteilnahme' (motiv1 motiv2 motiv3 motiv4
  motiv5 andere (1))
  /FREQUENCIES=$Motive .
```

## 12.2 Kreuztabellen

Wenn wir uns für Geschlechtsunterschiede bei der Zustimmung zu den fünf konkreten Motiven interessieren (z.B.: *Wer interessiert sich mehr für Statistik?*), sind genau *fünf* (2×2)-Tabellen zu analysieren. Über den aus Abschnitt 10 bekannten Menübefehl **Analysieren > Deskriptive Statistiken > Kreuztabellen** erhalten wir z.B. für das Statistik-Motiv (Nummer fünf) folgendes Ergebnis:

Interesse an Statistik \* Geschlecht Kreuztabelle

			Geschlecht		Gesamt
			Frau	Mann	
Interesse an Statistik	Nein	Anzahl	15	5	20
		% von Interesse an Statistik	75,0%	25,0%	100,0%
		% von Geschlecht	62,5%	83,3%	66,7%
	Ja	Anzahl	9	1	10
		% von Interesse an Statistik	90,0%	10,0%	100,0%
		% von Geschlecht	37,5%	16,7%	33,3%
Gesamt	Anzahl	24	6	30	
	% von Interesse an Statistik	80,0%	20,0%	100,0%	
	% von Geschlecht	100,0%	100,0%	100,0%	

Weil die Motiv-Variablen nur zwei Ausprägungen haben, sind die Ergebnisse zur Nein-Kategorie überflüssig. Es genügt zu wissen, dass 37,5% von den 24 Frauen und 16,7% von den sechs Männern ein Interesse an Statistik angegeben haben. Durch Verzicht auf die redundanten Zeilen erhält man eine sehr kompakte Darstellung der (2×2)-Tabellen zu Geschlechtsunterschieden bei den Kursmotiven:

Kreuztabelle \$Motive\*geschl

			Geschlecht		Gesamt
			Frau	Mann	
Motive zur Kursteilnahme	Eigene Studie	Anzahl	19	4	23
		Innerhalb \$Motive %	82,6%	17,4%	
		Innerhalb geschl %	79,2%	66,7%	
	Bewerbung um Stelle	Anzahl	1	0	1
		Innerhalb \$Motive %	100,0%	,0%	
		Innerhalb geschl %	4,2%	,0%	
	Bewerbung um HIWI-Job	Anzahl	0	1	1
		Innerhalb \$Motive %	,0%	100,0%	
		Innerhalb geschl %	,0%	16,7%	
	Interesse an der EDV	Anzahl	3	2	5
		Innerhalb \$Motive %	60,0%	40,0%	
		Innerhalb geschl %	12,5%	33,3%	
	Interesse an Statistik	Anzahl	9	1	10
		Innerhalb \$Motive %	90,0%	10,0%	
		Innerhalb geschl %	37,5%	16,7%	
	Andere Motive	Anzahl	1	0	1
		Innerhalb \$Motive %	100,0%	,0%	
		Innerhalb geschl %	4,2%	,0%	
	Gesamt	Anzahl	24	6	30

Prozentsätze und Gesamtwerte beruhen auf den Befragten.

a. Dichotomy group tabulated at value 1.

Beachten Sie bitte: Dies ist **nicht eine** (6×2)-Kontingenztafel, **sondern dies sind sechs** (2×2)-Kontingenztafeln. In der vorletzten Zeile befindet sich etwa die Essenz der MOTIV5 × GESCHL - Kontingenztafel.

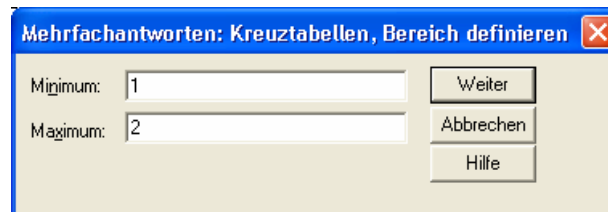
Sie können die Dialogbox zu obiger Kombitabelle anfordern mit

**Analysieren > Mehrfachantwort > Kreuztabellen**

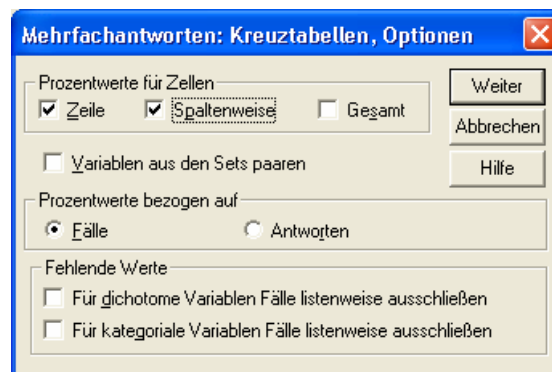
Befördern Sie das **Mehrfachantworten-Set** \$MOTIVE in die **Zeilen**:



Für die **Spalten**-Variable GESCHL müssen Sie noch den folgenden **Bereich definieren**:



Die spalten- und zeilenrelativierten Prozentangaben werden in der **Optionen**-Subdialogbox angefordert:



Auch bei dieser Kontingenzanalyse ist die in Abschnitt 12.1 gerügte MD-Konzeption der SPSS-Mehrfachwahl-Auswertung zu beachten. Wäre nicht die Variable ANDERE Mitglied im Set \$MOTIVE, dann würde SPSS in der Kombitabelle nur noch diejenigen Fälle berücksichtigen, die mindestens ein konkretes Motiv bejaht haben.

### 12.3 Ein sparsames Set kategorialer Variablen expandieren

In Abschnitt 1.4.2.3 wurde das sparsame Set aus kategorialen Variablen für Mehrfachwahlfragen mit sehr vielen Antwortmöglichkeiten zur Vereinfachung der Erfassung empfohlen. Zwar ist diese Datenstruktur kein Nachteil bei den Analyseprozeduren, die in den Abschnitten 12.1 und 0 beschrieben wurden, doch sind Auswertungen denkbar, die ein vollständiges Set aus dichotomen Variablen erfordern. In dieser Situation kann man das sparsame Set mit Hilfe der SPSS-Kommandosprache „expandieren“. Die folgenden Kommandos erzeugen zu unseren Variablen METH1 bis METH3 die acht dichotomen Variablen STAT1 bis STAT8, die für jeweils eine bestimmte statistische Methode festhalten, ob sie genannt worden ist (Wert Eins) oder nicht (Wert Null):

```
do repeat stat = stat1 to stat8 /n = 1 to 8.
  do if (meth1 = n) or (meth2 = n) or (meth3 = n).
    compute stat = 1.
  else.
    compute stat = 0.
  end if.
end repeat.
execute.
```

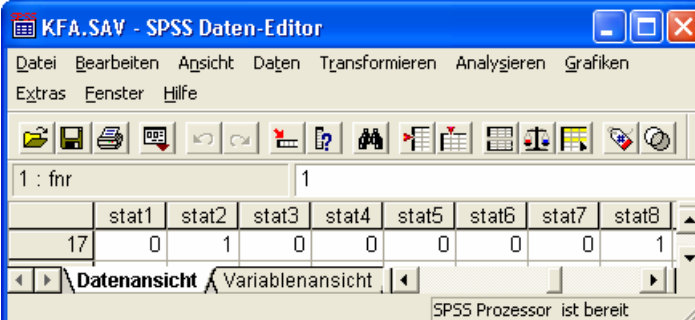
Die Variable STAT2 steht z.B. für die Regressionsanalyse, weil gemäß Kodierplan bei einer der Variablen METH1 bis METH3 eine 2 zu notieren war, wenn ein Fall im Fragebogenteil 4b die Regressionsanalyse genannt hatte.

Beim Fall Nr. 17 wurden die genannten Methodenwünsche 8 (= logistische Regression) und 2 (= Regressionsanalyse) folgendermaßen mit dem sparsamen Set kategorialer Variablen METH1 bis METH3 erfasst:



1 : fnr	smg	meth1	meth2	meth3	dekade
17	1	8	2	0	1

Daraus ergeben sich folgende Werte für die Variablen STAT1 bis STAT8:



1 : fnr	stat1	stat2	stat3	stat4	stat5	stat6	stat7	stat8
17	0	1	0	0	0	0	0	1

In obiger Syntax werden zwei ausgesprochen nützliche Kontrollstrukturen der SPSS-Kommandosprache verwendet:



**Schleife für strukturgleiche Transformationen**

Die (DO REPEAT - END REPEAT) - Schleife wird achtmal ausgeführt, wobei im  $i$ -ten Umlauf die beiden Stellvertreter STAT und N gerade mit den  $i$ -ten Elementen der zugehörigen Listen identisch sind.

**Fallunterscheidung**

Beim Ausführen der (DO IF - ELSE - END IF) - Struktur passiert in Abhängigkeit vom Wahrheitswert des logischen Ausdruck folgendes:

<b>Wert des logischen Ausdrucks</b>	<b>Aktion</b>
wahr, z.B. im ersten Schleifenumlauf bei METH1 = 1, METH2 = 2, METH3 = SYSMIS	Das erste COMPUTE-Kommando wird ausgeführt.
falsch, z.B. im ersten Schleifenumlauf bei METH1 = 3, METH2 = 5, METH3 = 8	Das zweite COMPUTE-Kommando wird ausgeführt.
unbestimmt, z.B. im ersten Schleifenumlauf bei METH1=SYSMIS,METH2=SYSMIS,METH3=SYSMIS	Die neuen Variablen STAT1 bis STAT8 behalten den Initialisierungswert SYSMIS.

---

## 13 Datendateien im Textformat einlesen

Gelegentlich sind Daten auszuwerten, die in Textdateien vorliegen. In Abschnitt 3.1.2 wurden zwei Dateiformate beschrieben, die uns dabei begegnen können:

- positionierte Daten (feste Breite)
- separierte Daten (mit Trennzeichen).

Zum Importieren von Textdatendateien stellt SPSS einen leistungsfähigen Assistenten zur Verfügung, der mit

### **Datei > Textdaten lesen**

gestartet wird. Er kommt aber auch dann zum Einsatz, wenn Sie nach

### **Datei > Öffnen > Daten**

eine Textdatendatei wählen.

An der im Vorwort vereinbarten Stelle finden Sie die Dateien **kfar-kv-pos.txt** und **kfar-kv-sep.txt** mit positionierten bzw. separierten KFA-Daten von 77 Fällen. Es bietet sich an, diese Daten einzulesen, um die in Abschnitt 8 durch graphische Datenexploration gewonnene Moderatorversion der differentialpsychologischen Hypothese anhand einer unabhängigen Stichprobe zu überprüfen.

### **13.1 Import von positionierten Textdaten (feste Breite)**

In der Datei **kfar-kv-pos.txt** sind die Werte eines Falles auf zwei Zeilen verteilt, und jede Variable hat eine feste Position im Datensatz eines Falles (z.B. Variable AERGO in Zeile 2, Spalten 5-6), so dass auch ihre Breite fixiert ist.

```
11 177115848
12  6 6 431214542432 110000
21 177115955
22  4 8 343335442442 110010
31 174416048
32  3 8 433224443342 100010
41 175116578
42  2 2 553125544531 100100
. . . . . . . . .
. . . . . . . . .
```

Die für uns relevanten Variablen haben folgende Positionen:

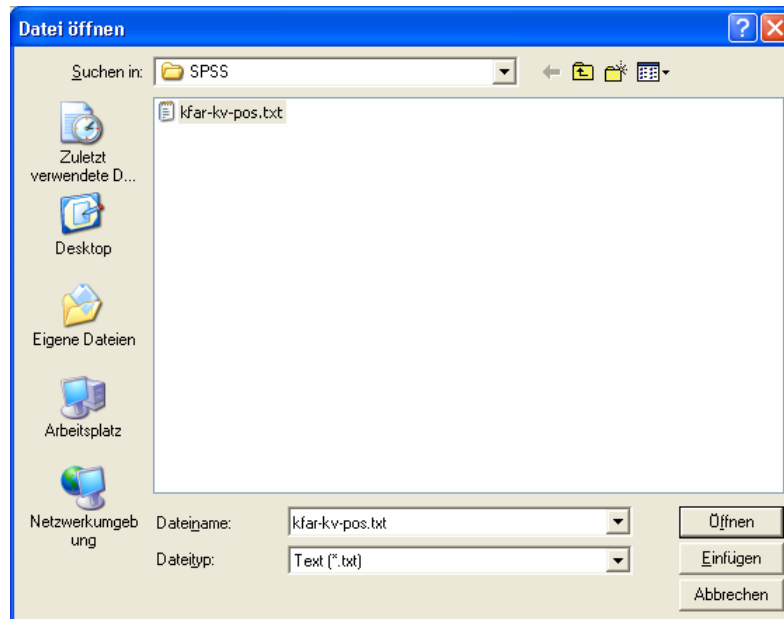
Variable	Datenzeile	Spalten
GESCHL	1	5
AERGO	2	5-6
AERGM	2	7-8
LOT01-LOT12	2	10-21

Alle übrigen Variablen können wir ignorieren.

Gehen Sie folgendermaßen vor, um die relevanten Daten zu importieren:

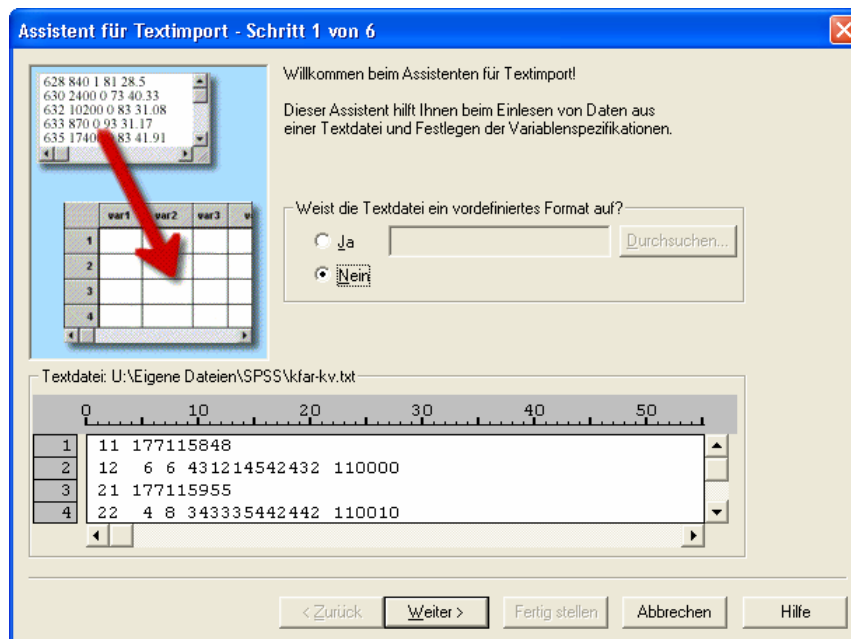
## Textimport-Assistenten starten und Datei auswählen

Nach dem Start des Textimport-Assistenten ist zunächst die Eingabedatei zu wählen:



### Schritt 1

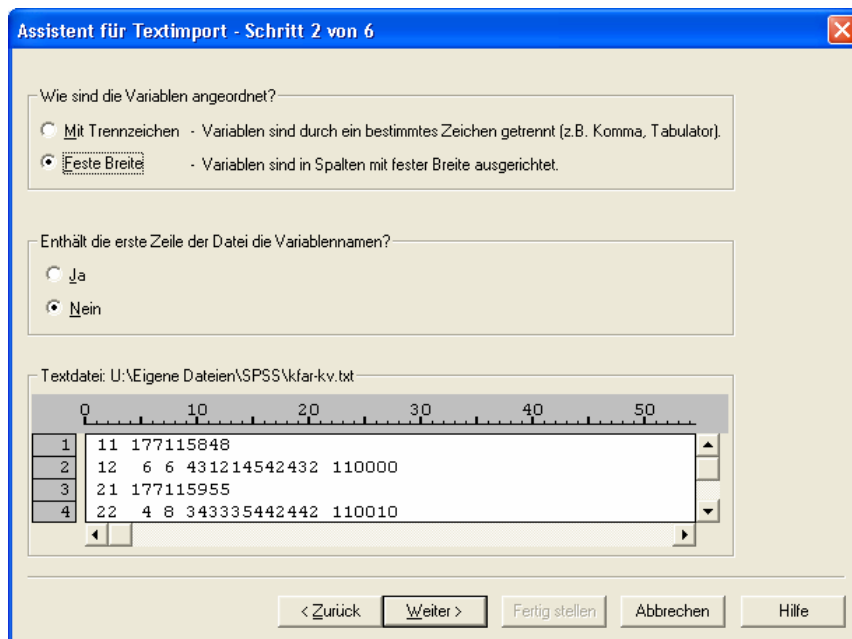
Im ersten Schritt zeigt der Assistent den Anfang unserer Datei und akzeptiert ggf. ein **vordefiniertes Format** aus früheren Assistenteneinsätzen, das die Dateistruktur beschreibt.



Da wir auf eine solche Vorarbeit nicht zurückgreifen können, machen wir **weiter**.

## Schritt 2

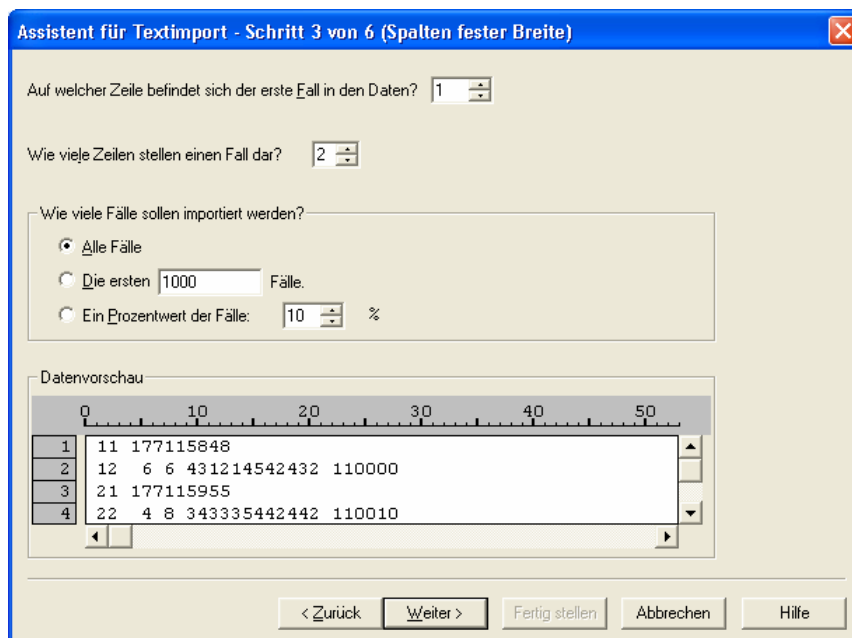
Im zweiten Schritt teilen wir mit, dass die Variablen in unserer Eingabedatei feste Positionen bzw. eine **feste Breite** besitzen:



Von der Möglichkeit, in der **ersten Zeile der Datei die Variablennamen** zu transportieren, wird in unserem Beispiel kein Gebrauch gemacht.

## Schritt 3

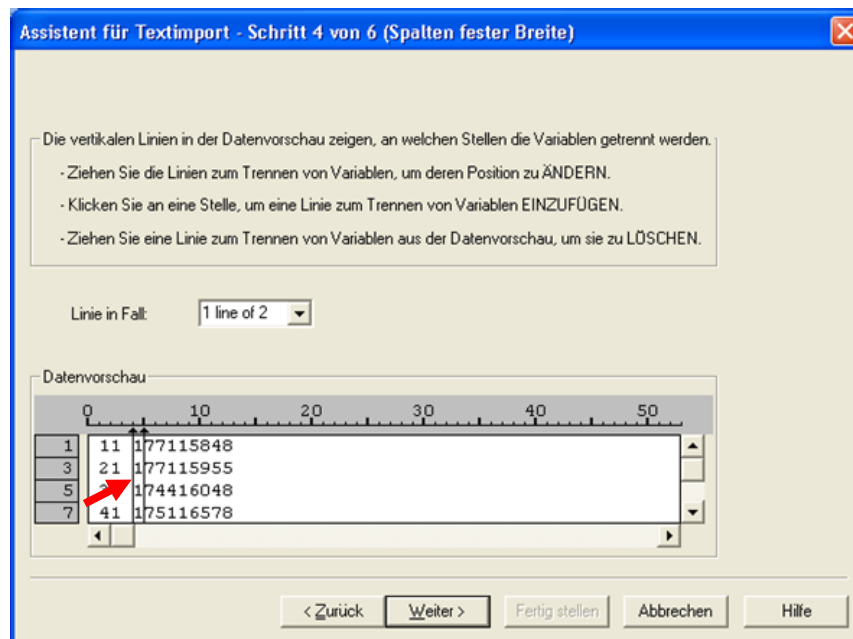
Da unsere Datei keinen Vorspann enthält, **befindet sich der erste Fall** in Zeile 1. Allerdings befindet er sich dort nicht komplett, weil jeweils zwei **Zeilen einen Fall darstellen**:



## Schritt 4

Nun müssen wir die Positionen der einzulesenden Variablen festlegen, wobei der Assistent nur wenig Hilfestellung geben kann, wenn Variablen nicht separiert sind.

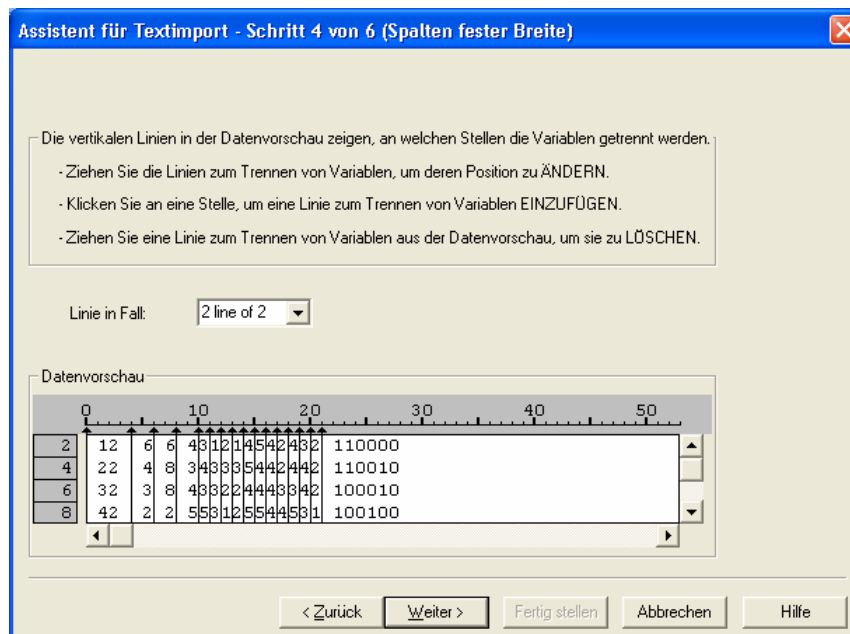
Von der ersten Datenzeile (**1 line of 2** in der Aufklappliste **Linie in Fall**) benötigen wir nur die Variable GESCHL, die wir durch zwei senkrechte Linien abgrenzen:



Hinweise zur Benutzung der Trennlinien:

- Neue Trennlinie einfügen  
Klicken Sie innerhalb der Datenzone auf die gewünschte Spaltenposition (siehe Pfeil in obigem Bildschirmphoto).
- Trennlinie verschieben  
Klicken Sie innerhalb der Datenzone auf die Trennlinie und verschieben Sie diese bei fest gehaltener Maustaste.
- Trennlinie löschen  
Klicken Sie auf das Dreieck an der Spitze der Trennlinie.

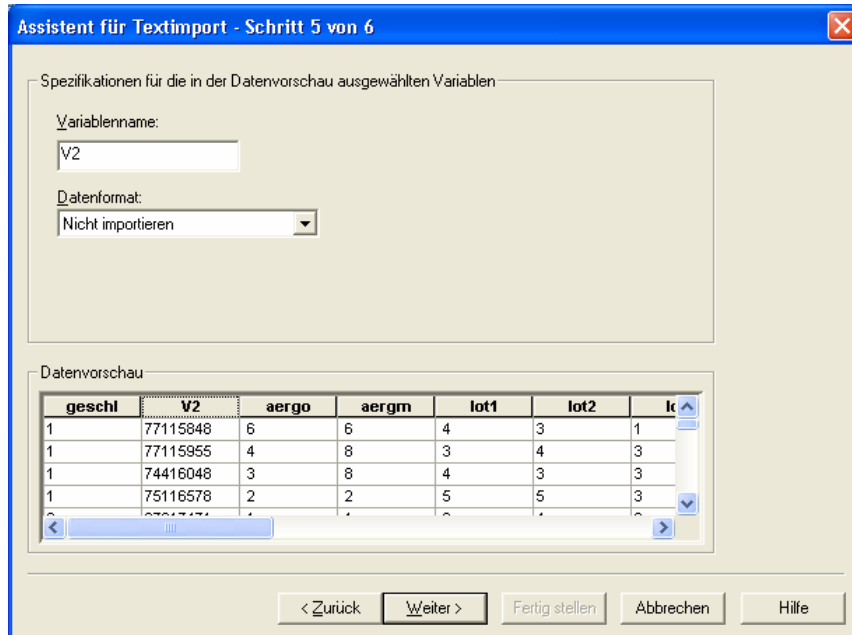
Auf der zweiten Datenzeile benötigen wir erheblich mehr Trennlinien:



## Schritt 5

Im fünften Assistentenschritt können wir die von SPSS vorgeschlagenen Variablennamen ändern und ein **Datenformat** festlegen. Zum Umbenennen ist jeweils genau eine Spalte zu markieren. Das Datenformat lässt sich auch für eine markierte Variablenliste wählen.

Mit dem speziellen Datenformat **Nicht importieren** können überflüssige Variablen ausgeschlossen werden:



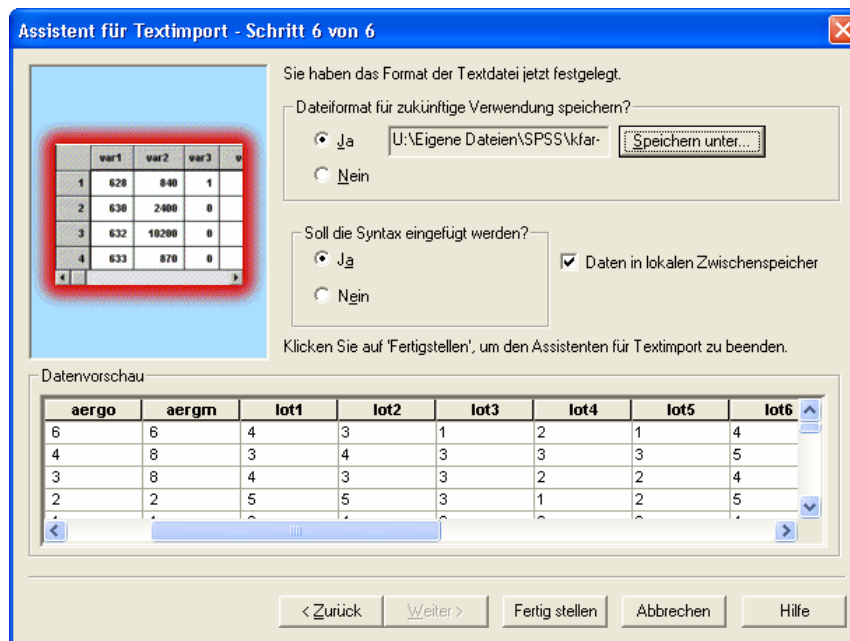
Zumindest bei den LOT-Variablen ist echte Fleißarbeit zu leisten, so dass wir nach Schritt 5 noch **weiter** machen, um unsere Arbeit zu konservieren.

## Schritt 6

Der Assistent bietet zwei Möglichkeiten zum Konservieren einer Dateispezifikation:

- **Dateiformat für zukünftige Verwendung speichern?**  
Es entsteht eine Textassistenten-Formatdatei (Erweiterung **.tpf**), die bei einem späteren Assistenteneinsatz im ersten Schritt angegeben werden kann (siehe oben).
- **Soll die Syntax eingefügt werden?**  
Das für den Datenimport verantwortliche GET DATA – Kommando wird in ein Syntaxfenster geschrieben. Es bietet sich an, zusätzliche Kommandos zu ergänzen, z.B. zum Deklarieren von MD-Indikatoren, die in den Textdaten vorhanden sind. Später kann mit Hilfe des entstandenen SPSS-Programms der Import mit allen erforderlichen Zusatzmaßnahmen automatisiert werden.

Es spricht nichts dagegen, beide Konservierungsoptionen zu verwenden:



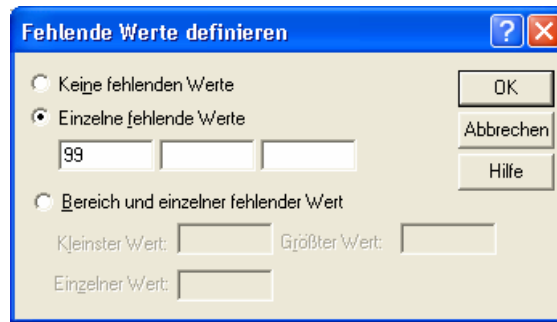
Das vom Textimport-Assistenten erzeugte GET DATA – Kommando verblüfft etwas mit einer Spaltenzählung ab 0:

```
GET DATA  /TYPE = TXT
/FILE = 'U:\Eigene Dateien\SPSS\kfar-kv-pos.txt '
/FIXCASE = 2
/ARRANGEMENT = FIXED
/FIRSTCASE = 1
/IMPORTCASE = ALL
/VARIABLES =
/1  geschl 4-4 F1.0
V2 5-12 8X
/2 aergo 4-5 F2.1
aergm 6-7 F2.1
lot1 8-9 F2.1
lot2 10-10 F1.0
lot3 11-11 F1.0
lot4 12-12 F1.0
lot5 13-13 F1.0
lot6 14-14 F1.0
lot7 15-15 F1.0
lot8 16-16 F1.0
lot9 17-17 F1.0
lot10 18-18 F1.0
lot11 19-19 F1.0
lot12 20-20 F1.0
V18 21-27 7X .
CACHE.
EXECUTE.
```

Nach dem Einlesen einer Textdatei dürfen Sie auf keinen Fall die Deklaration der dort eventuell verwendeten **MD-Indikatoren** vergessen. Studieren Sie also sorgfältig den hoffentlich vorhandenen Kodierplan, der in unserem Fall vorschreibt:

Variable	MD-Indikator
GESCHL	9
AERGO	99
AERGM	99
LOT1-LOT12	9

Die Deklaration kann in der Variablenansicht des Dateneditors erfolgen (siehe Abschnitt 3.2.2). Bei der Variablen AERGO ist z.B. für die Spalte **Fehlende Werte** einzutragen:



Das Kommando MISSING VALUES erlaubt allerdings eine rationellere (und automatisierbare) MD-Deklaration:

```
missing values geschl (9) /aergo aergm (99) /lot1 to lot12 (9).
```

### 13.2 Import von separierten Daten Textdaten

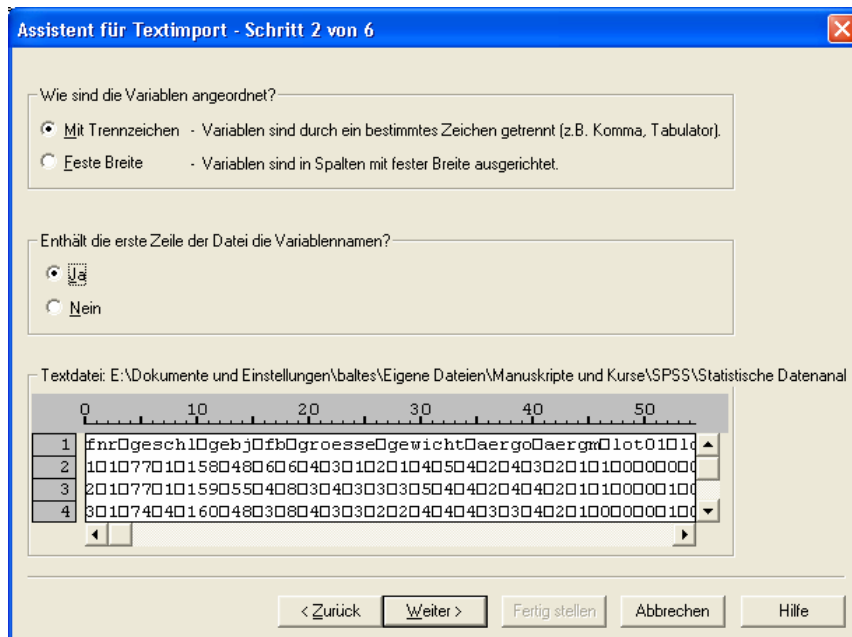
Separierte Textdaten lassen sich erheblich bequemer importieren als positionierte, zumal sie üblicherweise durch eine Zeile mit den Variablennamen eingeleitet werden. Die Datei **kfar-kv-sep.txt** enthält dieselben KFA-Daten, die in Abschnitt 0 aus einer positionierten Datei gelesen wurden:

FNR	GESCHL	GEBJ	FB	GROESSE	GEWICHT	AERGO	AERGM	LOT1	LOT2	...
1	1	77	1	158	48	6	6	4	3	...
2	1	77	1	159	55	4	8	3	4	...
3	1	74	4	160	48	3	8	4	3	...
4	1	75	1	165	78	2	2	5	5	...
.	.	.	.	.	.	.	.	.	.	...
.	.	.	.	.	.	.	.	.	.	...

Beim Import der separierten KFA-Textdaten informieren wir den über

#### Datei > Textdaten lesen

gestarteten Assistenten im zweiten Schritt darüber, dass **Trennzeichen** für Ordnung in der Datei sorgen, und dass die erste Zeile die **Variablennamen** enthält:





### Schritt 3

Der erste Fall befindet sich in der zweiten Zeile der Datei (hinter der einleitenden Zeile mit den Variablennamen). Jeder Fall belegt genau eine Zeile:

Assistent für Textimport - Schritt 3 von 6 (Trennzeichen)

In welcher Zeile befindet sich der erste Fall in den Daten?

Wie sind die Fälle dargestellt?

- Jede Zeile stellt einen Fall dar
- Folgende Anzahl von Variablen stellt einen Fall dar:

Wie viele Fälle sollen importiert werden?

- Alle Fälle
- Die ersten  Fälle.
- Zufälliger Prozentwert der Fälle (ungefähr):  %

Datenvorschau

	0	10	20	30	40	50
1	1	0	1	0	7	7
2	2	0	1	0	7	7
3	3	0	1	0	7	7
4	4	0	1	0	7	7

< Zurück Weiter > Fertig stellen Abbrechen Hilfe

### Schritt 4

In der Datei **kfar-kv-sep.txt** kommt als Trennzeichen nur der **Tabulator** zum Einsatz:

Assistent für Textimport - Schritt 4 von 6 (Trennzeichen)

Welches Zeichen trennt die Variablen?

- Tabulator
- Leerzeichen
- Komma
- Semikolon
- Anderes:

Was ist das Texterkennungszeichen?

- Keins
- Hochkommata
- Anführungszeichen
- Anderes:

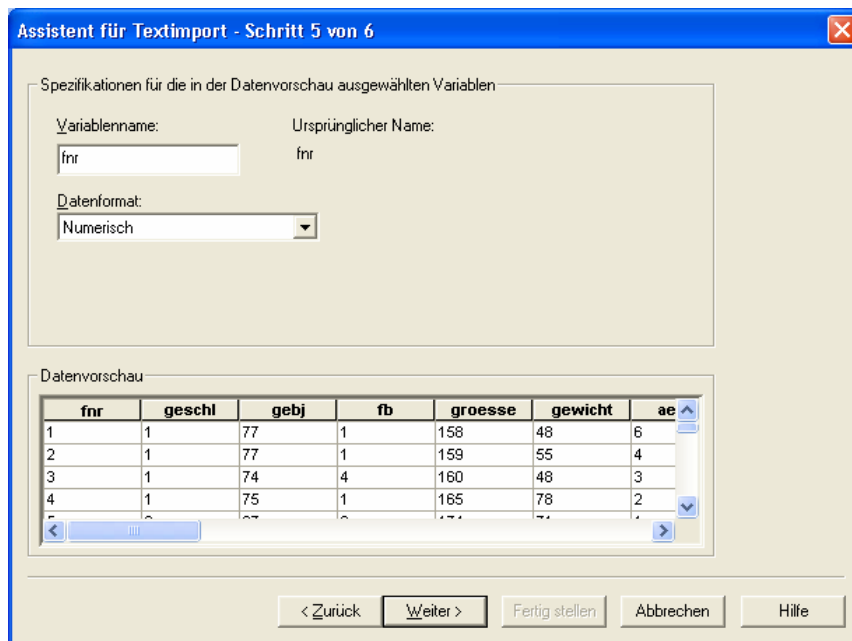
Datenvorschau

	fnr	geschl	gebj	fb	groesse	gewicht	ae
1	1	1	77	1	158	48	6
2	1	1	77	1	159	55	4
3	1	1	74	4	160	48	3
4	1	1	75	1	165	78	2

< Zurück Weiter > Fertig stellen Abbrechen Hilfe

### Schritt 5

Im fünften Assistentenschritt müssen wir nur prüfen, ob die automatische Erkennung des **Datenformats** erfolgreich war:



### Schritt 6

Im letzten Assistentendialog werden die schon in Abschnitt 0 vorstellten Optionen zum Konservieren der Importspezifikation angeboten.

Auch nach dem Einlesen von separierten Textdaten dürfen Sie auf keinen Fall die Deklaration der eventuell vorhandenen **MD-Indikatoren** vergessen.

### 13.3 Überprüfung der revidierten differentialpsychologischen Hypothese

Um mit den in Abschnitt 13.1 bzw. Abschnitt 13.2 importierten Daten die revidierte differentialpsychologische Hypothese prüfen zu können, sind zunächst einige Datentransformationen erforderlich, wobei wir uns die erforderlichen Kommandos teilweise aus dem Transformationsprogramm **kfat.sps** besorgen können:

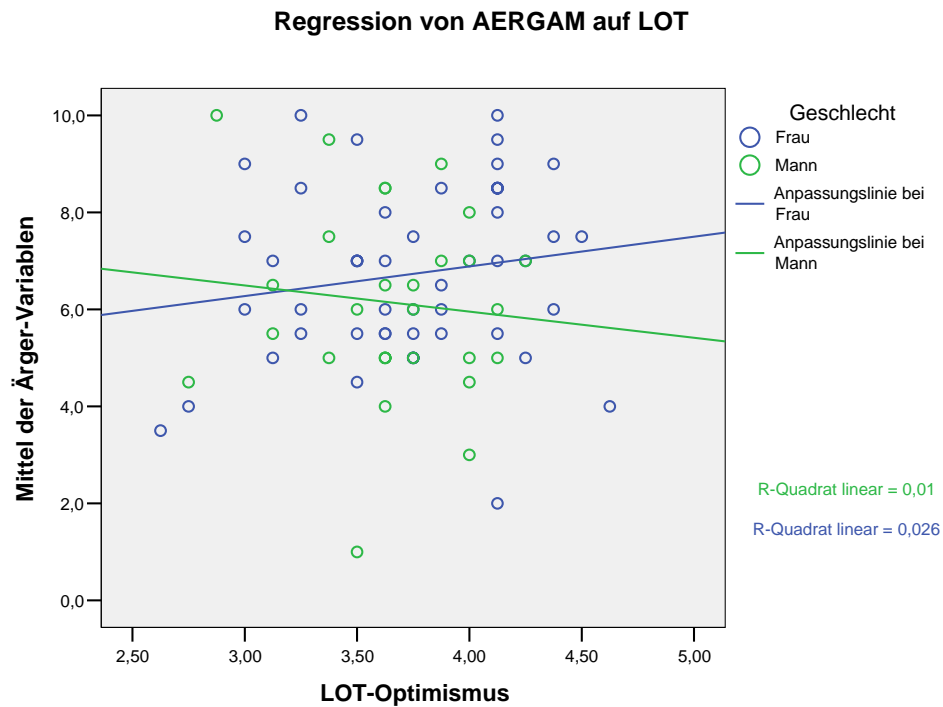
```
* LOT-Fragen umkodieren.
RECODE
  lot3 lot4 lot5 lot12 (5=1) (4=2) (2=4) (1=5) .
EXECUTE .

* LOT berechnen.
COMPUTE lot = MEAN.6(lot1,lot3,lot4,lot5,lot8,lot9,lot11,lot12) .
VARIABLE LABELS lot 'LOT-Optimismus' .
EXECUTE .

* AERGAM berechnen.
COMPUTE aergam = (aergo + aergm)/2 .
VARIABLE LABELS aergam 'Mittel der Ärger-Variablen' .
EXECUTE .

* Produktvariable für die Moderatorhypothese.
COMPUTE geslot = geschl * lot.
VARIABLE LABELS geslot 'GESCHL * LOT'.
EXECUTE .
```

Auch in der neuen Stichprobe scheint das Geschlecht die Regression von AERGAM auf LOT im erwarteten Sinn zu moderieren:



Allerdings wird der Interaktionseffekt *nicht* signifikant ( $p = 0,307$ ):

**Koeffizienten<sup>a</sup>**

Modell		Nichtstandardisierte Koeffizienten		Standardisierte Koeffizienten	T	Signifikanz
		B	Standardfehler	Beta		
1	(Konstante)	,773	5,562		,139	,890
	Geschlecht	3,670	4,130	,949	,889	,377
	LOT-Optimismus	1,761	1,493	,413	1,180	,242
	GESCHL * LOT	-1,150	1,118	-1,120	-1,029	,307

a. Abhängige Variable: Mittel der Ärger-Variablen

Weitere Versuche zur Rettung der differentialpsychologischen Hypothese könnten sich z.B. auf eventuelle Mängel bei der Operationalisierung der theoretischen Begriffe (Ärger und Optimismus) konzentrieren.

Allerdings muss auch die theoretische Fundierung kritisch hinterfragt werden.

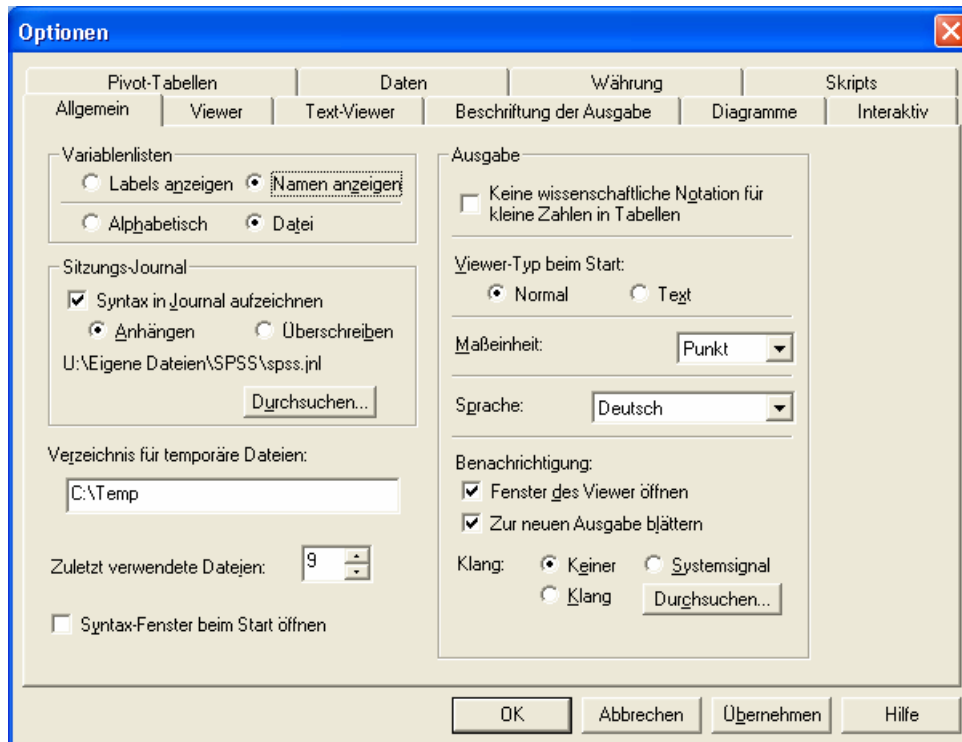
## 14 Einstellungen modifizieren

Das Standardverhalten von SPSS für Windows lässt sich auf vielfältige Weise den individuellen Bedürfnissen anpassen, was wir bei passender Gelegenheit auch schon getan haben.

Über den Menübefehl

### **Bearbeiten > Optionen**

erhalten Sie die folgende Dialogbox mit Optionen zur SPSS-Konfiguration:



Auf dem Registerblatt **Allgemein** sind u.a. folgende Optionen von Relevanz:

### **Variablenlisten**

Bei den Listen auswählbarer Variablen in Dialogboxen verwendet SPSS folgende Voreinstellungen:

- SPSS präsentiert die Variablen durch ihre Labels (falls vorhanden). Dabei werden die Variablenlisten aufgrund des begrenzten Platzangebotes oft recht unübersichtlich. Ein 50-stelliges Label, das auf ca. 20 Zeichen gekürzt werden musste, ist in der Regel weniger informativ als der vollständig sichtbare Variablenname. Mit der Option **Namen anzeigen** im Bereich **Variablenlisten** kann man auf die kompaktere Darstellung umschalten.
- Die Variablen sind angeordnet wie in der Arbeitsdatei, was in der Regel ein bequemes Arbeiten erlaubt. Gemeinsam zu analysierende und damit in Dialogboxen auszuwählende Variablen stehen nämlich oft in der Arbeitsdatei hintereinander. Bei der Arbeit mit einer unbekanntem Datendatei findet man (namentlich bekannte) Variablen jedoch leichter bei alphanumerischer Sortierung. Im Rahmen **Variablenlisten** kann bei Bedarf das Sortierkriterium gewechselt werden.

## Sitzungs-Journal

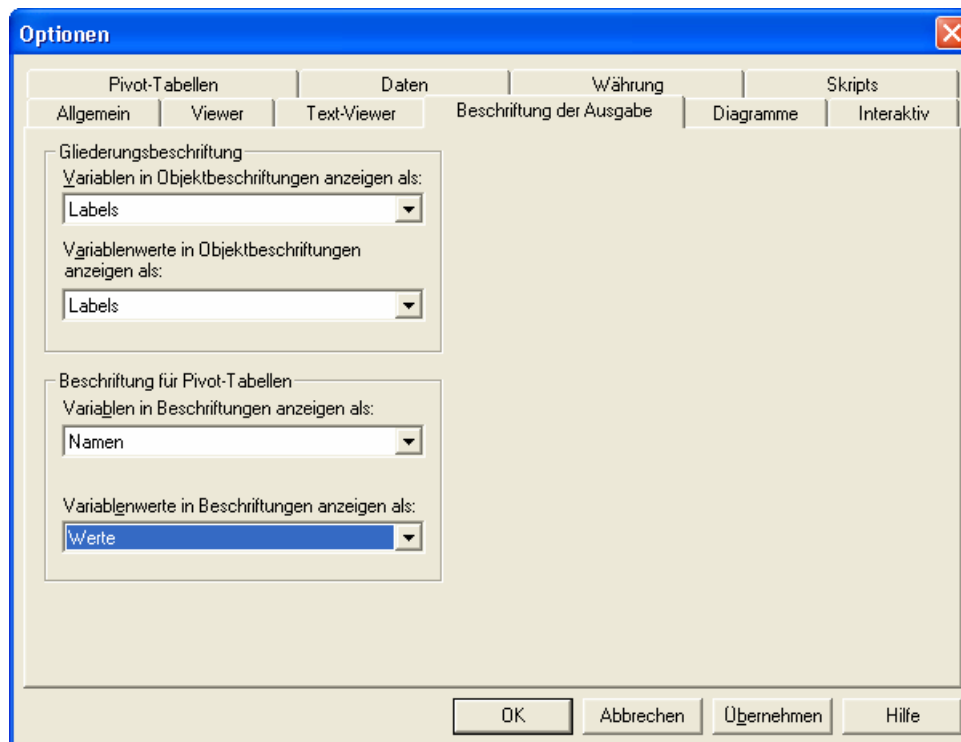
Per Voreinstellung protokolliert SPSS alle Kommandos, die Sie während einer Sitzung per Dialogbox oder via Syntaxfenster abschicken, in einer so genannten **Journaldatei**. Bei den Pool-PCs an der Universität Trier ist dies in der Regel die Datei:

**U:\Eigene Dateien\SPSS\spss.jnl**

Diese Journaldatei kann für Anwender(innen) mit „Mut zur SPSS-Syntax“ z.B. nach einem SPSS-Programmabsturz sehr nützlich sein, weil sie die Kommando-Äquivalente zu praktisch allen Arbeiten der verunglückten Sitzung enthält.

Per Voreinstellung wird beim Start einer SPSS-Sitzung eine vorhandene Journaldatei *nicht* überschrieben, sondern die neuen Kommandos werden am Ende angehängt. Falls die Datei zu groß wird, muss sie gelegentlich verkleinert oder gelöscht werden. Man kann aber auch im Rahmen **Sitzungs-Journal** der Karteikarte **Allgemein** den voreingestellten Öffnungsmodus **Anhängen** abändern auf **Überschreiben**. Dann wird die Journaldatei zu Beginn jeder Sitzung neu erstellt, wobei gegebenenfalls der alte Inhalt überschrieben wird.

Auf dem Registerblatt **Beschriftung der Anzeige** können Sie z.B. veranlassen, dass in Pivot-Tabellen vorhandene Wertelabels ignoriert und stattdessen die Werte selbst angezeigt werden:



## 15 Anhang

### 15.1 Weitere Hinweise zur SPSS-Kommandosprache

In Abschnitt 5 wurden nur sehr oberflächliche Hinweise zur SPSS-Kommandosprache gegeben. Diese sollten genügen für Anwender(innen), die nicht frei programmieren, sondern nur gelegentlich ein von SPSS automatisch erzeugtes Kommando modifizieren wollen.

Der aktuelle Abschnitt ist für ambitionierte Anwender(innen) gedacht, die bereit sind, SPSS-Programme zu schreiben, ...

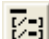
- um auch die ausschließlich per Syntax verfügbaren SPSS-Leistungen nutzen zu können,
- um rationeller mit SPSS zu arbeiten.

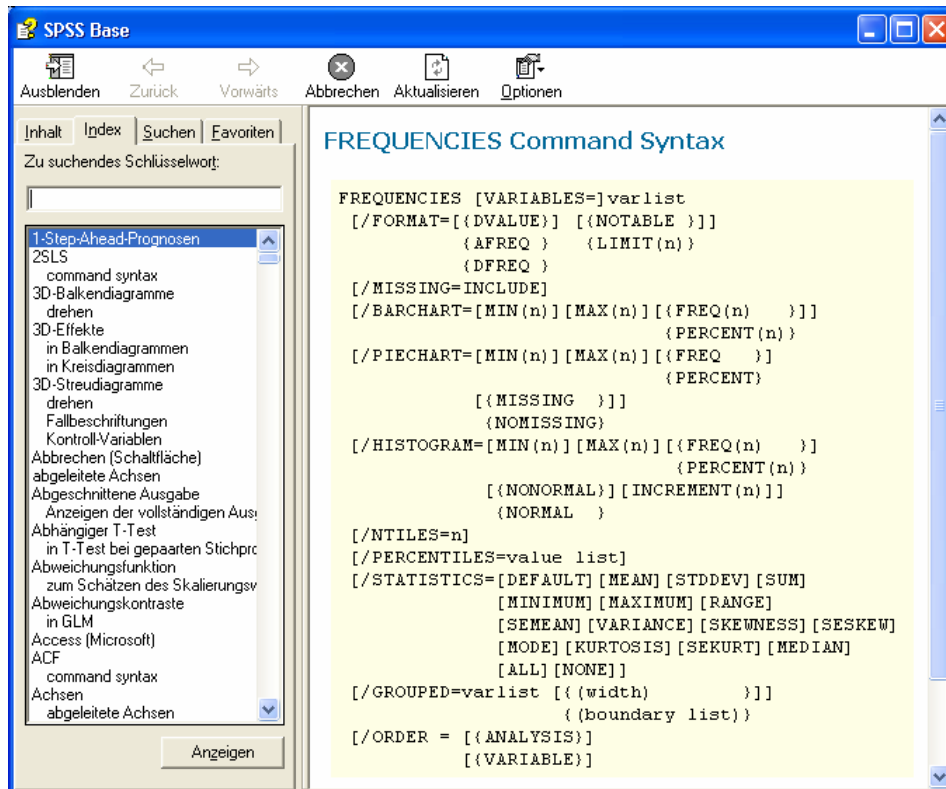
#### 15.1.1 Hilfsmittel für das Arbeiten mit der SPSS-Kommandosprache

Das wichtigste Hilfsmittel für das Arbeiten mit der SPSS-Kommandosprache ist die *Command Syntax Reference*, die als PDF-Dokument über das Hilfesystem verfügbar ist:

##### Hilfe > Command Syntax Reference

Hier findet man ausführliche Beschreibungen der SPSS-Kommandos mit zahlreichen Beispielen und wertvollen Literaturhinweisen zu den realisierten statistischen Methoden.

Die Syntaxfenster bieten ein einfaches Verfahren, das **Syntaxdiagramm** zu einem konkreten Kommando einzusehen: Setzen Sie die Schreibmarke auf das Kommando, und klicken Sie dann auf das Symbol . Zum FREQUENCIES-Kommando, das der **Häufigkeiten**-Dialogbox zugrunde liegt, erscheint z.B. das folgende Hilfefenster:



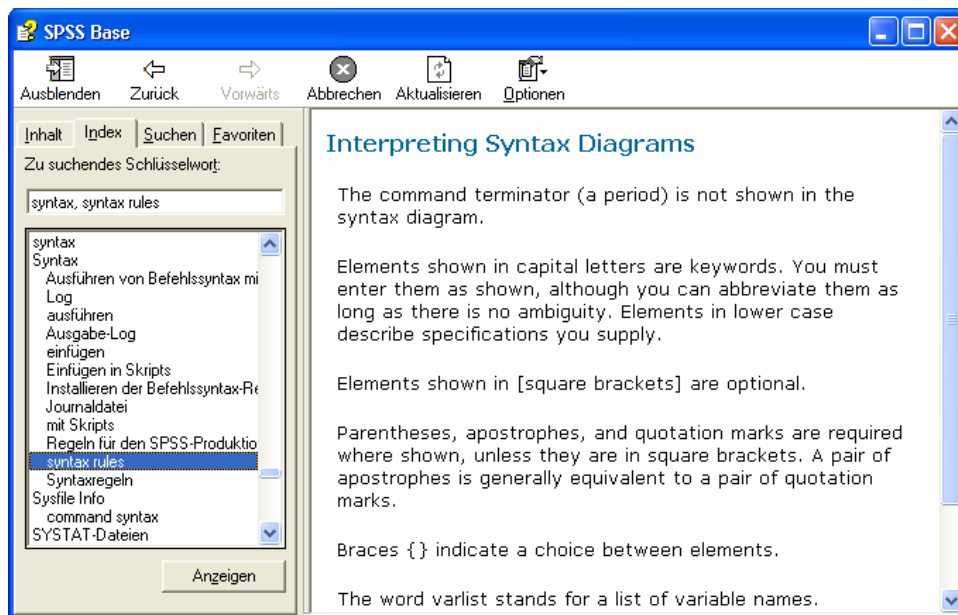
#### 15.1.2 Interpretation von Syntaxdiagrammen

Mit dem Syntaxdiagramm wird die allgemeine Form eines Kommandos definiert und somit festgelegt, wie konkrete Beispiele gebildet werden müssen. Solche Syntaxdiagramme werden auch

im weiteren Verlauf dieses Abschnitts benutzt, um Bestandteile der SPSS-Kommandosprache zu erläutern. In den Syntaxdiagrammen treten einige Metazeichen auf (z.B. "[", "{"), die nicht zur Kommandosprache selbst gehören, sondern diese Sprache beschreiben. Die Bedeutung dieser Metazeichen müssen Sie kennen, um Syntaxdiagramme richtig interpretieren zu können. Im Hilfesystem finden Sie eine Erklärung, indem Sie nach

### Hilfe > Themen > Index

den Suchbegriff *syntax* in das aktive Textfeld eintippen und dann einen Doppelklick auf den Eintrag **syntax rules** setzen:



### 15.1.3 Aufbau von SPSS-Programmen

Welche Kommandos SPSS für das Erstellen von Programmen bereithalten muss, ergibt sich aus unseren Zielvorstellungen: Wir möchten SPSS anweisen, unsere empirischen Daten zu lesen, gegebenenfalls aus den gelesenen Variablen interessantere neue Variablen zu berechnen und schließlich statistische Verfahren mit den eingelesenen oder neu erstellten Variablen zu rechnen. Darüber hinaus haben wir gelegentlich Sonderwünsche hinsichtlich der Arbeitsweise von SPSS.

Orientiert an den gerade skizzierten Teilaufgaben kann man die verfügbaren SPSS-Kommandos in folgende Gruppen einteilen:

- **Dateidefinitions-Kommandos**  
Sie dienen zum Einlesen von Daten in die Arbeitsdatei. Als Beispiel haben wir bereits das GET-Kommando kennen gelernt.  
Wenn ein Programm kein Dateidefinitions-Kommando enthält, wenn es also nicht selbst für das Einlesen seiner Daten sorgt, kann es natürlich nur ausgeführt werden, wenn zuvor eine Arbeitsdatei erzeugt worden ist.
- **Transformations-Kommandos**  
Diese Kommandos dienen zur Veränderung oder Neuberechnung von Variablen bzw. zur Auswahl von Fällen für die weitere Verarbeitung.
- **Prozedur-Kommandos**  
Damit werden statistische Analysen, graphische Präsentationen oder Dateibearbeitungen (z.B. Sortieren der Fälle) angefordert. Ein Beispiel ist das FREQUENCIES-Kommando.

- **Dienst-Kommandos**

Damit kann man u.a. die Arbeitsweise von SPSS beeinflussen (z.B. Startwert des Pseudo-zufallszahlengenerators setzen) und verschiedene Informationen anfordern.

In folgendem SPSS-Programm treten Kommandos aus allen Gruppen auf:

<code>comment Größe und Gewicht.</code>	Dienst-Kommando
<code>get file = 'kfar.sav'.</code>	Dateidef.-Kommando
<code>frequencies var = groesse gewicht /statistics = all /histogram = normal.</code>	Prozedur-   Kommando
<code>compute ideal = groesse - 100.</code>	Transformations-   Kommando
<code>t-test pairs = gewicht ideal.</code>	Prozedur-   Kommando

SPSS-Programme können flexibel gestaltet werden:

- Für die Reihenfolge der SPSS-Kommandos gilt lediglich die selbstverständliche Regel, dass auf eine Variable erst dann Bezug genommen werden darf, nachdem sie im Rahmen einer Dateidefinition oder durch ein Transformations-Kommando eingeführt worden ist.
- In einem Programm dürfen beliebig viele Prozedur-Kommandos auftreten. Manche Anwender leben in dem Irrglauben, pro SPSS-Programm sei nur eine einzige Statistik-Prozedur erlaubt, und verstreuen daher zusammenhängende Auswertungen über unübersichtlich viele Mini-Programme. Andere haben den falschen Ehrgeiz, ihr gesamtes Projekt in einem einzigen Programm abzuwickeln, und erstellen dabei ein unpraktisches Monster-Programm mit mehreren hundert Zeilen. Wie so oft im Leben ist auch hier der gesunde Mittelweg zu empfehlen: Für abgrenzbare Aufgabenpakete sollte jeweils ein eigenes Programm erstellt werden (z.B. mit allen Prozeduren zur Datenprüfung).
- Auch *nach* einer Prozedur dürfen Datentransformationen vorgenommen werden.
- Man kann nach einer Prozedur sogar weitermachen mit der Definition einer neuen Arbeitsdatei, welche dann die alte ersetzt.

#### 15.1.4 Aufbau eines einzelnen SPSS-Kommandos

Die wichtigsten Regeln für SPSS-Befehle:

- Ein Kommando besteht aus seinem Namen und den zugehörigen Spezifikationen:

<i>kommandoname spezifikationen</i>
-------------------------------------

- Der **Kommandoname** kann aus einem Wort bestehen oder aus mehreren Wörtern.  
Beispiele:       - FREQUENCIES  
                  - GET DATA



- Die **Spezifikationen** dürfen enthalten:
  - Schlüsselwörter (z.B. VARIABLES)
  - Variablennamen
  - Zahlen
  - Zeichenfolgen (z.B. Variablenlabel)
  - Operatoren (z.B. "+")
  - spezielle Begrenzungszeichen: / ( ) = ' "

Zwischen diesen Elementen ist mindestens ein Leerzeichen erforderlich. Ausnahme:

Die speziellen Begrenzungszeichen, die arithmetischen Operatoren und manche Vergleichsoperatoren (z.B. ">") sind selbstbegrenzend, d.h. davor und danach sind keine Leerzeichen nötig (aber erlaubt).

Statt eines Leerzeichens darf man meist verwenden:

- beliebig viele Leerzeichen,
- ein Komma,
- einen Zeilenwechsel.

Dies ermöglicht eine übersichtliche Programmgestaltung.

- *Innerhalb* eines Kommandos sind keine Leerzeilen erlaubt.
- Jedes Kommando muss in einer neuen Zeile beginnen und mit einem Punkt enden. Die Kommandos müssen dabei keinesfalls in der ersten *Spalte* beginnen, sondern dürfen eingerückt werden. Von dieser Möglichkeit sollte man z.B. bei Schleifen-Konstruktionen Gebrauch machen.  
Beispiel: 

```
do repeat mc=mc001 to mc100.
    compute mc=normal(1).
end repeat.
```

 Hier werden 100 unabhängige, normalverteilte Zufallsvariablen erzeugt. Durch das Einrücken wird deutlich gemacht, dass die COMPUTE-Anweisung innerhalb der DO REPEAT - Schleife steht.
- In SPSS für Windows brauchen Sie keine maximale Länge für Programmzeilen zu beachten. Manche andere SPSS-Versionen, unter denen Ihr Programm möglicherweise auch laufen soll, haben jedoch eine Beschränkung auf 80 Spalten.
- Ein Kommando kann sich über beliebig viele Fortsetzungszeilen erstrecken.
- Die Verwendung von Groß- oder Kleinbuchstaben ist beliebig.
- Schlüsselwörter dürfen meist bis auf die ersten drei Zeichen abgekürzt werden.  
Beispiel: "fre" für "frequencies"
- Bei den meisten Kommandos sind die Spezifikationen in Subkommandos unterteilt. Diese beginnen mit einem Subkommando-Namen, meist gefolgt von einem Gleichheitszeichen, und sind durch Schrägstriche voneinander getrennt.  
Beispiel: 

```
frequencies var=lot01 /format=notable
/statistics=all.
```

Merken Sie sich aus dieser Liste für den Anfang vor allem:

**JEDES KOMMANDO MUSS IN EINER NEUEN ZEILE BEGINNEN UND MIT EINEM PUNKT ENDEN.**

### 15.1.5 Regeln für Variablenlisten

#### 15.1.5.1 Abkürzende Spezifikation einer Serie von Variablen

In Transformations- oder Prozedur-Kommandos soll häufig eine Folge **bereits existierender** und **in der Arbeitsdatei hintereinander liegender** Variablen angesprochen werden. Dies ermöglicht das **aufrufende TO**, dessen Syntax im Folgenden erläutert wird:

```
vara TO varb
```

*vara, varb*                    Namen bereits vorhandener Variablen, wobei *vara* in der Arbeitsdatei vor *varb* stehen muss.

Beispiele:                    - `frequencies var=alter to beruf.`  
                                   Für alle Variablen, die in der Arbeitsdatei von ALTER bis BERUF positioniert sind, werden Häufigkeitstabellen erstellt.  
                                   - `frequencies var=frage1 to frage3.`  
                                   Wenn in der Arbeitsdatei zwischen FRAGE1 und FRAGE3 1500 beliebig benannte Variablen stehen, dann bewirkt dieses Kommando 1502 Häufigkeitstabellen.

#### 15.1.5.2 Der Platzhalter varlist

In folgendem Syntaxdiagramm wird der in SPSS-Kommandos häufig auftretende Platzhalter *varlist* definiert:

```
{varname | varname_1 TO varname_2} [{...}]
```

*varname,*  
*varname\_1,*  
*varname\_2*                    Variablennamen

Beispiel:                    `missing values nieder01 to hoehe ozon mess1 to mess4 (9).`  
                                   Hier wird mit dem MISSING VALUES - Kommando für alle aufgelisteten Variablen die 9 als MD-Indikator vereinbart.

---

## Literaturverzeichnis

- Backhaus, K., Erichson, B., Plinke, W. & Weiber, R. (2006). *Multivariate Analysemethoden* (11. Aufl.). Berlin: Springer.
- Baltes-Götz, B. (2006). *Lineare Regressionsanalyse mit SPSS*. Online-Dokumentation: <http://www.uni-trier.de/urt/user/baltes/docs/linreg/linreg.pdf>
- Bortz, J. (1977). *Lehrbuch der Statistik*. Berlin: Springer.
- Bortz, J. & Döring, N. (1995). *Forschungsmethoden und Evaluation*. Berlin: Springer.
- Cohen, J., Cohen, P., West, S.G. & Aiken, L. (2003). *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences* (3rd ed.). Mahwah: Lawrence Erlbaum Associates.
- Erdfelder, E., Faul, F., & Buchner, A. (1996). GPOWER: A general power analysis program. *Behavior Research Methods, Instruments & Computers*, 28, 1-11.
- Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (in press). G\*Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*.
- Hartung, J. (1989). *Statistik* [7. Auflage]. München: Oldenbourg.
- Kahneman, D. & Miller, D.T. (1986) Norm theory: comparing reality to its alternatives. *Psychological Review*, 93, 136-153.
- Mehta, C.R., Patel, N.R. (1996). *SPSS Exact Tests 7.0 for Windows*. Chicago, IL: SPSS Inc.
- Norušis, M.J. (2005). *SPSS 14.0. Statistical Procedures Companion*. Upper Saddle River, NJ: Prentice Hall.
- Norušis, M.J. (2005). *SPSS 14.0. Advanced Statistical Procedures Companion*. Upper Saddle River, NJ: Prentice Hall.
- Pedhazur, E.J. & Pedhazur Schmelkin L. (1991). *Measurement, design, and analysis. An integrated approach*. Hillsdale, NJ: Lawrence Erlbaum.
- Raudenbush, S. W. & Bryk, A. S. (2002). *Hierarchical Linear Models* (2nd ed.). Thousand Oaks, CA: Sage.
- Scheier, M.F. & Carver, C.S. (1985). Optimism, Coping, Health: Assessment and implications of generalized outcome expectancies. *Health Psychology*, 4, 219-247.
- Schnell, R., Hill, P. B. & Esser, E. (2005). *Methoden der empirischen Sozialforschung* (7. Aufl.). München: Oldenbourg.
- Siegel, S. (1976). *Nichtparametrische statistische Methoden*. Frankfurt: Fachbuchhandlung für Psychologie
- Tabachnik, B.G. & Fidell, L.S. (2007). *Using multivariate statistics* (5<sup>th</sup> ed.). Boston: Pearson.
- Stevens, J. (1996). *Applied Multivariate Statistics for the Social Sciences* (3<sup>rd</sup> ed.). Mahwah: Lawrence Erlbaum.
- Wallis, W.A. & Roberts, H.V. (1956). *Statistics, a new approach*. Glencoe, Ill.: The Free Press.
- Wentura, D. (2004). Ein kleiner Leitfaden zur Teststärke-Analyse. Online-Dokument: <http://www.uni-saarland.de/fak5/excops/download/POWER.pdf>

---

## Stichwortregister

### A

Ablehnungsbereich	110
Achsenteilstriche	134
Alpha-Fehler	3, 7, 109
Alphanumerische Variablen	18
Alternativhypothese	1, 108
AND-Operator	100
Anwärterliste	58
Arbeitsdatei	40, 50, 58
speichern	50
Assistent	
zum Textimport	166
Ausblenden	
von Kategorien	130
Ausgabeblock	61
Ausgabefenster	27, 60, 127
designiertes	75
Mehrere verwenden	75
Neues anfordern	75
Ausreißer	114
Ausrichtung	44
Automatisierte Datenerfassung	34

### B

Balkendiagramm	63
Bedingte Datentransformation	97, 141
Benutzerberatung an der Universität Trier	33
Benutzerschnittstelle	76
Beobachtungseinheit	2
Berechnen	91
Beta-Fehler	4, 7, 110
BMP	74
Boxplot	114

### C

CGM	74
Chi-Quadrat-Statistik	150
COMMENT-Kommando	81
COMPUTE-Kommando	91
COUNT-Kommando	102

### D

Data Entry	39
Dateidefinitions-Kommandos	179
Daten suchen	70
Datendatei	
öffnen	57
Dateneditor	13, 40

Dateneditorfenster	27
Dateneingabe	53
Datenerfassung	34
automatisierte	34
manuelle	22, 36
per Datenbankprogramm	38
per SPSS-Dateneditor	40
per Texteditor	37, 56
Datenfenster	40
Neu	78
Datenmatrix	13, 40
Datenschutz	14
Datensicherheit	83
Datentransformation	5, 82
bedingte	97
Datumsvariablen	18
Deklarationsteil	41
Demographische Merkmale	10
Deskriptive Statistik	1
Dezimalstellen	42
in Pivot-Tabellen	71
Dezimaltrennzeichen	96
Dienst-Kommandos	180
Differentialpsychologische Hypothese	138
Diskriminanzanalyse	18
DO IF - Kommando	165
DO REPEAT - Kommando	165
Drucken	
Viewer-Dokumente	61

### E

Eigenschaftsfenster	134
Einfügen	
Fall	54
Variable	47
Einfügen-Schaltfläche	77
Einfügen-Schaltfläche	76
Einscannen	36
Einseitige Hypothesen	
für (2 × 2)-Tabellen	155
Einstellungen modifizieren	176
Ein-Stichproben-t-Test	97
EMF	74
EPS	74
Erfassungsfehler	56
Exact Tests - Modul	153
Exakte Tests	153
EXECUTE-Kommando	87, 89
Explorative Datenanalyse	114, 115

Explorative Verfahren	1	SQRT	93
Exportieren	73	statistische	93
Exzeß	67	SUM	94
<b>F</b>		UNIFORM	95
Fall	13	VALUE	94
einfügen	54	Fußzeile	62
erschieben	55	<b>G</b>	
löschen	54	Generalisierbarkeit	64
Fälle		GET DATA - Kommando	171
auflisten	141	GET-Kommando	78
ausfiltern	140	Gitterlinien	72
gewichten	157	<b>GlobalPark</b>	34
Fälle auswählen	140	GPower 3	7, 122
Fallidentifikation	14	GRAPH-Kommando	132
Falls-Subdialogbox	97	Gruppeneinteilung	85
Fallstudien	31	Gruppierungen	
Fallweiser Ausschluss fehlender Werte	124	in einer Pivot-Tabelle	128
Fehlende Werte	19, 94	<b>H</b>	
deklarieren	44	Handbücher	32
fallweiser Ausschluss	124	Häufigkeitsanalyse	58, 59
paarweiser Ausschluss	124	Hauptausgabefenster	75
Rechenregeln für ...	96	Hilfesystem	29
Fehler		Homogenitätshypothese	149
erster Art	3, 109	Homoskedastizität	113
zweiter Art	4, 110	HTML	73
Fertigdatendatei	52, 82	Hypothesen	2, 3
Festes Format	37	Hypothesentests	1, 108
Filter	140, 141	<b>I</b>	
Filterfragen	39	ICR	36
Fishers exakter Test	111, 155	IGRAPH-Kommando	132
Fokus		Inferenzstatistik	1, 108
im Ausgabefenster	61	Initialisierung numerischer Variablen	84
FORMATS-Kommando	106	INPUT II	40, 57
FREQUENCIES-Kommando	76, 78	Internet	33, 34
Funktionen	93	Intervallschätzungen	1
ABS	93	Intervallskalenqualität	6
arithmetische	93	<b>J</b>	
EXP	93	Journaldatei	177
für fehlende Werte	94	JPG	74
LG10	93	<b>K</b>	
LN	93	Kategorien	
MAX	93	ausblenden	130
MEAN	93	KFA-Hypothese	6
MIN	93	Kodierplan	4, 13, 24
MOD	93	Kodierung	4, 13, 18
NMISS	94	Kolmogorov-Smirnov - Test	116, 118
NORMAL	95		
Pseudozufallszahlengeneratoren	95		
RND	93		
SD	94		

Kommandosprache	76, 80, 164, 178	Nominalskala	143
Kommentare in SPSS-Programmen	81, 106	Nominalskalenniveau	18
Konfirmatorische Verfahren	1	Normalitätsannahme	113
Kontinuitätskorrektur nach Yates	156	Normalverteilungsannahme	112
Kopfzeile	62	Normalverteilungsannahme	118
Kreuztabellen	143	Normalverteilungstests	116, 118
Kritischer Wert	109	NOT-Operator	100
Künstliche Gruppenbildung	85	Nullhypothese	1, 108
Kurtosis	67	Numerische Funktionen	<i>Siehe Funktionen</i>
<b>L</b>		Numerische Variablen	18
Leerzeilen	106	Numerischer Ausdruck	93
Lernprogramm	30	Auswertungsprioritäten	95
Life Orientation Test	9	<b>O</b>	
Likelihood-Quotienten-Test	für	OCR	36
Kreuztabellen	152	Offene Fragen	17
Linearitätsannahme	112	dynamisches Set aus kateg. Variablen	17
Logischer Ausdruck	99, 100, 140	Offene Transformationen	90
Abarbeitungsreihenfolge	101	Öffnen	
unbestimmter	99	Datendatei	57
Wahrheitstafeln	100	Viewer-Dokumente	62
Logischer Operator	100	OMR	36
Löschen		Online-Datenerhebung	34
Fall	54	Operationalisierung	3, 6
Variable	48	Ordinalskalenniveau	18
LOT	89	Ordinatenabschnitt	113
<b>M</b>		OR-Operator	100
Macintosh	23	<b>P</b>	
Mantel-Haenszel-Statistik	153	Paarweiser Ausschluss fehlender Werte	124
MD-Indikator	19	PCT	74
Mehrfachantworten-Set	15, 16	PDF-Export	62, 74
Mehrfachwahl		Pearsons Chi-Quadrat-Statistik	150
Häufigkeiten	159	Pivot-Editor	70, 127
Kreuztabellen	161	Plausibilitätsprüfungen	39
Mehrfachwahlfragen	158	PNG	74
sparsames Set aus kateg. Variablen	16	Population	1
vollständiges Set aus dichot. Variablen	15	Positionierte Daten	37, 166
Mehrfachwahl-Fragen	15	Positiv semidefinit	124
Mehrfachwahl-Set		Power	111
definieren	159	t-Test für die Pearson-Korrelation	122
speichern	161	Poweranalyse	
Menüzeile	28	Post hoc	122
Messniveau	44	Programm-orientierte Arbeitsweise	77
MISSING VALUES - Kommando	172	Prozedur-Kommandos	179
Missing-Data-Indikator	19	Prüfstatistik	108, 151
Moderatoreffekt	137	Pseudozufallszahlengenerator	95
<b>N</b>		<b>R</b>	
Navigationsbereich	61, 73	Ratingskalen	6
NMISS	103	RECODE-Kommando	85

Regressionsanalyse	120, 121	Stichprobenumfang	7
Repräsentativität der Stichprobe	143	String-Variablen	18
Rohdatendatei	51, 82	Strukturierung	4, 13, 14
Rückgängig-Befehl im Datenfenster	55	Subkommando	181
<b>S</b>		Suchen	
SamplePower	7	Begriffe	29
SAV-Dateien	51	Daten	70
SAVE-Kommando	104	Symbolleisten	28
SCALE	107	Syntaxdiagramm	178
Schätzmethoden	1	Syntaxfenster	76, 80, 178
Schiefe	66	aktivieren	80
Schreibschutz	83	designiertes	80
SEED	95	Kommandos ausführen	79
SELECT IF	91	neu erstellen	80
Separierte Daten	38, 172	öffnen	80
Shapiro-Wilk - Test	116	schließen	80
Shapiro-Wilk - Test	118	speichern	79
Skalenniveau	3, 18, 44	Syntax-Regeln	80
Sortierung bei Variablenlisten	176	SYSMIS	19, 53, 54, 70, 96
Spaltenbreite	72	Systemdefiniert fehlend	19
Spaltenformat	42	System-Missing	19, 87
Speichern		<b>T</b>	
Arbeitsdatei	50	Tabellenvorlagen	72
Syntax	79	Teilausgabe	61
Viewer-Dokumente	62	Teilnehmerliste	58
SPSS		Teleform	36
Kommandosprache	164	Testproblem	
Lizenzen	26	zweiseitiges	111
Mietlizenzen	26	Teststärke	111, 122
Module	26	t-Test für die Pearson-Korrelation	122
SPSS-		Textdatendateien	166
Benutzerschnittstelle	76	Textimport-Assistent	166
SPSS-		TIF	74
Prozessor	76	TO	94
SPSS-		TO-Schlüsselwort	182
Syntax	80	Transformations-Kommandos	179
SPSS im Internet	33	Transformationsprogramm	52, 76, 82, 104
SPSS-Datendatei	50	Transformieren	
SPSS-Kommandosprache	76, 80	Berechnen	91
SPSS-Programm	52, 76	Umkodieren	85
dialogunterstützte Erstellung	77	Zählen	102
SPSS-Usenet-Diskussionsgruppe	33	t-Test	
Standardfehler		für abhängige Stichproben	7
der Schiefe	67	für eine Stichprobe	97
Startassistent	27	für gepaarte Stichproben	109, 112
Statistik-Assistent	31	t-Verteilung	109
StatTransfer	37	<b>U</b>	
Statuszeile	28	Überschreitungswahrscheinlichkeit	109
Stichprobe	4	Umkodieren	85
Stichprobenmodell	108, 150		

Umlaute		Vergleich	100
in Variablenamen	23	Verschieben	
Unabhängigkeit	108	Fall	55
von Residuen	3	Variable	48
Unabhängigkeit der Residuen	113	Versuchsplanung	3
Unabhängigkeitshypothese	149	Verteilungsfreier Lagevergleich	118
Undo-Funktion im Datenfenster	55	Viewer	27, 60, 127
Untersuchungsdesign	3	Vorlagen	
Untersuchungsplanung	2, 6	Grafiken	138
<b>V</b>		Vorzeichentest	118, 125
Variable	13	<b>W</b>	
einfügen	47	Wahrheitstafeln	100
löschen	48	Wahrheitswert	100
verschieben	48	Wertelabels	43, 46
VARIABLE LABELS - Kommando	87	WMF	74
Variablen		WRITE-Kommando	37
abgeleitete	15	<b>Z</b>	
Variablenattribute	42	Zählen von Werten	102
Variablendefinition	41	Zelleneigenschaften	71
Variablenlabel	42	Zellenmarkierung	53
Variablenlisten	176, 182	Zufällige Teilstichprobe ziehen	141
Variablenamen	14, 23	Zufallszahlengenerator	95
Variablentypen	18, 42	Zweiseitiges Testproblem	111
Varianzhomogenität	113	Zwischenablage	62
Varlist	182		
Verfälschter Test	111		