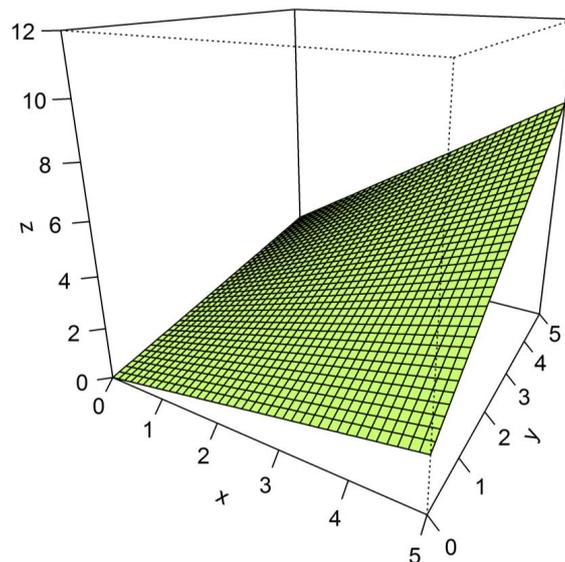




B. Baltes-Götz

R als Ergänzung zu SPSS



Herausgeber: Zentrum für Informations-, Medien- und Kommunikationstechnologie (ZIMK)
an der Universität Trier
Universitätsring 15
D-54286 Trier
WWW: zimk.uni-trier.de

Copyright: © ZIMK 2017
Autor: Bernhard Baltes-Götz (E-Mail : baltes@uni-trier.de)

Vorwort

In diesem Manuskript geht es um die Nutzung der freien Statistik-Entwicklungsumgebung **R** als Erweiterung zu IBM SPSS Statistics (ab jetzt kurz bezeichnet als *SPSS*). Wer die Benutzerfreundlichkeit und Funktionsvielfalt von SPSS schätzt, aber auch in **R** realisierte Lösungen nutzen und vielleicht sogar eigene Lösungen in **R** entwickeln möchte, kann sich über die Kooperationsbereitschaft der beiden Programme freuen und hat ein leistungsfähiges Gespann zur Verfügung. Im Manuskript werden SPSS 24 und **R** 3.2 unter Windows 10 verwendet.

Die aktuelle Version des Manuskripts ist als PDF-Dokument zusammen mit den im Kurs benutzten Dateien auf dem Webserver der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend folgendermaßen zu finden:

[IT-Services \(ZIMK\) > Downloads & Broschüren >
Statistik > R als Ergänzung zu SPSS](#)

Kritik und Verbesserungsvorschläge zum Manuskript werden dankbar entgegen genommen (z.B. unter der Mail-Adresse baltes@uni-trier.de).

Trier, August 2017

Bernhard Baltes-Götz

Inhaltsverzeichnis

1	EINLEITUNG	8
2	SPSS-ERWEITERUNGEN AUF R- UND/ODER PYTHON-BASIS INSTALLIEREN	11
2.1	Python- und R-Essentials	11
2.1.1	Python-Essentials	11
2.1.2	R-Essentials	11
2.1.3	Installationsarbeiten im Überblick	12
2.2	Erweiterungsbundles	12
2.2.1	Inhalt	12
2.2.2	Erstellung	13
2.2.3	Installation	13
2.3	SPD-Dateien	16
3	EXEMPLARISCHE BESCHREIBUNG EINIGER ERWEITERUNGSBUNDLES	18
3.1	Robuste Regression	18
3.2	Breusch-Pagan - Heteroskedastizitäts-Test	22
3.3	Polychorische und polyserielle Korrelationen	23
3.4	Tobit-Regression	25
4	R-FUNKTIONEN ÜBER DAS SPSS-SYNTAXFENSTER NUTZEN	28
4.1	SPSS-Variablen an R übergeben	28
4.1.1	Übergabe der kompletten Arbeitsdatei	28
4.1.2	Eine Auswahl von SPSS-Variablen übergeben	30
4.1.3	Variablen in einer R-Datentabelle ansprechen	30
4.1.4	Persistenz und Löschen von R-Objekten	31
4.1.5	Kategoriale SPSS-Variablen als (ordinale) Faktoren an R übergeben	31
4.1.6	Indikatoren für fehlende Werte	32
4.2	R-Auswertungsfunktionen verwenden und Ausgaben im SPSS-Viewer anzeigen	33
4.3	SPSS-Variablen mit R erstellen	35
5	R ALS STATISTIKORIENTIERTE PROGRAMMIERUMGEBUNG	38
5.1	RGui zur direkten Interaktion mit R	38
5.1.1	Arbeitsverzeichnis	39
5.1.2	Workspace und Anweisungsgedächtnis	40
5.1.3	Sichern und Laden einzelner Datenobjekte im Binärformat von R	41
5.1.4	Konfigurationsoptionen	41
5.1.5	Initialisierungsdateien	42

5.2 Pakete	43
5.2.1 Pakete laden	43
5.2.2 Pakete installieren	44
5.2.3 Installierte Pakete aktualisieren.....	46
5.2.4 Task Views.....	47
5.2.5 Pakete entladen.....	47
5.2.6 Pakete zitieren	47
5.3 Elementare Eigenschaften der Programmiersprache R	47
5.3.1 Hilfe und Dokumentation.....	47
5.3.1.1 Hilfe aufrufen	47
5.3.1.2 Beispiele aus den Hilfetexten ausführen lassen	49
5.3.1.3 Elektronische Handbücher	49
5.3.2 Bezeichner und Kommentare	50
5.3.3 Funktionen	51
5.3.3.1 Regeln für den Aufruf von Funktionen.....	51
5.3.3.2 Elementare Funktionen	52
5.3.4 Datentypen	54
5.3.4.1 Datentypbezogene Funktionen	54
5.3.4.2 Vektor	55
5.3.4.3 Faktor.....	59
5.3.4.4 Matrix	60
5.3.4.5 Array	63
5.3.4.6 Liste	64
5.3.4.7 Datentabelle	66
5.3.5 Fehlende Werte	70
5.3.6 Indezzugriff.....	72
5.3.6.1 Zugriff auf einzelne Elemente	72
5.3.6.2 Zugriff auf einen Zeilen- oder Spaltenvektor aus einer Matrix oder Datentabelle	73
5.3.6.3 Indexvektoren	73
5.3.6.4 Indexmatrizen	75
5.3.7 Operatoren.....	75
5.3.7.1 Arithmetische Operatoren.....	75
5.3.7.2 Vergleichsoperatoren	76
5.3.7.3 Logische Operatoren.....	76
5.3.7.4 Sequenzoperator	77
5.3.7.5 Recycling-Regel	77
5.3.7.6 Zuweisungsoperatoren	78
5.3.8 Anweisungen.....	78
5.3.8.1 Blockanweisung.....	79
5.3.8.2 if - Anweisungen	79
5.3.8.3 if-else - Anweisung.....	79
5.3.8.4 Wiederholungsanweisungen	79
5.4 Mit Skripten arbeiten	80
5.5 Generische Funktionen und Ausgabenverwaltung	81
5.6 Eigene Funktionen	82
6 BEDIENUNGSERLEICHTERUNGEN FÜR R	84
6.1 Dateneditor	84
6.2 R Commander	86
6.2.1 Datentabelle anlegen, definieren und füllen	86
6.2.2 Datenverwaltung	89

7	DATENVERWALTUNG UND -TRANSFORMATION MIT R	91
7.1	Beispieldaten in R-Paketen nutzen.....	91
7.2	Daten in Fremdformaten lesen und schreiben	92
7.2.1	Textdateien mit separierten Daten.....	92
7.2.1.1	Lesen.....	93
7.2.1.2	Schreiben	96
7.2.2	SPSS-Datendateien lesen	97
7.2.3	Datenexport an SPSS	98
7.2.4	Dateiauswahl per Dialogbox	99
7.3	Variablen erstellen oder modifizieren.....	99
7.3.1	Umkodieren.....	99
7.3.1.1	Lösung per recode().....	100
7.3.1.2	Lösung mit logischen Indexvektoren.....	100
7.3.1.3	Lösung mit ifelse().....	100
7.3.1.4	Metrische Variable kategorisieren	101
7.3.2	Berechnen.....	101
7.4	Zufallszahlen erzeugen	102
7.4.1.1	Normalverteilte Zufallszahlen	102
7.4.1.2	χ^2 -verteilte Zufallszahlen.....	102
7.4.1.3	Binomialverteilte Zufallszahlen.....	103
7.5	Auswahl von Fällen und/oder Variablen	104
7.5.1	Auswahl von Fällen.....	104
7.5.2	Auswahl von Variablen.....	104
7.6	Detail- und Master-Daten kombinieren.....	105
7.7	Daten aggregieren.....	106
7.8	Sekundärstichproben ziehen.....	106
8	GRAFIK-OPTIONEN IN R	107
8.1	Ausgabeformate und -geräte.....	107
8.1.1	Verfügbare Ausgabegeräte	107
8.1.2	Grafikfenster	108
8.1.3	Ausgabe in eine Datei	109
8.1.4	R-Diagramme an Microsoft Office übergeben.....	110
8.1.5	Ausgabegeräte verwalten	112
8.1.6	R-Diagramme im SPSS-Ausgabefenster	112
8.2	Das traditionelle R-Grafiksystem.....	113
8.2.1	High- und Low-Level - Grafikfunktionen	113
8.2.2	Grafikparameter	114
8.2.3	Beschriftungen	116
8.2.4	Die generische Funktion plot().....	116
8.2.4.1	Argumente	116
8.2.4.2	Ergänzende Low-Level - Grafikfunktionen	119
8.2.5	Wichtige Diagrammtypen	122
8.2.5.1	Liniendiagramm.....	122
8.2.5.2	Streudiagramm.....	122
8.2.5.3	Boxplot	132
8.2.5.4	Histogramm mit Dichteschätzung	134
8.2.5.5	Mehrere Diagramme kombinieren.....	135
8.2.5.6	Funktionsplots	136

8.3	Das Grafikpaket ggplot2	138
8.3.1	Grammatik eines ggplot2-Plots	139
8.3.1.1	Plot-Objekte, Schichten und Geome	139
8.3.1.2	Aesthetics	142
8.3.1.3	Statistische Transformationen	143
8.3.1.4	Skalen, Achsen und Legenden	144
8.3.1.5	Positionsanpassungen	154
8.3.1.6	Facetten	157
8.3.1.7	Themes	158
8.3.2	ggplot2 - Diagramm in eine Datei sichern	161
8.3.3	Streudiagramme	161
8.3.3.1	Plot-Objekt anlegen	161
8.3.3.2	Einfaches Streudiagramm	162
8.3.3.3	Einfaches Streudiagramm mit Konfidenzzone	163
8.3.3.4	Gruppiertes Streudiagramm	163
8.3.3.5	Schichtaufbau mit qplot() starten	166
8.3.3.6	Dichtedarstellung bei großen Stichproben	167
8.3.4	Weitere Diagrammtypen	169
8.3.4.1	Histogramm und Dichteschätzung	169
8.3.4.2	Boxplot	174
8.3.4.3	Balkendiagramme	176
8.3.4.4	Liniendiagramme	185
9	STATISTISCHE DATENANALYSE MIT R	187
9.1	Einfache univariate Verteilungsbeschreibung	187
9.1.1	Univariate Verteilungsbeschreibung für metrische Variablen	187
9.1.1.1	Kompakte Verteilungsbeschreibung	187
9.1.1.2	Statistische Funktionen für numerische Vektoren	187
9.1.2	Absolute und relative Häufigkeiten für kategoriale Variablen ausgeben	188
9.2	Modellformulierung	189
9.3	Lücken im SPSS-Statistikangebot füllen	190
9.3.1	Gerichtete t-Tests und einseitige Vertrauensintervalle	190
9.3.2	Soziale Netzwerkanalyse mit igraph	191
9.3.2.1	Visualisierung	192
9.3.2.2	Analyse	195
10	MATHEMATISCHE PROBLEME MIT R BEWÄLTIGEN	197
10.1	Mengenlehre	197
10.2	Lineare Algebra	198
LITERATUR		199
STICHWORTVERZEICHNIS		201

1 Einleitung

R ist eine freie *Programmierungsumgebung für Datenanalyse und Grafik* (dt. Übersetzung des Untertitels von Venables et al. 2016), die als Implementation der statistik-orientierten Programmiersprache **S** (Chambers 1998) entstanden ist.

Dank einer von der Firma IBM SPSS unter dem Namen **R-Essentials** gelieferten Integrationslösung kann **R** von SPSS aus genutzt werden, z.B. ...

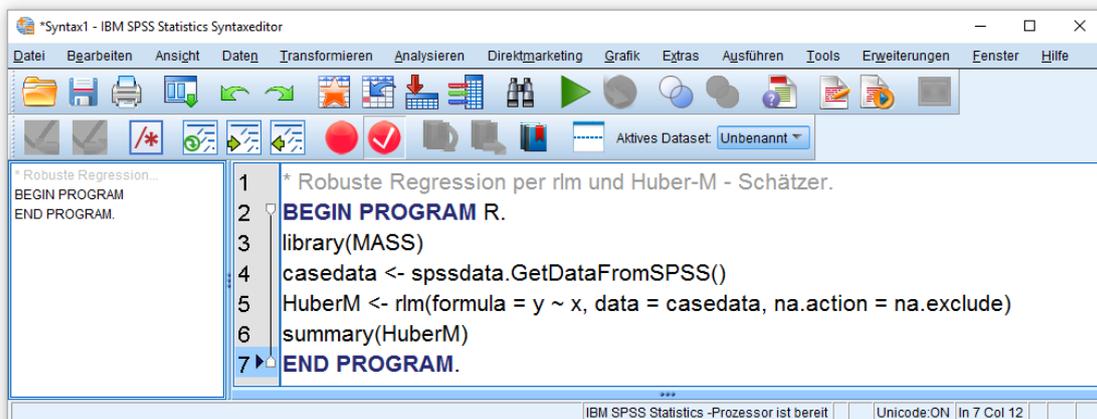
- zur Verwendung der zahlreichen Analysefunktionen in **R**-Paketen
So werden Analyseoptionen zugänglich, die in SPSS fehlen, z.B. Gini-Koeffizient, polychorische Korrelation, Rasch-Itemanalyse, robuste Regression.
- zur Realisation von Algorithmen in der Programmiersprache **R**

Dabei ist es problemlos möglich, Variablen aus einem SPSS-Datenblatt in einen **R** - Data Frame zu kopieren und in **R**-Funktionen über ihre SPSS-Namen anzusprechen.

Die mit **R** erzeugten Ergebnisse können wiederum an SPSS übergeben werden:

- Ausgaben
Normale **R**-Ausgaben landen als Text im SPSS-Ausgabefenster. Durch Verwendung von Funktionen des SPSS-Erweiterungspakets für **R**, das zur Integrationslösung gehört, ist es aber auch möglich, SPSS-Pivot-Tabellen mit **R** zu erstellen.
- Variablen
Aus den mit **R** generierten oder modifizierten Variablen lässt sich ein SPSS-Datenblatt erstellen.

Um **R**-Funktionen auf SPSS-Variablen anzuwenden, erstellt man im Normalfall in einem SPSS-Syntaxfenster einen Block mit **R**-Syntax, eingerahmt durch die SPSS-Kommandos `BEGIN PROGRAM R` und `END PROGRAM`. Im folgenden Beispiel wird für die SPSS-Variablen `y` und `x` eine robuste Regressionsanalyse mit Hilfe der **R**-Funktion `rlm()` gerechnet (vgl. Abschnitt 3.1):

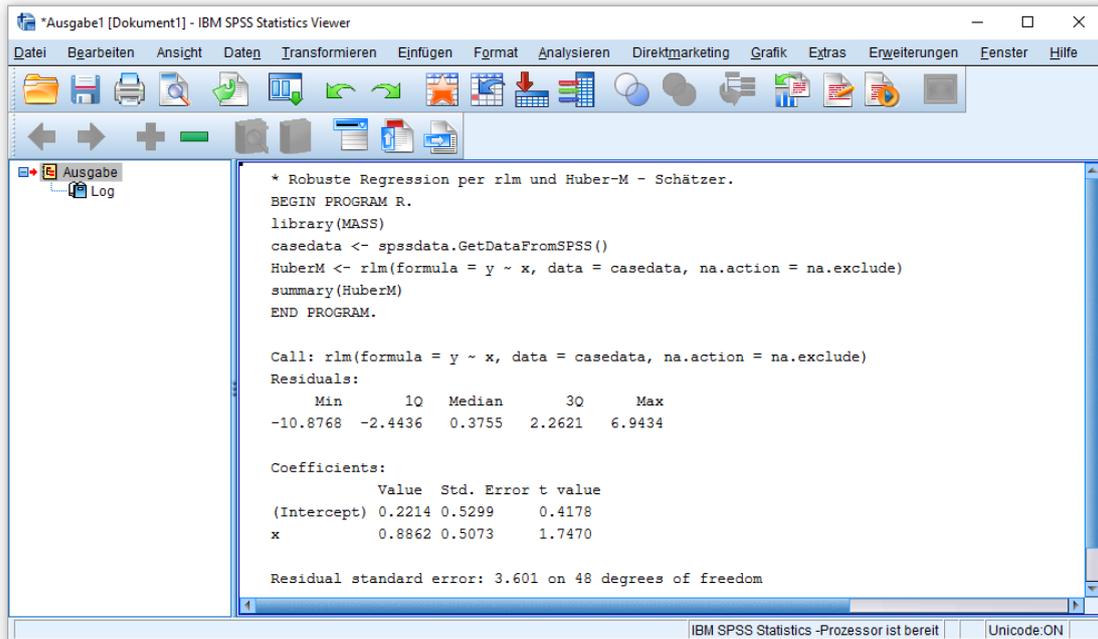


```

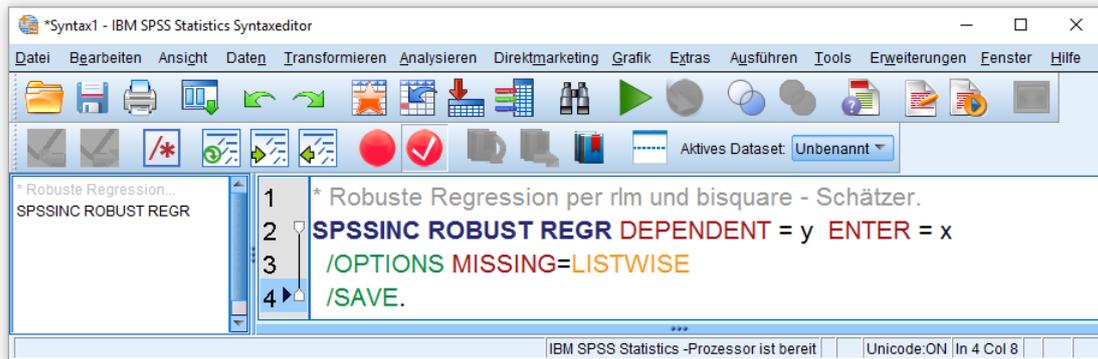
*Robuste Regression per rlm und Huber-M - Schätzer.
1  *Robuste Regression per rlm und Huber-M - Schätzer.
2  BEGIN PROGRAM R.
3  library(MASS)
4  casedata <- spssdata.GetDataFromSPSS()
5  HuberM <- rlm(formula = y ~ x, data = casedata, na.action = na.exclude)
6  summary(HuberM)
7  END PROGRAM.

```

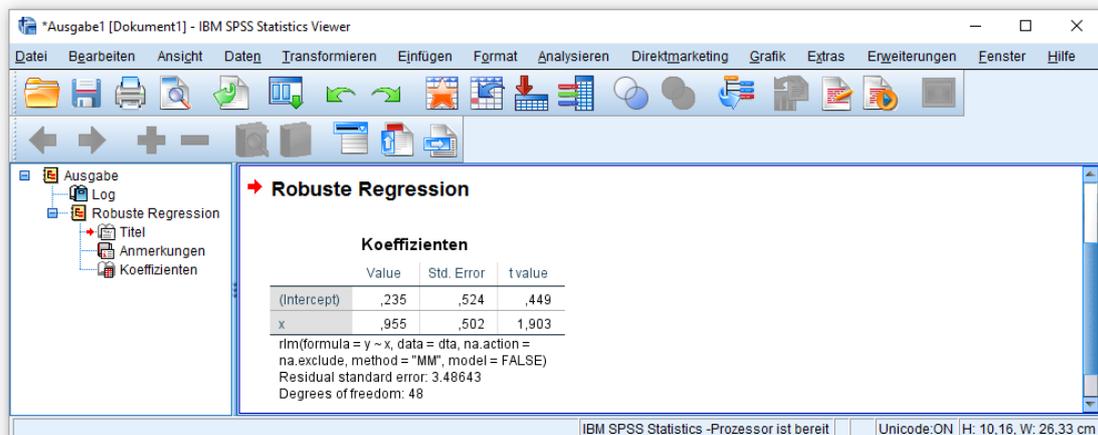
Man erhält die Ausgabe der **R**-Funktion `summary()` als Text im SPSS-Ausgabefenster (Viewer):



Im Rahmen der **R**-Integrationslösung zu SPSS ist es aber auch möglich, ein **SPSS-Erweiterungskommando** zu erstellen, das wie gewöhnliche SPSS-Syntax zu benutzen ist und dabei im Hintergrund mit **R** arbeitet. Für das obige Beispiel ist diese Arbeit von Entwicklern der Firma IBM SPSS erledigt worden, so dass SPSS-Anwender die robuste Regression mit **R** über vertraute Syntax anfordern können:



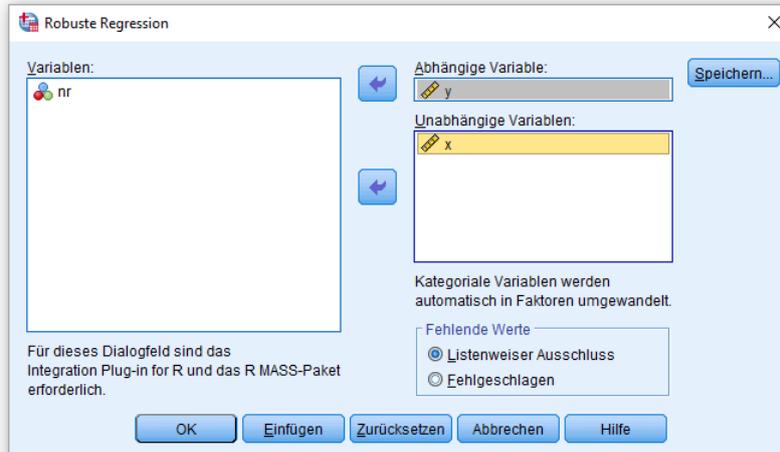
Dass im Beispiel die **R**-Ausgaben in eine **SPSS-Pivot-Tabelle** gewandelt werden, steigert den Benutzerkomfort noch (vgl. Abschnitt 3.1):



Alternativ oder ergänzend zu einem Erweiterungskommando kann ein **benutzerdefinierter Dialog** definiert und in das SPSS-Menüsystem integriert werden, so dass die mit Hilfe von **R** implementierte Funktionalität auch ohne Syntaxfenster nutzbar ist. Im Beispiel erhält man über den Menübefehl

Analysieren > Regression > Robuste Regression

den folgenden Dialog:



Dieser Dialog bildet zusammen mit dem eben vorgestellten Erweiterungskommando und der realisierenden **R**-Syntax ein **Erweiterungsbundle**, das in einer Datei mit der Namensweiterung **SPE** angeboten wird. Es handelt sich um das zusammen mit den **R-Essentials** installierte Bundle **SPSSINC ROBUST REGR** zur Unterstützung der robusten Regression.

Manche Erweiterungsbundles (SPE-Dateien) enthalten nur einen benutzerdefinierten Dialog *ohne* begleitendes Erweiterungskommando. Für Dialoge ohne Erweiterungskommando kann auch ein älteres Dateiformat mit der Namensweiterung **SPD** verwendet werden.

2 SPSS-Erweiterungen auf R- und/oder Python-Basis installieren

Um die R-basierten Funktionserweiterungen für SPSS 24 nutzen zu können, sind die folgenden Installationen auszuführen:

- **SPSS 24**
Dabei sollten (über eine Option des Installationsassistenten) auch die **Python Essentials zu SPSS 24** installiert werden.
- **R 3.2.2 oder R 3.2.5**
IBM SPSS empfiehlt explizit die R-Version 3.2.2, doch spricht erfahrungsgemäß nichts gegen die Verwendung der R-Version 3.2.5 (letzte R-Version aus der 3.2 - Generation).
- **R Essentials zu SPSS 24**
In den R-Essentials sind bereits etliche Erweiterungsbundles enthalten (z.B. das in der Einleitung vorgestellte Bundle SPSSINC ROBUST REGR).
- **Spezielle Erweiterungsbundles oder benutzerdefinierte Dialoge**

2.1 Python- und R-Essentials

2.1.1 Python-Essentials

Python ist eine attraktive Skriptsprache, die für Automatisierungszwecke in SPSS verwendet werden kann. Das mit den R-Essentials gelieferte Erweiterungskommando für heterogene Korrelationen (siehe Abschnitt 3.3) benötigt Python, so dass Sie auch die **Python-Essentials** installieren sollten. Dies kann bequem über eine Option im Installations-Assistenten von SPSS 24 geschehen. Es ist also kein separater Download erforderlich.

Wenn Sie bei der Installation von SPSS 24 die Python-Option wählen, werden die folgenden Bestandteile der Python-Essentials eingerichtet:

- Python 2.7.6 und Python 3.4.3
- Python-Integration-Plugin für SPSS 24
- Python-basierte Erweiterungsbundles
Es werden etliche mit Hilfe von Python implementierte SPSS-Erweiterungskommandos samt zugehöriger benutzerdefinierter Dialoge mitgeliefert, z.B. zum automatischen Erstellen der Kodiervariablen zu einer nominalskalierten Variablen (Menübefehl: **Transformieren > Dummy-Variablen erstellen**).

Sollten Sie bei der Installation von SPSS 24 auf die Python-Option verzichtet haben, können Sie deren Installation so nachholen:

- Installationsmenü zu SPSS 24 erneut starten per Doppelklick auf **setup.exe** im auspackten Installationsordner oder im Hauptverzeichnis der DVD
- Mausklick auf das Item **IBM SPSS Statistics 24 installieren**
- Im Dialog **Programmverwaltung** wählen: **Programm ändern**
- Im Dialog **IBM SPSS Statistics - Essentials for Python** der Installation zustimmen

2.1.2 R-Essentials

Die R-Essentials zu SPSS 24 werden über die folgende Webseite angeboten:

https://github.com/IBMPredictiveAnalytics/R_Essentials_Statistics/releases

Als Ergänzung zu SPSS 24 auf einem Windows-System mit 64-Bit - Architektur erhält man z.B. die folgende Installationsdatei:

SPSS_Statistics_R_Essentials_24_win64.exe

In dem auf derselben Webseite unter dem Titel **Installation Instructions for Windows** angebotenen Archiv

SPSS_Statistics_R_Essentials_Installation_Documents_24_win.zip

findet sich die folgende PDF-Datei mit Installationshinweisen:

Essentials for R Installation Instructions.pdf

In den **R**-Essentials zu SPSS 24 sind enthalten:

- Das **R**-Integration-Plugin für SPSS 24
- Das **R**-Paket **spss240**
- **R**-basierte Erweiterungsbundles für SPSS 24
Es werden etliche mit Hilfe von **R** implementierte Erweiterungsbundles installiert, die aus einem Erweiterungskommando und einen benutzerdefinierten Dialog bestehen (siehe Abschnitt 3).

Im **R** Essentials - Installationspaket zu SPSS 24 ist die vorausgesetzte **R**-Version *nicht* enthalten. Diese muss aus anderer Quelle beschafft und *vor* den **R**-Essentials installiert werden. Das eben genannte PDF-Dokument empfiehlt, die **R**-Version 3.2.2 zu verwenden und nennt u.a. den folgenden Download-Link:

<ftp://ftp.stat.math.ethz.ch/Software/CRAN/bin/windows/base/old/>

Alternativ kann auch die **R**-Version 3.2.5 verwendet werden (letzte **R**-Version aus der 3.2 - Generation).

2.1.3 Installationsarbeiten im Überblick

Gehen Sie also folgendermaßen vor, um SPSS 24, **R** 3.2 und die zugehörigen **R**-Essentials zu installieren:

- **SPSS 24 installieren (inkl. Python-Essentials)**
Falls noch nicht geschehen, müssen Sie zuerst SPSS 24 installieren, wobei Sie nicht auf die Option der **Python-Essentials** verzichten sollten (vgl. Abschnitt 2.1.1). Anschließend sollten Sie das aktuelle FixPack zu SPSS 24 installieren.
- **R 3.2.2 oder R 3.2.5 installieren**
Wenn Sie **R** auf einem Rechner mit 64-Bit - Windows installieren, müssen Sie unbedingt die 32-Bit-Version von **R** in die Installation einbeziehen (= Voreinstellung).
- **R-Essentials für SPSS 24 installieren**
Zusammen mit einem FixPack zu SPSS werden eventuell auch die **R**-Essentials aktualisiert. Es ist also darauf zu achten, das korrekte **R**-Essentials - Installationsprogramm zu verwenden. Zu SPSS 24 passt auf einem Windows-Rechner mit 64-Bit - Architektur das Installationsprogramm **SPSS_Statistics_R_Essentials_24_win64.exe**.

2.2 Erweiterungsbundles

2.2.1 Inhalt

Erweiterungsbundles ergänzen die SPSS-Funktionalität durch zusätzliche Verfahren zur Datenanalyse und/oder -verwaltung, die durch SPSS-interne Programmieroptionen (z.B. Makrotechnik, Programmiersprache MATRIX) oder externe Programmieroptionen (z.B. **R** oder Python) realisiert werden. Zur Nutzung der erweiterten Funktionalität bietet ein Bundle ein **Erweiterungskommando** und/oder einen **benutzerdefinierten Dialog**. Das in der Einleitung vorgestellte Bundle **SPSSINC ROBUST REGR** unter-

stützt *beide* Optionen. Speziell bei den *nicht* von IBM SPSS entwickelten Bundles ist oft die Beschränkung auf einen benutzerdefinierten Dialog anzutreffen. So bietet z.B. Hans Grüner vom Rechenzentrum der FU Berlin diverse benutzerdefinierte Dialoge mit **R**-Implementierung an (vgl. Abschnitt **Fehler! Verweisquelle konnte nicht gefunden werden.**), die er neuerdings in Erweiterungsbundles verpackt, um die bequeme Distribution, Installation und Aktualisierung innerhalb von SPSS 24 zu ermöglichen (siehe Abschnitt 2.2.3).

2.2.2 Erstellung

Wer ein Erweiterungsbundle mit **R** erstellen möchte, muss den Implementierungscode mit **R** verfassen und die Bedienung per Syntax und/oder Dialog unterstützen. Über die Erstellung eines Erweiterungskommandos informiert das Buch *Programming and Data Management for IBM SPSS Statistics 24* (IMB Corp. 2016a). Wer ein benutzerdefiniertes Dialogfeld erstellen möchte, findet eine Anleitung im *IBM SPSS Statistics 24 Core-System Benutzerhandbuch* (IBM Corp. 2016b, Kapitel 20). Die beiden genannten Bücher sind kostenlos im PDF-Format verfügbar.

Um die erstellten Dateien in ein Erweiterungsbundle zu verpacken, startet man mit dem Menübefehl

Erweiterungen > Extras > Erweiterungsbundle erstellen

Es resultiert eine Erweiterungsbundle-Datei mit der Namensweiterung **SPE**.

2.2.3 Installation

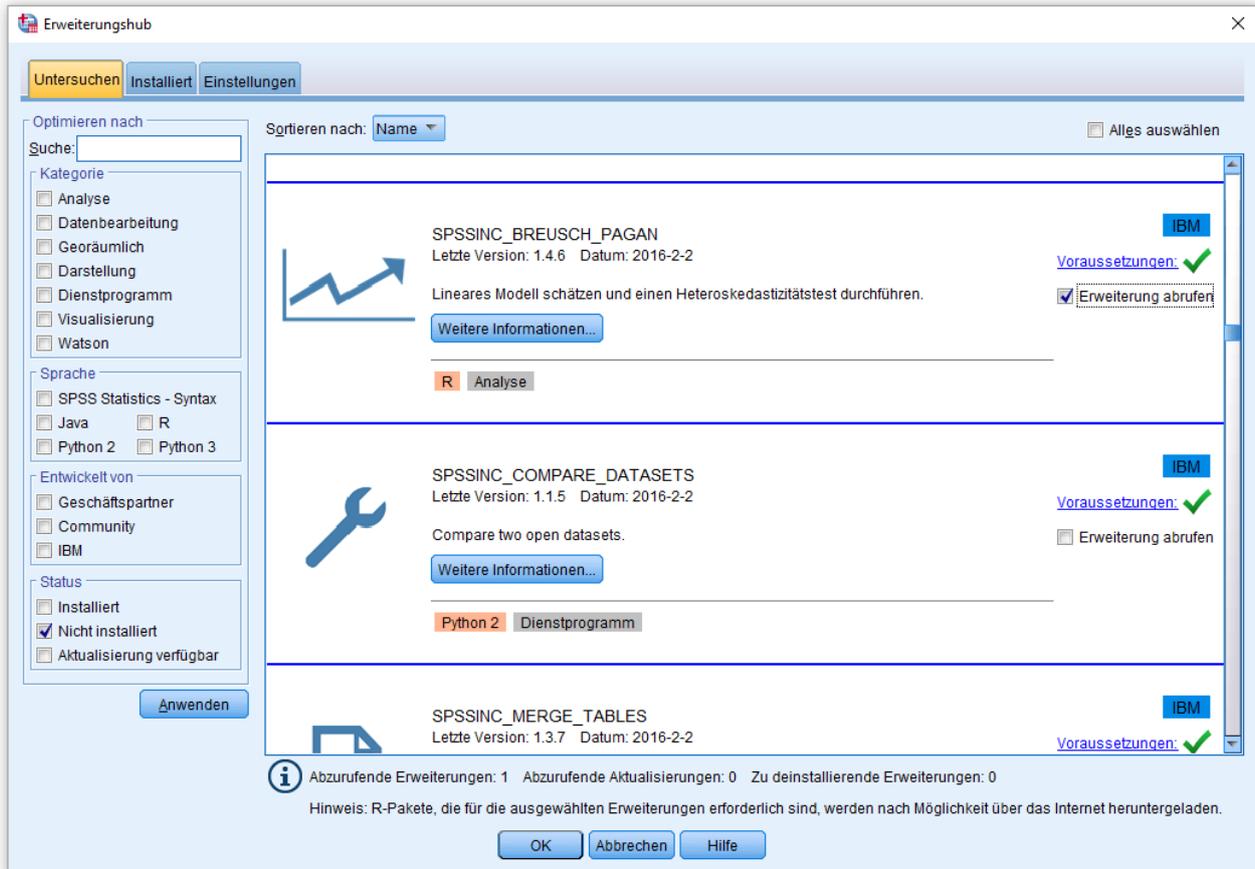
Mit den **R**-Essentials für SPSS 24 werden ca. 50 mit Hilfe von **R** und/oder Python implementierte Erweiterungsbundles installiert, die auch in das Menüsystem integriert sind, z.B.:

- **Analysieren > Korrelation > Heterogene Korrelationen**
- **Analysieren > Regression > Robuste Regression**
- **Analysieren > Regression > Tobit-Regression**
- **Analysieren > Regression > Quantilregression**
- **Analysieren > Skala > Rasch erweitert**
- **Analysieren > Vorhersage > GARCH-Modelle**

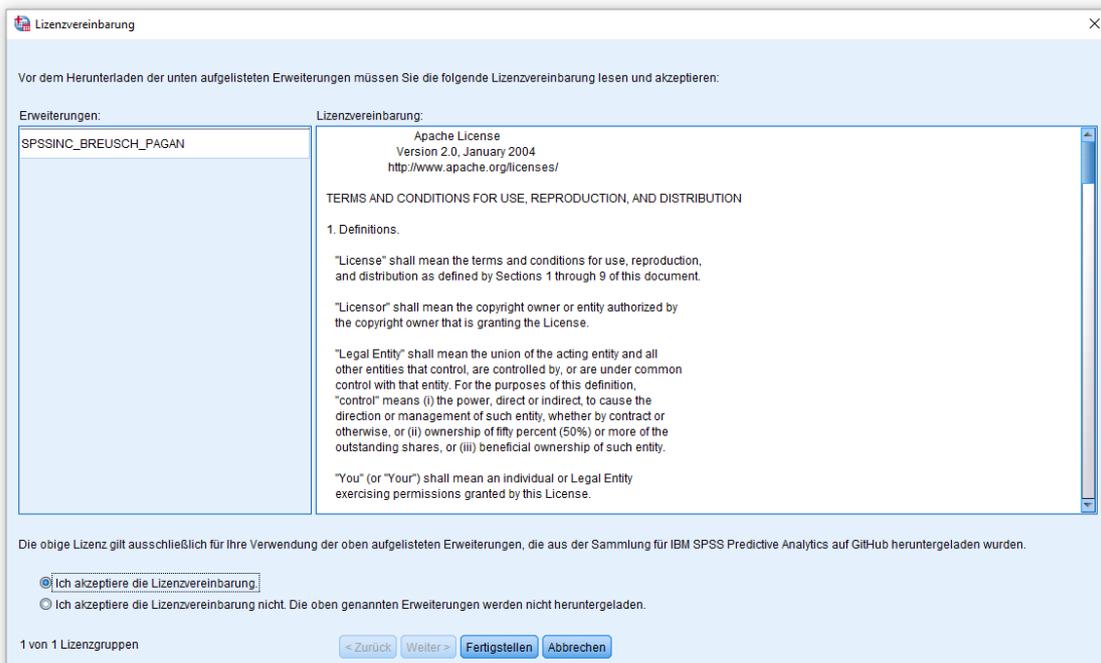
Seit der Version 22 kann SPSS Erweiterungsbundles selbständig aus dem Internet beziehen, um sie zu installieren oder zu aktualisieren, so dass der vorherige separate Download entfällt. Ob die Erweiterung anschließend für alle Benutzer oder nur für den Installateur verfügbar ist, hängt davon ab, ob SPSS mit Administratorrechten ausgeführt wird (siehe unten). In der Version 24 erscheint nach dem Menübefehl

Erweiterungen > Erweiterungshub

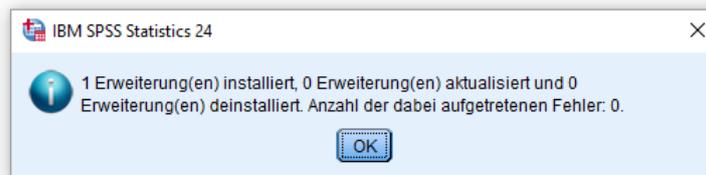
ein Dialog, der die installierten, die verfügbaren und die aktualisierbaren Erweiterungen samt Versionsstand anzeigt, z.B.:



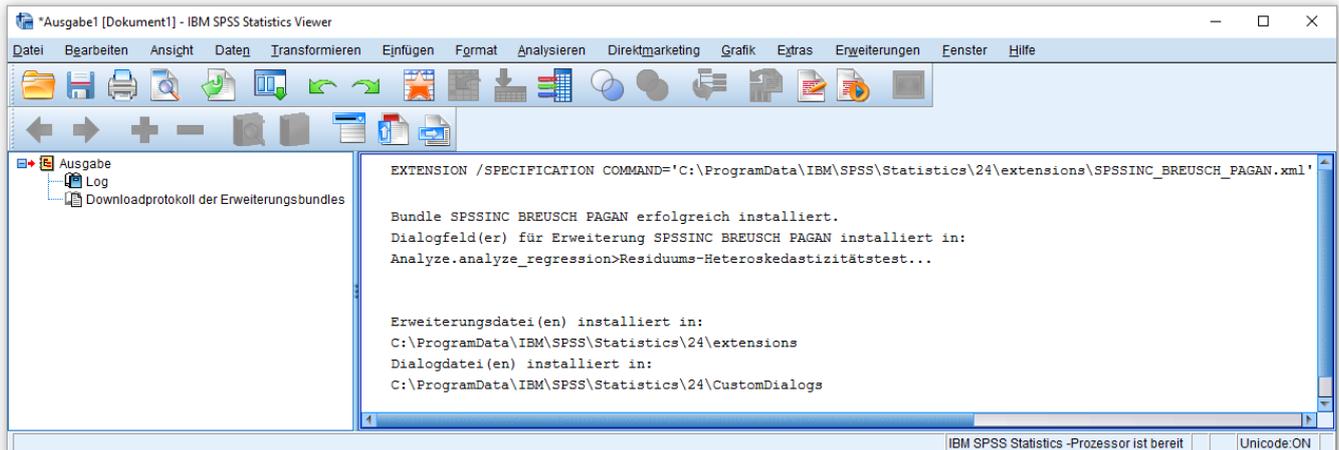
Um eine Erweiterung zu installieren oder zu aktualisieren, markiert man das Kontrollkästchen **Erweiterung abrufen** bzw. **Aktualisierung abrufen** in der zugehörigen Zeile und klickt auf **OK**. Als Beispiel installieren wir die Erweiterung **SPSSINC_BREUSCH_PAGAN**, die einen Heteroskedastizitätstest zum Linearen Modell durchführt und leider (abweichend von der früheren Praxis) nicht mehr automatisch mit den **R**-Essentials installiert wird. Im folgenden Dialog akzeptiert man die Lizenzvereinbarung und fordert die **Fertigstellung** an:



Nach einer erfolgreichen Installation



wird im Ausgabefenster protokolliert, wo das zum Bundle gehörige benutzerdefinierte Dialogfeld installiert worden ist, z.B.:



Wer ein Erweiterungsbundle per SPE-Datei bezogen hat und es nun installieren möchte, wählt folgenden Menübefehl:

Erweiterungen > Lokales Erweiterungsbundle installieren

Spätestens nach einem Neustart von SPSS stehen das Erweiterungskommando und das Dialogfeld allen Benutzern zur Verfügung. Die voreinstellten Installationsorte unter Windows sind:

- Erweiterungskommandos landen zusammen mit den Dateien zur Bundle-Konfiguration in einem Unterordner von
C:\Program Files\IBM\SPSS\Statistics\24\extensions
- Benutzerdefinierte Dialoge landen in einem Unterordner von
C:\ProgramData\IBM\SPSS\Statistics\24\CustomDialogs

Auf einem Pool-PC an der Universität Trier ist die Installation einer Erweiterung von räumlich begrenzter Wirkung. Die Installation landet auf einem einzelnen Pool-PC, kann sich also nicht auf andere Pool-PCs auswirken. Zur Lösung des Problems sollten Sie die benutzereigenen Windows-Umgebungsvariablen **SPSS_CDIALOGS_PATH** und **SPSS_EXTENSIONS_PATH** definieren über

Systemsteuerung > Benutzerkonten > Benutzerkonten > Eigene Umgebungsvariablen ändern

Dabei sind folgende Regeln zu beachten:

- SPSS berücksichtigt nur Umgebungsvariablen, die schon beim Programmstart vorhanden waren.
- Als Name für die Umgebungsvariable muss **SPSS_EXTENSIONS_PATH** bzw. **SPSS_CDIALOGS_PATH** verwendet werden.
- Als Inhalt sollten Sie den Namen eines bereits vorhandenen (!) Ordners auf Ihrem persönlichen Laufwerk **U:** angeben, weil dieses Laufwerk auf jedem Pool-PC verfügbar ist, z.B.
 - bzw. **U:\Eigene Dateien\SPSS\Extensions**
 - bzw. **U:\Eigene Dateien\SPSS\CustomDialogs**

Wenn Sie anschließend ein Erweiterungsbundle (mit benutzerdefiniertem Dialog) installieren, landen die Dateien in den vereinbarten benutzereigenen Ordnern, und die Erweiterung steht dauerhaft auf jedem Pool-PC zur Verfügung.

Das Dialogfeld zu der als Beispiel installierten Breusch-Pagan - Erweiterung ist über den folgenden Menübefehl erreichbar:

Analysieren > Regression > Residuums-Heteroskedastizitätstest

Benötigt ein Erweiterungsbundle zusätzliche **R**-Pakete, versucht der Bundle-Installer, diese aus dem Internet herunter zu laden und zu installieren, was auch mit normalen Benutzerrechten gelingt, weil **R** seine Pakete system- oder benutzerbezogen installieren kann (vgl. Abschnitt 5.2.2).

Eine kurze Funktionsbeschreibung der von IBM SPSS erstellten Erweiterungsbundles bietet Cohen (2016).

2.3 SPD-Dateien

Neben dem per SPE-Datei verteilten Erweiterungsbundle, das einen benutzerdefinierten Dialog und/oder ein Erweiterungskommando enthalten kann, existiert noch ein kleineres (und älteres) Distributionsformat, das lediglich einen benutzerdefinierten Dialog enthält und in einer Datei mit der Namensendung SPD steckt. Die per Dialog bequem verfügbare Funktionalität wird entweder durch eine SPSS-interne Programmieroption (z.B. Makrotechnik, Programmiersprache MATRIX) oder durch eine externe Programmieroption (z.B. **R** oder Python) realisiert.

Wer einen benutzerdefinierten Dialog per SPD-Datei bezogen hat und nun installieren möchte, wählt den Menübefehl:

**Erweiterungen > Extras >
Benutzerdefiniertes Dialogfeld installieren (Kompatibilitätsmodus)**

Anschließend steht das benutzerdefinierte Dialogfeld allen Benutzern zur Verfügung. Per Voreinstellung erfolgt die Installation unter Windows in einen Unterordner von

C:\ProgramData\IBM\SPSS\Statistics\24\CustomDialogs

Auf einem Pool-PC an der Universität Trier ist die Installation eines benutzerdefinierten Dialogfelds von räumlich begrenzter Wirkung. Die Installation landet auf einem einzelnen Pool-PC, kann sich also nicht auf andere Pool-PCs auswirken. Zur Lösung des Problems sollten Sie die benutzereigene Windows-Umgebungsvariable **SPSS_CDIALOGS_PATH** definieren über

**Systemsteuerung > Benutzerkonten > Benutzerkonten >
Eigene Umgebungsvariablen ändern**

Dabei sind folgende Regeln zu beachten:

- SPSS berücksichtigt nur Umgebungsvariablen, die schon beim Programmstart vorhanden waren.
- Als Name für die Umgebungsvariable muss **SPSS_CDIALOGS_PATH** verwendet werden.
- Als Inhalt sollten Sie den Namen eines bereits vorhandenen (!) Ordners auf Ihrem persönlichen Laufwerk **U:** angeben, weil dieses Laufwerk auf jedem Pool-PC verfügbar ist, z.B. **U:\Eigene Dateien\SPSS\CustomDialogs**

Wenn Sie anschließend ein benutzerdefiniertes Dialogfeld installieren, landet es im vereinbarten Ordner und steht dauerhaft auf jedem Pool-PC zur Verfügung.

Benötigt ein benutzerdefiniertes Dialogfeld zusätzliche **R**-Pakete, wird versucht, diese aus dem Internet herunter zu laden und zu installieren, was auch mit normalen Benutzerrechten gelingt, weil **R** seine Pakete system- oder benutzerbezogen installieren kann (siehe Abschnitt 2.2).

3 Exemplarische Beschreibung einiger Erweiterungsbundles

In diesem Abschnitt demonstrieren wir einige Erweiterungsbundles zu SPSS 24, die allesamt auch über ein Dialogfeld verfügen. Damit stehen wichtige Funktionsergänzungen bei größtmöglicher Bequemlichkeit zur Verfügung, und wir können uns in diesem Abschnitt auf die Forschungsmethodik konzentrieren.

3.1 Robuste Regression

Bei der OLS-Regressionsanalyse (*Ordinary Least Squares*) können problematische Einzelwerte die Schätzung und Testung der Modellparameter verfälschen:

- **Große Residuen**
Hier geht es um Fälle mit extremen Residuen, die nicht auf Erhebungs- oder Erfassungsfehler zurückgehen und auch nicht überzeugend als außerhalb der Betrachtung liegend (z.B. zu anderen Populationen gehörig) entfernt werden können (Ausreißer).
- **Starke Hebel**
Fälle mit extremen Werten bei unabhängigen Variablen verfügen über eine große Hebelwirkung und damit über einen oft unerwünscht starken Effekt auf die Schätz- und Testergebnisse.

Der Gesamteinfluss eines Falls auf die Ergebnisse ist im Wesentlichen ein Produkt aus der Hebelwirkung und der absoluten Größe des Residuums.

Im vorzustellenden Erweiterungsbundle sowie in der zugrunde liegenden **R**-Funktion **rlm()** aus dem **R**-Paket **MASS** kommen so genannte **M-Schätzer** zum Einsatz, die ...

- zwar den Einfluss von Fällen mit großen Residuen begrenzen,
- jedoch *nicht* robust sind gegenüber Fällen mit unauffälligen Residuen, aber große Hebelwirkung.

Die Verfahren mit M-Schätzern arbeiten als WLS-Regression (*Weighted Least Squares*), wobei das Einflussgewicht eines Falles sinkt, je größer sein Residuum betragsmäßig ausfällt. Zur Schätzung der Fallgewichte und der Regressionskoeffizienten wird ein iteratives Verfahren eingesetzt (siehe z.B. IDRE UCLA 2014a):

- Aufgrund der Residuen basierend auf den aktuellen Parameterschätzungen werden neue Fallgewichte berechnet (je größer der Residualbetrag, desto kleiner das Gewicht).
- Mit den neuen Fallgewichten erhält man aktualisierte Parameterschätzungen und neue Residuen.

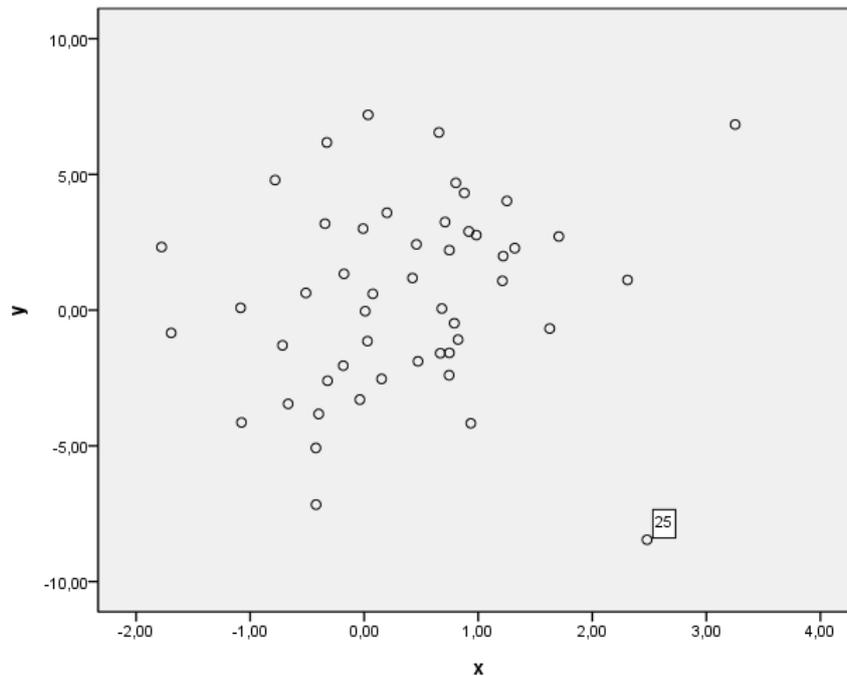
Wenn sich die Parameterschätzer nicht mehr ändern, stoppt das Verfahren.

Eine ausführliche Behandlung der robusten Regression bietet z.B. Ryan (1997).

Wir betrachten als Anwendungsbeispiel synthetische Daten mit einer bivariaten Regression und einem Ausreißer, der die Schätzung und Testung empfindlich stört. Das wahre Modell lautet:

$$Y_i = 1 \cdot X_i + \varepsilon_i, \quad \text{mit } \text{Var}(\varepsilon_i) = 9, \text{Cov}(\varepsilon_i, \varepsilon_j) = 0$$

In der Stichprobe mit $n = 50$ befindet sich ein Fall (Nr. 25) mit dem „Störimpuls“ -13 im Y -Messwert, was im Streudiagramm gut zu erkennen ist:



In der OLS-Regression mit SPSS erhält man folgende Schätz- und Testergebnisse:

Modellzusammenfassung^b

Modell	R	R-Quadrat	Korrigiertes R-Quadrat	Standardfehler des Schätzers
1	,172 ^a	,030	,010	3,49338

a. Einflußvariablen : (Konstante), x

b. Abhängige Variable: y

Koeffizienten^a

Modell		Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten		95,0% Konfidenzintervalle für B		
		Regressionskoeffizient B	Standardfehler	Beta	T	Sig.	Untergrenze	Obergrenze
1	(Konstante)	,245	,528		,465	,644	-,816	1,306
	x	,613	,505	,172	1,213	,231	-,403	1,628

a. Abhängige Variable: y

Die Schätzung zum Steigungskoeffizienten ist deutlich gemindert (0,61 statt 1,00), und der Signifikanztest erlaubt es nicht, die falsche Nullhypothese zu verwerfen ($p = 0,231 > 0,05$). Als Schätzung der Effektstärke erhalten wir ein sehr kleines korrigiertes R^2 von 0,01.

In der Tabelle mit den Cook-Distanzen zur Beurteilung der Einflussstärke von Einzelfällen¹ erreicht der Fall 25 den größten Wert (zur Definition der Cook-Distanzen siehe z.B. Baltes-Götz 2016, Abschnitt 3.1.3):

¹ In SPSS per Syntax anzufordern mit dem folgenden Subkommando der Regressionsprozedur:
/RESIDUALS outliers(cook)

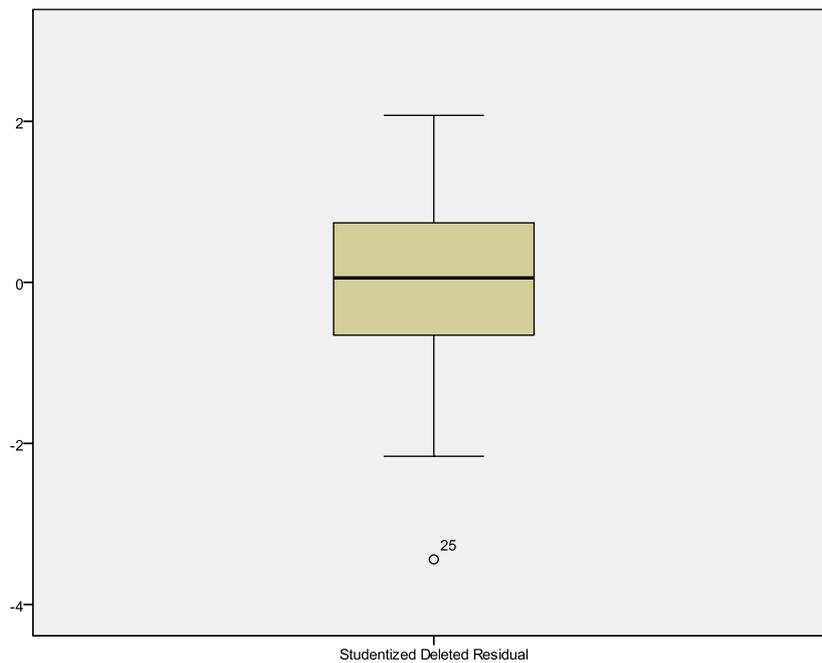
Ausreißerstatistik^a

	Fallnummer	Statistik	Sig. F
Cook-Distanz	1	,617	,544
	2	,259	,773
	3	,074	,929
	4	,061	,941
	5	,054	,947
	6	,049	,952
	7	,046	,955
	8	,041	,960
	9	,037	,964
	10	,032	,968

a. Abhängige Variable: y

Als kritische Cook-Distanzen werden in der Literatur die Werte 1 (z.B. Weisberg 1985) und $4/n$ (z.B. Gordon 2010, S. 367) genannt, wobei wohl die zweite, strengere Grenze herangezogen werden sollte. In unserem Beispiel ergibt sich der Grenzwert 0,08, den der kritische Fall 25 deutlich überschreitet.

Das ausgelassen-studentisierte Residuum von -3,441 für Fall 25 deutet auf einen ernst zu nehmenden Ausreißer hin (zur Definition der ausgelassen-studentisierten Residuen siehe z.B. Baltes-Götz 2016, Abschnitt 1.7.2):



Entfernt man den Fall 25 aus der Analyse, liefert die OLS-Regression ein deutlich verschiedenes Ergebnisbild:

Modellzusammenfassung^b

Modell	R	R-Quadrat	Korrigiertes R-Quadrat	Standardfehler des Schätzers
1	,323 ^a	,104	,085	3,15529

a. Einflußvariablen : (Konstante), x

b. Abhängige Variable: y

Koeffizienten^a

Modell		Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten	T	Sig.	95,0% Konfidenzintervalle für B	
		Regressionskoeffizient B	Standardfehler	Beta			Untergrenze	Obergrenze
1	(Konstante)	,289	,477		,606	,548	-,670	1,248
	x	1,122	,480	,323	2,339	,024	,157	2,087

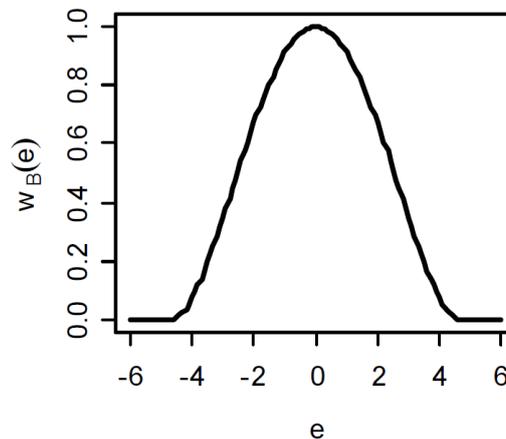
a. Abhängige Variable: y

Wir erhalten eine recht präzise Schätzung für den Regressionskoeffizienten (wahrer Wert: 1) mit einem signifikanten Testergebnis ($p = 0,024$) und einem „stattlichen“ korrigiertes R^2 von 0,085.

Die SPSS-Erweiterung zur robusten Regression ist nach der Installation der R-Essentials verfügbar über den Menübefehl¹

Analysieren > Regression > Robuste Regression

und verwendet die Funktion `rlm()` aus dem **R**-Paket **MASS**. Dabei wählt SPSS von den verfügbaren robusten Schätzmethoden den so genannten *biquadratischen Schätzer* (engl.: *bisquare estimator*), der folgende Gewichtungsfunktion verwendet (Abbildung übernommen aus Fox 2002, S. 3):



Der Einfluss eines Falles wird mit wachsendem Betrag seines Residuums reduziert, wobei auch kleine Abweichungsbeträge bereits zu einer Minderung führen.

Für die Beispieldaten (*inkl.* Fall 25) resultiert die folgende Tabelle mit `rlm()` - Ergebnissen:

Koeffizienten

	Value	Std. Error	t value
(Intercept)	,235	,524	,449
x	,955	,502	1,903

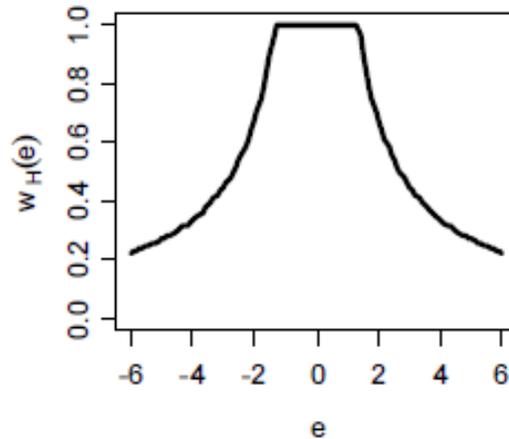
```
rlm(formula = y ~ x, data = dta, na.action =
na.exclude, method = "MM", model = FALSE)
Residual standard error: 3.48643
Degrees of freedom: 48
```

Der Steigungskoeffizient wird sehr präzise geschätzt (wahrer Wert: 1). Leider werden zu den Regressionskoeffizienten mangels Vertrauen in die Verteilung der Prüfstatistik (berechnet als Quotient aus dem Schätzer und seinem Standardfehler) keine Überschreitungswahrscheinlichkeiten (p -Werte) geliefert.

¹ Bis zur Version 1.2.1 des Erweiterungsbundles zur robusten Regression hieß das Menüitem **Solide Regression**.

Ebenso fehlen Vertrauensintervalle. In Abschnitt 4.2 wird demonstriert, wie man mit Hilfe von **R**-Syntax approximative p -Werte ermitteln kann.

Ist an Stelle des biquadratischen Schätzers, den SPSS im Erweiterungskommando verwendet, Hubers M-Schätzer mit der folgenden Gewichtungsfunktion (Abbildung übernommen aus Fox 2002, S. 3)



gewünscht, kann man in einem **R**-Syntax - Block die Variablen der SPSS-Arbeitsdatei an **R** übergeben und die **R**-Funktion `rlm()` direkt verwenden (siehe Abschnitt 4.2). Bei der Huber-Gewichtung muss der Residuumsbetrag eine Grenze überschreiten, bevor die Minderung des Einflussgewichts einsetzt. Im Unterschied zum biquadratischen Schätzer bleibt der Einfluss auch bei Fällen mit sehr großem Residuumsbetrag größer 0.

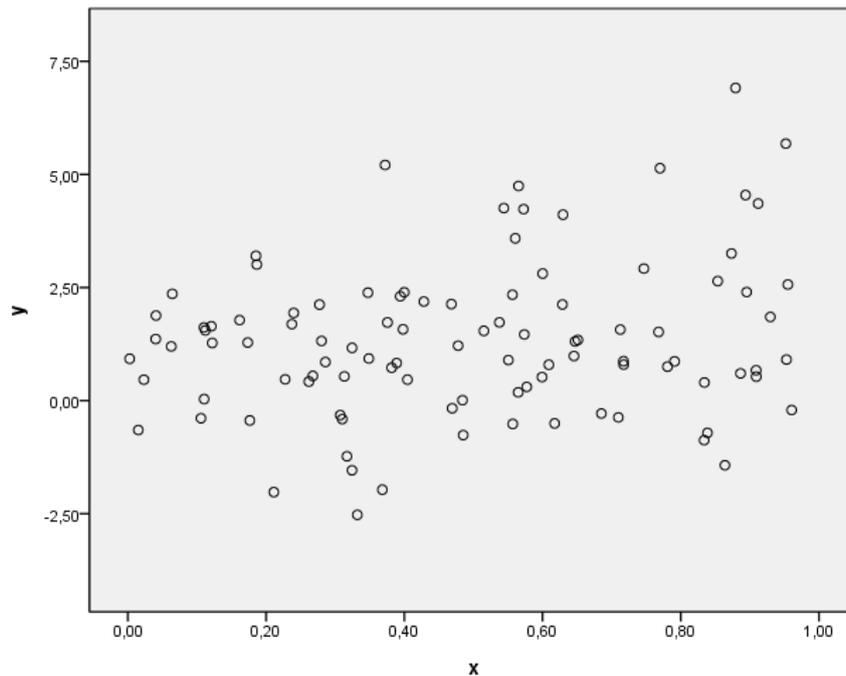
3.2 Breusch-Pagan - Heteroskedastizitäts-Test

Der Breusch-Pagan - Test prüft für lineare Modelle die Nullhypothese homogener Fehlervarianzen.

Wir betrachten als Anwendungsbeispiel synthetische Daten mit einer bivariaten Regression und Fehlervarianzen mit ausgeprägter Abhängigkeit vom Wert des Regressors:

$$Y_i = 2 \cdot X_i + \varepsilon_i, \quad \text{mit } \text{Var}(\varepsilon_i) = (1 + X_i)^2$$

In der Stichprobe mit $n = 100$ ist die mit dem Regressor wachsende Fehlervarianz gut zu erkennen:



Nach der Installation des Erweiterungsbundles **SPSSINC_BREUSCH_PAGAN** (siehe Abschnitt 2.2.3) ist in SPSS über den Menübefehl

Analysieren > Regression > Residuums-Heteroskedastizitätstest

ein Dialog zur robusten Regression verfügbar, wobei im Hintergrund die Funktion **ncvTest** aus dem **R**-Paket **car** zum Einsatz kommt.

Im Beispiel wird die Heteroskedastizität erkannt:

Non-constant Variance Score Test

	ChiSquare	D.f	Sig.
Test Result	6,352	1,000	,012

Variance model: fitted values
Computed by R ncvTest function

Nach einem signifikanten Heterogenitätstest muss man übrigens den Versuch einer linearen Modellierung nicht aufgeben. Mittlerweile sind Verfahren zur robusten Inferenzstatistik trotz Heteroskedastizität entwickelt worden, die in **R** wie in SPSS zur Verfügung stehen (siehe z.B. Baltes-Götz 2016, Abschnitt 1.7.3.4).

3.3 Polychorische und polyseriale Korrelationen

In der statistisch-empirischen Forschung sind häufig *ordinale* Variablen im Sinne einer vergrößernden Messung von latenten Variablen mit metrischer Skalenqualität und approximativer Normalverteilung anzutreffen. Oft ist man an den Korrelationen zwischen den latenten Merkmalen interessiert, hat aber nur die ordinalen Indikatoren zur Verfügung. Mit der Pearson-Formel zur Korrelationsberechnung erhält man aus den ordinalen Maßen mehr oder weniger stark verzerrte Schätzungen (siehe unten).

Stimmt für zwei ordinale Indikatoren die Annahme, dass sie durch vergrößernde Messung aus einer bivariaten Normalverteilung der eigentlich interessierenden latenten Variablen hervorgegangen sind, dann kommt das **polychorische** Schätzverfahren der gesuchten Korrelation näher als die Pearson-Formel.

Ist von den beiden zu korrelierenden Merkmalen nur eines vergrößernd gemessen worden, sodass ein intervallskaliertes Merkmal auf ein ordinal trifft, dann erbringt die **polyseriale** Schätzformel eine analoge Reparaturleistung.

Als Anwendungsbeispiel betrachten wir latente Variablen X und Y mit einer wahren Korrelation von 0,5:¹

$$Y_i = 0,5 \cdot X_i + \varepsilon_i \quad \text{mit } \varepsilon_i \sim N(0;0,75), \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \text{ und } X_i \sim N(0;1)$$

Daraus entstehen durch vergrößerndes Messen die ordinalen Indikatoren X3 und Y3 mit jeweils 3 Ausprägungen, z.B. mit Hilfe der folgenden SPSS-Syntax:

```
RECODE X (LO THRU -1 = 1) (-1 THRU 1 = 2) (ELSE = 3) INTO X3.
RECODE Y (LO THRU 0 = 1) (0 THRU 2 = 2) (ELSE = 3) INTO Y3.
VARIABLE LEVEL X3, Y3 (ORDINAL).
```

Verwendet man die SPSS-Prozedur zur Berechnung von Pearson-Korrelationen, dann resultieren deutlich geminderte Schätzungen für die interessierende Korrelation von X und Y durch die Korrelationen von X3 und Y3 (beide Variablen ordinal) bzw. X und Y3 (nur Y3 ordinal):

		X	Y	X3	Y3
X	Korrelation nach Pearson	1	,526	,848	,389
	Signifikanz (2-seitig)		,000	,000	,000
	N	500	500	500	500
Y	Korrelation nach Pearson	,526	1	,446	,827
	Signifikanz (2-seitig)	,000		,000	,000
	N	500	500	500	500
X3	Korrelation nach Pearson	,848	,446	1	,343
	Signifikanz (2-seitig)	,000	,000		,000
	N	500	500	500	500
Y3	Korrelation nach Pearson	,389	,827	,343	1
	Signifikanz (2-seitig)	,000	,000	,000	
	N	500	500	500	500

Nach der Installation der **R**-Essentials ist in SPSS ein Dialog zur Berechnung von polychorischen und polyserialen Korrelationen mit Hilfe von **R** (und Python) verfügbar über den Menübefehl²

Analysieren > Korrelation > Heterogene Korrelationen

Im Beispiel erhalten wir im polychorischen und im polyserialen Fall deutlich verbesserte Schätzungen in der Nähe des wahren Wertes:

$$^1 \text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{\text{Cov}(X, 0,5X + \varepsilon)}{\sqrt{\text{Var}(0,5X + \varepsilon)}} = \frac{0,5}{\sqrt{0,25 + \text{Var}(\varepsilon)}} = \frac{0,5}{\sqrt{0,25 + 0,75}} = 0,5$$

² **Achtung:** Das vom Erweiterungsbundle für heterogene Korrelationen benötigte **R**-Paket **polycor** wird leider (im Dezember 2016) von den CRAN-Spiegelservern (*Comprehensive R Archive Network*) zu der von SPSS 24 benötigten **R**-Version 3.2 nicht ausgeliefert. Folglich scheitert die Berechnung von polychorischen und polyserialen Korrelationen. Allerdings kann man (am 05.12.2016) das Paket **polycor** für die **R**-Version 3.0 von den CRAN-Spiegelservern beziehen und nach dem Kopieren des Paketordners `...|library|polycor` in das Verzeichnis von **R** 3.2.5 auch zusammen mit dieser **R**-Version und damit zusammen mit dem Erweiterungsbundle zu SPSS 24 verwenden. Ein auf diese Weise beschaffter Paketordner **polycor** ist an der im Vorwort beschriebenen Stelle neben den Beispieldateien zu den heterogenen Korrelationen zu finden.

Pearson-, polyseriale und polychorische Korrelationen

Variablen	Statistiken	Variablen	
		X3	Y3
X3	Korrelation	1,000	,494
	Standardfehler	,000	,052
	H	500,000	500,000
Y3	Korrelation	,494	1,000
	Standardfehler	,052	,000
	H	500,000	500,000

Durch das R-Hetcor-Paket berechnete Korrelationen

Korrelationstypen

Variablen	Variablen	
	X3	Y3
X3	--	Polychorisch
Y3	Polychorisch	--

Pearson-, polyseriale und polychorische Korrelationen

Variablen	Statistiken	Variablen	
		X	Y3
X	Korrelation	1,000	,469
	Standardfehler	,000	,043
	H	500,000	500,000
Y3	Korrelation	,469	1,000
	Standardfehler	,043	,000
	H	500,000	500,000

Durch das R-Hetcor-Paket berechnete Korrelationen

Korrelationstypen

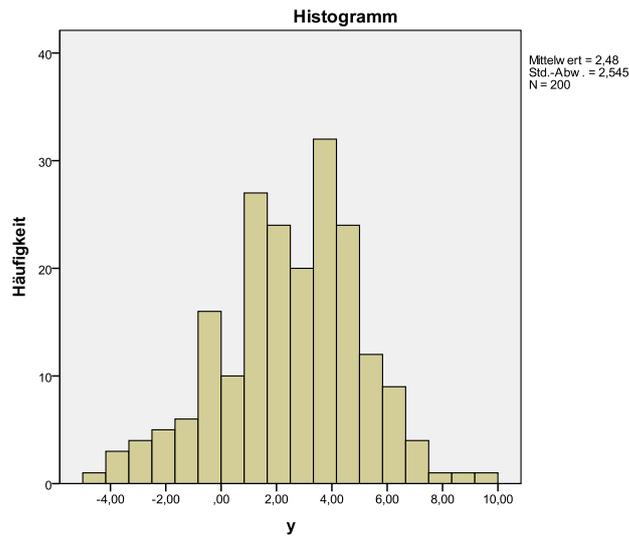
Variablen	Variablen	
	X	Y3
X	--	Polyserial
Y3	Polyserial	--

3.4 Tobit-Regression

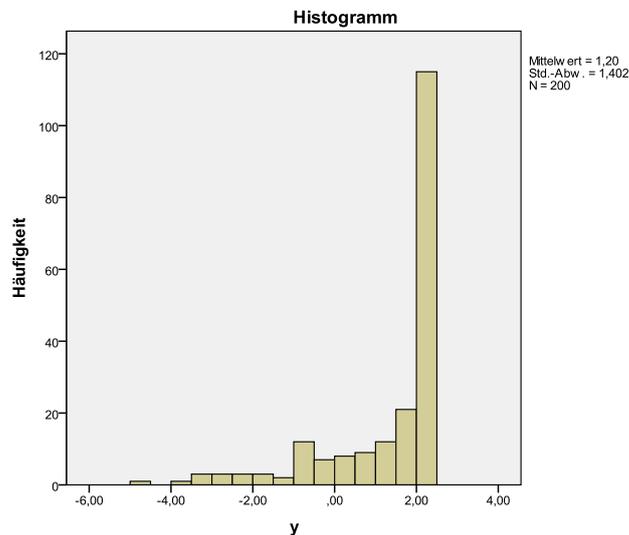
Die Tobit-Regression kommt als Alternative zur OLS-Regression bei zensierten Daten (siehe z.B. IRDE UCLA 2014b) in Frage. Als Anwendungsbeispiel betrachten wir eine bivariate Regression mit dem folgenden wahren Modell (mit Regressionsgewicht 1):

$$Y_i = 1 \cdot X_i + \varepsilon_i \quad \text{mit } \varepsilon_i \sim N(0;4)$$

Aus der Kriteriumsvariablen mit der folgenden Stichprobenverteilung ($n = 200$)



entsteht eine rechts-zensierte Variante, indem alle Werte größer oder gleich 2 auf den Randwert 2 gesetzt werden:



Bei einer OLS-Regression (*Ordinary Least Squares*) mit dem zensierten Kriterium ist die Schätzung des Regressionskoeffizienten erheblich verzerrt:

Koeffizienten^a

Modell		Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten	T	Sig.	95,0% Konfidenzintervalle für B	
		Regressionskoeffizient B	Standardfehler	Beta			Untergrenze	Obergrenze
1	(Konstante)	,098	,181		,539	,590	-,260	,456
	x	,438	,063	,443	6,959	,000	,314	,562

a. Abhängige Variable: y

Nach der Installation der R-Essentials kann in SPSS über den Menübefehl

Analysieren > Regression > Tobit-Regression

ein Dialog zur Berechnung einer Tobit-Regression mit Hilfe von **R**-Funktionen angefordert werden.¹
Für die Beispieldaten resultiert eine Schätzung nahe beim wahren Wert 1:

Koeffizienten

	Koeffizient	Standardfehler	z-Wert	Sig.
(Intercept)	-,101	,353	-,285	,776
x	1,068	,151	7,092	,000
Log(scale)	,806	,084	9,614	,000

Lower bound: Keine, Upper bound: 2
 tobit(formula = y ~ x, left = -Inf, right = 2, dist = "gaussian", data =
 dta, na.action = na.exclude)
 Scale: 2.2380
 Residual d.f.: 197
 Log likelihood: -245.733 D.f.: 3
 Wald statistic: 50.290 D.f.: 1

Neben dem Regressionsgewicht wird unter der Bezeichnung **Scale** noch ein Koeffizient mitgeteilt, der die Standardabweichung der Residuen schätzt und im konkreten Fall nahe bei dem Wert liegt, den eine OLS-Regression für die *unzensierten* Daten ermittelt (Bezeichnung **Standardfehler des Schätzers** in der Tabelle **Modellzusammenfassung**):

Modellzusammenfassung^b

Modell	R	R-Quadrat	Korrigiertes R-Quadrat	Standardfehler des Schätzers
1	,586 ^a	,344	,341	2,06634

a. Einflußvariablen : (Konstante), x

b. Abhängige Variable: y

¹ In **R** sind verschiedene Tobit-Lösungen vorhanden (z.B. in den Paketen **AER**, **VGAM** und **censReg**). SPSS verwendet die Funktion **tobit()** aus dem Paket **AER**.

4 R-Funktionen über das SPSS-Syntaxfenster nutzen

Um **R**-Funktionen in SPSS zu nutzen, die *nicht* über Erweiterungskommandos oder benutzerdefinierte Dialoge erschlossen sind, ist **R**-Syntax zu verfassen. Die Übergabe von SPSS-Variablen an **R** und der Rücktransport von **R**-Ergebnissen fallen aber leicht. Man erstellt in einem SPSS-Syntaxfenster einen Block mit **R**-Syntax, eingerahmt durch die SPSS-Kommandos `BEGIN PROGRAM R` und `END PROGRAM`. In der **R**-Syntax spielen Funktionen aus dem **R**-Paket `spss240`, das als Bestandteil der **R**-Essentials installiert wird, eine wesentliche Rolle, z.B. die Funktion `spssdata.GetDataFromSPSS()` für den Zugriff auf die Variablen der SPSS-Arbeitsdatei. Eine Dokumentation dieser **R**-Funktionen zur SPSS-Unterstützung finden Sie unter der Überschrift

Integration Plug-in for R Help

in der mit dem Menübefehl

Hilfe > Themen

zu startenden SPSS-Hilfe.

Das **R**-Paket zur SPSS-Unterstützung wird bei der **R**-Nutzung via SPSS automatisch geladen, so dass ein explizites Laden mit der **R**-Funktion `library()` (vgl. Abschnitt 5.2.1) *nicht* erforderlich ist.

Bei der in diesem Abschnitt vorgestellten Arbeitsweise sind nur wenige Kenntnisse der **R**-Syntax erforderlich, so dass wir eine Behandlung von **R** als Programmiersprache bis zum Abschnitt 5 aufschieben.

Sollte eine **R**-Anweisung nicht in eine Zeile passen, kann sie auf weiteren Zeilen fortgesetzt werden, was im folgenden Beispiel beim Aufruf der Funktion `boxplot()` geschieht:

```
BEGIN PROGRAM R.  
dta <- spssdata.GetDataFromSPSS(factorMode="labels")  
boxplot(dta$aergo ~ dta$geschlecht, col="lightblue3", varwidth=TRUE, boxwex=0.75,  
        xlab="Geschlecht", ylab="Ärger ohne KFA")  
END PROGRAM.
```

4.1 SPSS-Variablen an R übergeben

4.1.1 Übergabe der kompletten Arbeitsdatei

Mit Hilfe der Funktion `spssdata.GetDataFromSPSS()`, die im **R**-Paket `spss240` enthalten ist, kann man in **R** auf die Variablen der SPSS-Arbeitsdatei zugreifen. Im folgenden Beispiel aus IBM Corp. (2016a, S. 248) werden *alle* Variablen der Arbeitsdatei übertragen:

```

1 DATA LIST FREE /age (F4) income (F8.2) car (F8.2) employ (F4).
2 BEGIN DATA.
3 55 72 36,20 23
4 56 153 76,90 35
5 28 28 13,70 4
6 END DATA.
7
8 BEGIN PROGRAM R.
9 casedata <- spssdata.GetDataFromSPSS()
10 print(casedata)
11 END PROGRAM.

```

In **R** resultiert ein so genannter **Data Frame** (deutsch: *Datentabelle*, vgl. Abschnitt 5.3.4.7), der einem SPSS-Datenblatt weitgehend entspricht.

Der **Data Frame** wird mit dem Operator "<-" der **R**-Variablen (man sagt auch: *dem R-Objekt*) `casedata` zugewiesen.

Über die **R**-Funktion `print()` kann man die in einer Datentabelle vorhandenen Variablen ausgeben lassen.¹ Im Beispiel erhält man im SPSS-Ausgabefenster eine Protokollausgabe der SPSS- und **R**-Kommandos sowie das `print()` - Ergebnis:

```

DATA LIST FREE /age (F4) income (F8.2) car (F8.2) employ (F4).
BEGIN DATA.
55 72 36,20 23
56 153 76,90 35
28 28 13,70 4
END DATA.

BEGIN PROGRAM R.
casedata <- spssdata.GetDataFromSPSS()
print(casedata)
END PROGRAM.

  age income  car employ
1  55    72 36.2    23
2  56   153 76.9    35
3  28    28 13.7     4

```

¹ Die Funktion `print()` wird hier der Deutlichkeit halber explizit notiert. Der Variablenname `casedata` genügt eigentlich, weil er von **R** als impliziter Aufruf der Funktion `print()` verstanden wird.

Deutsche Umlaute in Variablennamen oder Wertbeschriftungen machen bei der Datenübergabe an **R** *keinen* Ärger.

4.1.2 Eine Auswahl von SPSS-Variablen übergeben

Um eine Auswahl der SPSS-Variablen in der aktuellen Arbeitsdatei an **R** zu übergeben, verwendet man in der **R**-Funktion `spssdata.GetDataFromSPSS()` das Argument `variables`, dem ein Vektor mit Variablennamen zu übergeben ist. Diesen Vektor bildet man in der Regel über die **R**-Funktion `c()`. Weil diese Funktion sehr oft benötigt wird, hat sie einen kurzen Namen erhalten, wobei das *c* für *combine* oder *concatenate* steht. Die folgende Syntax:

```
BEGIN PROGRAM R.
casedata <- spssdata.GetDataFromSPSS(variables=c("age","income", "car"))
print(casedata)
END PROGRAM.
```

liefert aufgrund der eben vorgestellten SPSS-Arbeitsdatei die Ausgabe:

```
  age income  car
1  55      72 36.2
2  56     153 76.9
3  28      28 13.7
```

In **R** ist die **Groß-/Kleinschreibung** signifikant und muss auch bei SPSS-Variablennamen beachtet werden. Daher würde das folgende Kommando zu einem Fehler führen:

```
casedata <- spssdata.GetDataFromSPSS(variables=c("Age","income", "car"))
```

Hingegen besteht einige syntaktische Freiheit beim Verfassen der Variablenliste, so dass folgende Varianten erlaubt sind:

a) Eine gemeinsame Zeichenfolge mit Leerzeichen zwischen den Variablennamen:

```
casedata <- spssdata.GetDataFromSPSS(variables=c("age income car"))
```

b) Mit dem Schlüsselwort **TO** bildet man eine Liste von Variablen, die in der Arbeitsdatei (dem aktiven Datenblatt) hintereinander stehen:

```
casedata <- spssdata.GetDataFromSPSS(variables=c("age to car"))
```

Beim Schlüsselwort **TO** ist die Groß-/Kleinschreibung beliebig. Selbstverständlich dürfen mit **TO** gebildete Sequenzen zusammen mit Einzelvariablen aufgelistet werden.

c) Anstelle der Namen kann man die Positionen der Variablen im SPSS-Datenblatt angeben:

```
casedata <- spssdata.GetDataFromSPSS(variables=c(0, 1, 2))
```

Es ist zu beachten, dass die Elemente mit 0 beginnend nummeriert werden.

d) Bei einer Liste von aufeinanderfolgenden Positionen ist der **R**-Sequenzoperator erlaubt (siehe Abschnitt 5.3.4.2.1):

```
casedata <- spssdata.GetDataFromSPSS(variables=0:2)
```

4.1.3 Variablen in einer R-Datentabelle ansprechen

Auf eine einzelne Variable in einer **R**-Datentabelle kann man über die per `$`-Zeichen verbundene Kombination aus dem Datentabellen- und dem Variablennamen zugreifen, z.B.:

```
BEGIN PROGRAM R.
print(casedata$age)
END PROGRAM.
```

4.1.4 Persistenz und Löschen von R-Objekten

Während SPSS große Datensätze nur teilweise im Arbeitsspeicher (RAM) hält, befindet sich ein **R** - Data Frame stets komplett im Speicher.

Das Beispiel in Abschnitt 4.1.3 zeigt außerdem, dass der **R**-Workspace (mit allen Variablen) innerhalb einer SPSS-Sitzung zwischen zwei **R**-Einsätzen erhalten bleibt, d.h.:

- Die in früheren Blöcken von **R**-Syntax angelegten Objekte stehen weiter zur Verfügung.
- Eventuell ist es sinnvoll, belegten Speicherplatz (im RAM) durch explizites Entfernen von Objekten mit der **R**-Funktion **rm()** (alias **remove()**) wieder frei zu geben, z.B.:

```
BEGIN PROGRAM R.
```

```
  . . .
  rm(casedata)
  gc()
  . . .
```

```
END PROGRAM.
```

Durch einen Aufruf der Funktion **gc()** (*Garbage Collector*, dt.: *Müllsammler*) wird sichergestellt, dass **R** den zuvor von gelöschten Objekten belegten Speicherplatz sofort an das Betriebssystem zurückgibt.

Weil SPSS nicht darauf angewiesen ist, alle Daten im Hauptspeicher (RAM) zu halten, müssen SPSS-Anwender nicht über das Einsparen von Speicherplatz nachdenken. **R** hingegen muss die zu analysierenden Daten komplett im Hauptspeicher halten, so dass es empfehlenswert ist (vgl. Muenchen 2011, S. 33),

...

- SPSS-seitig ein an **R** zu übergebendes Datenblatt möglichst klein zu halten (z.B. durch das Löschen von überflüssigen Variablen und/oder Fällen),
- überflüssig gewordene **R**-Objekte zu löschen.

4.1.5 Kategoriale SPSS-Variablen als (ordinale) Faktoren an R übergeben

Ein (ordinaler) Faktor in **R** (siehe Abschnitt 5.3.4.3) ist das Analogon zu einer kategorialen (nominalen oder ordinalen) SPSS-Variablen.

SPSS-Variablen mit dem Datentyp Zeichenfolge (String) werden generell in der **R**-Datentabelle als Faktoren abgelegt. Aus numerischen SPSS-Variablen resultieren per Voreinstellung in **R** numerische Vektoren. Das gilt auch für numerische SPSS-Variablen mit dem deklarierten Messniveau Nominal oder Ordinal.

Durch Verwendung des Arguments **factorMode** im Aufruf der Funktion **spssdata.GetDataFromSPSS()** veranlasst man, dass

- nominale numerische SPSS-Variablen in **R** zu Faktoren und
- ordinale numerische SPSS-Variablen in **R** zu *geordneten* Faktoren werden.

Mit dem zum Argument **factorMode** anzugebenden Wert legt man fest, ob bei einer kategorialen numerischen SPSS-Variablen die Werte oder die Wertetiketten in **R** zur Beschriftung der Faktorstufen verwendet werden sollen:

factorMode -Parameter	In R werden die Faktorstufen beschriftet über ...
labels	die Wertetiketten der SPSS-Variablen
levels	die Werte der SPSS-Variablen

Aus der folgenden Syntax

```

DATA LIST FREE /y (F2) gruppe (F2) stufe (F2).
BEGIN DATA.
8 1 1
7 1 1
3 2 2
1 2 3
END DATA.
VARIABLE LEVEL y (SCALE) /gruppe (NOMINAL) /stufe(ORDINAL).
VALUE LABELS gruppe 1 'Jung' 2 'Alt' /stufe 1 'A' 2 'B' 3 'C'.

```

```

BEGIN PROGRAM R.
casedata <- spssdata.GetDataFromSPSS(factorMode="labels")
class(casedata$gruppe)
class(casedata$stufe)
str(casedata$gruppe)
END PROGRAM.

```

resultiert die Ausgabe:

```

[1] "factor"
[1] "ordered" "factor"
Factor w/ 2 levels "Jung","Alt": 1 1 2 2

```

Per `class()` - Funktion erhält man die Auskunft, dass `casedata$gruppe` ein Faktor und `casedata$stufe` ein geordneter Faktor ist. Der `str()` - Ausgabe ist u.a. zu entnehmen, dass die Kategorien des **R**-Faktors `casedata$gruppe` gemäß `factorMode`-Spezifikation durch die SPSS-Wertelabels beschriftet sind.

4.1.6 Indikatoren für fehlende Werte

Bei numerischen SPSS-Variablen werden per Voreinstellung benutzerdefinierte Missing Data - Indikatoren (ab jetzt kurz als *MD-Indikatoren* bezeichnet) ebenso wie `SYSMIS` in der **R**-Datentabelle zum Wert **NaN** (*Not a Number*). Bei alphanumerischen SPSS-Variablen resultiert in **R** der Wert **NA** (*Not Available*). Die benutzerdefinierten MD-Indikatoren müssen ersetzt werden, weil es in **R** *nicht* möglich ist, beliebige Werte als Indikatoren für fehlende Werte zu deklarieren.

Das folgende Beispielprogramm

```

DATA LIST LIST (' ','') /numvar (F2) strvar (A1).
BEGIN DATA.
9,A
7,M
,B
END DATA.
MISSING VALUES numvar (9) /strvar ('M').

BEGIN PROGRAM R.
casedata <- spssdata.GetDataFromSPSS()
print(casedata)
END PROGRAM.

```

liefert die Ausgabe:

	numvar	strvar
1	NaN	A
2	7	<NA>
3	NaN	B

Setzt man, wie es z.B. von Muenchen (2011, S. 34) empfohlen wird, in der Funktion `spssdata.GetDataFromSPSS()` das optionale Argument `missingValueToNA` auf den Wert `TRUE`,

```
BEGIN PROGRAM R.
casedata <- spssdata.GetDataFromSPSS(missingValueToNA=TRUE)
print(casedata)
END PROGRAM.
```

dann resultiert bei fehlenden Werten von numerischen SPSS-Variablen in der **R**-Datentabelle der Wert **NA** (*Not Available*):

	numvar	strvar
1	NA	A
2	7	<NA>
3	NA	B

In der Regel werden in **R** bei numerischen Variablen die beiden Werte **NA** und **NaN** gleich behandelt, und die Funktion `is.na()` liefert für beide Werte ein `TRUE`. Wer gezielt auf den Wert **NaN** prüfen will, hat die **R**-Funktion `is.nan()` zur Verfügung, die bei **NaN** ein `TRUE` und bei **NA** ein `FALSE` liefert.

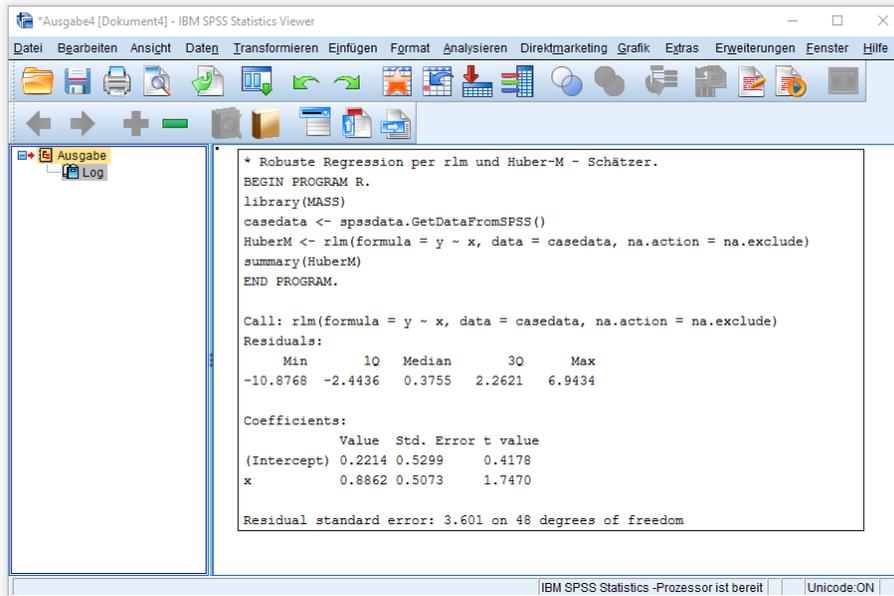
4.2 *R*-Auswertungsfunktionen verwenden und Ausgaben im SPSS-Viewer anzeigen

In Abschnitt 3.1 haben wir über den Dialog eines Erweiterungsbundles eine robuste Regressionsanalyse mit der **R**-Funktion `rlm()` aus dem Paket **MASS** durchgeführt und dabei den in der implementierenden **R**-Syntax verwendeten biquadratischen M-Schätzer akzeptiert. Wer stattdessen Hubers M-Schätzer bevorzugt, kann in einem expliziten **R**-Syntax - Block die Variablen der SPSS-Arbeitsdatei an **R** übergeben und die **R**-Funktion `rlm()` direkt verwenden, die per Voreinstellung mit dem M-Schätzer von Huber arbeitet:

```
BEGIN PROGRAM R.
library(MASS)
casedata <- spssdata.GetDataFromSPSS()
huber <- rlm(y ~ x, data = casedata)
summary(huber)
END PROGRAM.
```

Benötigte **R**-Funktionen befinden sich oft in Paketen, die zunächst über einen Aufruf der Funktion `library()` geladen werden müssen (siehe Beispiel).

Wie schon mehrfach zu sehen war, erscheint die Standardausgabe von **R**-Funktionen in einen **Log**-Abschnitt des SPSS-Ausgabefensters:



Im Programm wird die **rlm()** - Ausgabe in die Variable `huber` geleitet und von dort per **summary()** - Funktion in einiger Ausführlichkeit hervorgehoben. Auf direktem Weg

```
rlm(y ~ x, data = casedata)
```

erhält man nur eine spartanische Ausgabe:

Call:

```
rlm(formula = y ~ x, data = casedata, na.action = na.exclude)
```

Converged in 4 iterations

Coefficients:

```
(Intercept)      x
 0.2214148    0.8861699
```

Degrees of freedom: 50 total; 48 residual

Scale estimate: 3.6

Die zurückhaltende Originalausgabe von statistischen Auswertungsfunktionen und die Verwendung der Extraktionsfunktion **summary()** sind typisch für **R** (Muenchen 2011, S. 99f).

Die Koeffiziententabelle der **rlm()** - Ausgabe enthält (wie die in Abschnitt 3.1 präsentierte Pivot-Tabelle zum korrespondierenden Erweiterungskommando) *keine* Überschreitungswahrscheinlichkeiten (*p*-Values) zu den Signifikanztests. Bei genügend Vertrauen in die approximative Normalverteilung der Prüfstatistiken kann man mit folgender **R**-Syntax die Überschreitungswahrscheinlichkeiten für einen ein- bzw. zweiseitigen Test gegen den theoretischen Wert 0 bestimmen (siehe Bellio & Ventura 2005, S. 16):

```

BEGIN PROGRAM R.
pst <- 1.747
psingle <- min(1-pnorm(pst), pnorm(pst))
psingle
2 * psingle
END PROGRAM.
    
```

Im Beispiel lehnt der *einseitige* Test zum Regressor `x` (bei einem erwartungskonformen Vorzeichen des Koeffizienten) seine Nullhypothese ab, der *zweiseitige* hingegen nicht:

```
[1] 0.04031867
[1] 0.08063734
```

Manche **R**-Ausgaben lassen sich mit der Funktion `spsspivottable.Display()` aus dem **R**-Paket zur SPSS-Unterstützung, das mit den **R**-Essentials installiert wird, in eine SPSS-Pivot-Tabelle wandeln. Nach einer geringfügigen Erweiterung des **R**-Programms mit `rlm()` - Aufruf

```
BEGIN PROGRAM R.
library(MASS)
casedata <- spssdata.GetDataFromSPSS()
huber <- rlm(formula = y ~ x, data = casedata)
results <- summary(huber)
spsspivottable.Display(results$coefficients,
                        title = "Koeffizienten",
                        format = formatSpec.GeneralStat)

END PROGRAM.
```

erhalten wir als Ausgabe die folgende Pivot-Tabelle:

	Value	Std. Error	t value
(Intercept)	,221	,530	,418
x	,886	,507	1,747

Statistische Auswertungsfunktionen in **R** organisieren ihre Ausgabe in der Regel als Liste mit Komponenten unterschiedlichen Typs. Wir behandeln diesen flexiblen **R**-Datentyp in Abschnitt 5.3.4.6. Um zu erfahren, welche Komponenten in einer Listenausgabe vorhanden und für die Wandlung in eine Pivot-Tabelle verfügbar sind, kann man die `str()` - Funktion verwenden, z.B.:

```
str(results)
```

Im Beispiel hat sich herausgestellt, dass `results` eine Komponente namens `coefficients` enthält. Diese Komponente wird in obigem `spsspivottable.Display()` - Aufruf als erstes Argument verwendet. Ihm folgen ein frei wählbarer Titel und die Formatangabe `formatSpec.GeneralStat`, die in den meisten Fällen zu einem sinnvollen Ergebnis führt.

4.3 SPSS-Variablen mit R erstellen

Das folgende Programm nach IBM Corp. (2016a, S. 253) berechnet für eine SPSS-Variable mit Hilfe der **R**-Funktion `mean()` den Mittelwert aller Fälle und erstellt in SPSS ein um die Mittelwertsvariable erweitertes Datenblatt.

```
DATA LIST FREE /salary.
BEGIN DATA
1
2
3
4
END DATA.

VARIABLE LEVEL salary (SCALE).
```

```

BEGIN PROGRAM R.
dict <- spssdictionary.GetDictionaryFromSPSS()
casedata <- spssdata.GetDataFromSPSS()
varSpec <- c("meansal", "Mean Salary", "0", "F8.2", "scale")
dict <- data.frame(dict, varSpec, stringsAsFactors=FALSE)
spssdictionary.SetDictionaryToSPSS("results", dict)
casedata <- data.frame(casedata, mean(casedata$salary))
spssdata.SetDataToSPSS("results", casedata)
spssdictionary.EndDataStep()
END PROGRAM.

```

Mit der **R**-Funktion `spssdictionary.GetDictionaryFromSPSS()` werden die Variablendeklarationen aus der SPSS-Arbeitsdatei in eine **R**-Datentabelle mit dem Namen `dict` übernommen:

```
dict <- spssdictionary.GetDictionaryFromSPSS()
```

Um eine neue Variablendeklaration zu ergänzen, wird zunächst ein Vektor vom Datentyp **character** erstellt und im **R**-Objekt `varSpec` abgelegt:

```
varSpec <- c("meansal", "Mean Salary", "0", "F8.2", "scale")
```

Hier müssen 5 SPSS-Variablenattribute in vorgeschriebener Reihenfolge jeweils durch eine Zeichenkette festgelegt werden:

- **Variablenname**
Im Beispiel: "meansal"
- **Variablenbeschriftung**
Bei Verzicht auf eine Beschriftung ist eine leere Zeichenkette anzugeben ("").
Im Beispiel: "Mean Salary"
- **Variablentyp**
Für numerische Variablen ist eine 0 anzugeben, für Zeichenkettenvariablen die Länge (von 1 bis 32767).
Im Beispiel: "0"
- **Anzeigeformat**
Im Beispiel: "F8.2"
- **Messniveau**
Erlaubte Werte: **nominal**, **ordinal**, **scale**
Im Beispiel: "scale"

Im Beispielprogramm wird der Data Frame `dict` über die zum **R**-Kern gehörige Funktion `data.frame()` (vgl. Abschnitt 5.3.4.7.1) um die Variablendeklaration im **character**-Vektor `varSpec` erweitert:

```
dict <- data.frame(dict, varSpec, stringsAsFactors=FALSE)
```

Eine **R**-Datentabelle zum Speichern der Attribute zu den Variablen eines SPSS-Datenblatts enthält für jede SPSS-Variable einen Vektor mit Werten von Typ **character**. Wir erweitern das Datenlexikon `dict` um die Variablendeklaration im **character**-Vektor `varSpec` und verhindern mit dem Argument **FALSE** für das `data.frame()` - Argument `stringsAsFactors`, dass der **character**-Vektor bei der Aufnahme in einen Faktor (im Sinne von **R**) gewandelt wird.¹

Aus dem Ergebnis erstellt die Funktion `spssdictionary.SetDictionaryToSPSS()`

```
spssdictionary.SetDictionaryToSPSS("results", dict)
```

ein neues SPSS-Datenblatt mit dem Namen `results`.

¹ Die voreingestellte Wandlung wäre für die weitere Verarbeitung im konkreten Beispielprogramm kein Problem, doch sollte generell die Struktur einer **R**-Datentabelle zur Aufbewahrung von SPSS-Variablenattributen respektiert werden.

Die **R**-Datentabelle `casedata` mit den SPSS-Variablen wird im folgenden `data.frame()` - Aufruf um die per `mean()` gebildete Mittelwertvariable erweitert:¹

```
casedata <- data.frame(casedata, mean(casedata$salary))
```

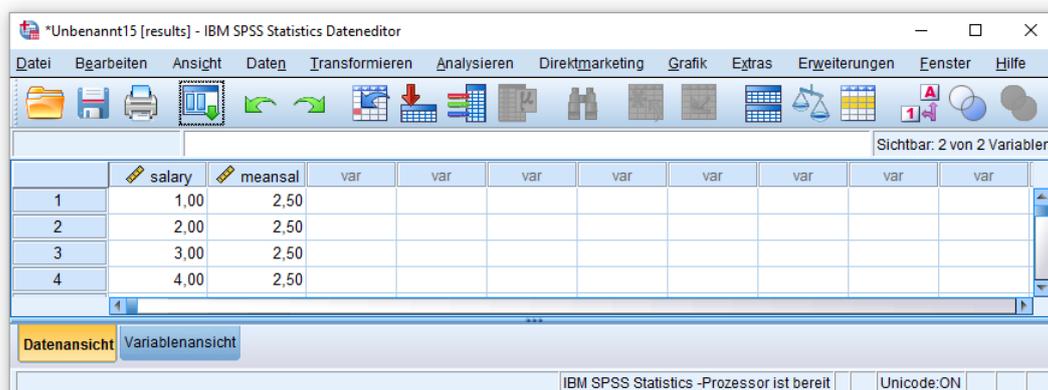
Den von der **R**-Funktion `mean()` gelieferten numerischen Vektor der Länge 1 erweitert **R** automatisch durch Wiederholung des vorhandenen Elements auf die passende Länge.

Mit der Funktion `spssdata.SetDataToSPSS()` werden die Variablen im erweiterten Data Frame `casedata` in das eben angelegte SPSS-Datenblatt geschrieben:

```
spssdata.SetDataToSPSS("results", casedata)
```

Am Ende der Datenblatt-Erstellung sollte die Funktion `spssdictionary.EndDataStep()` aufgerufen werden.

Das Ergebnis:



Mit der Mittelwertbildung wurde ein möglichst einfaches Beispiel gewählt, um die Erstellung von SPSS-Variablen durch **R** zu demonstrieren. Diese Aufgabe ist mit SPSS-internen Mitteln einfacher zu lösen:

```
AGGREGATE
  /OUTFILE=* MODE=ADDVARIABLES
  /meansal=MEAN(salary).
```

Es kann sich aber leicht die Situation ergeben, dass zur Lösung einer Datentransformationsaufgabe **R**-Syntax in Frage kommt, z.B. weil ein in **R** implementierter Algorithmus übernommen werden soll.

¹ Im Hinblick auf den Abschnitt 5.3.4.7 soll darauf hingewiesen werden, dass im Beispiel auf das explizite Benennen der neuen Variablen in der **R**-Datentabelle verzichtet wird, so dass eine automatische Namensvergabe stattfindet. Das Ergebnis verrät der **R**-Funktionsaufruf `str(casedata)`:

```
'data.frame': 4 obs. of 2 variables:
 $ salary      : num  1 2 3 4
 $ mean.casedata.salary.: num  2.5 2.5 2.5 2.5
```

5 R als statistikorientierte Programmierumgebung

Mit zunehmender Komplexität der **R**-Beispiele in den obigen Abschnitten ist vermutlich die Motivation der Leser(innen) gewachsen, sich mit *R als Programmierumgebung für Datenanalyse und Grafik* (dt. Übersetzung des Untertitels von Venables et al. 2016) systematisch zu beschäftigen, um Sicherheit bei der Anwendung zu gewinnen.

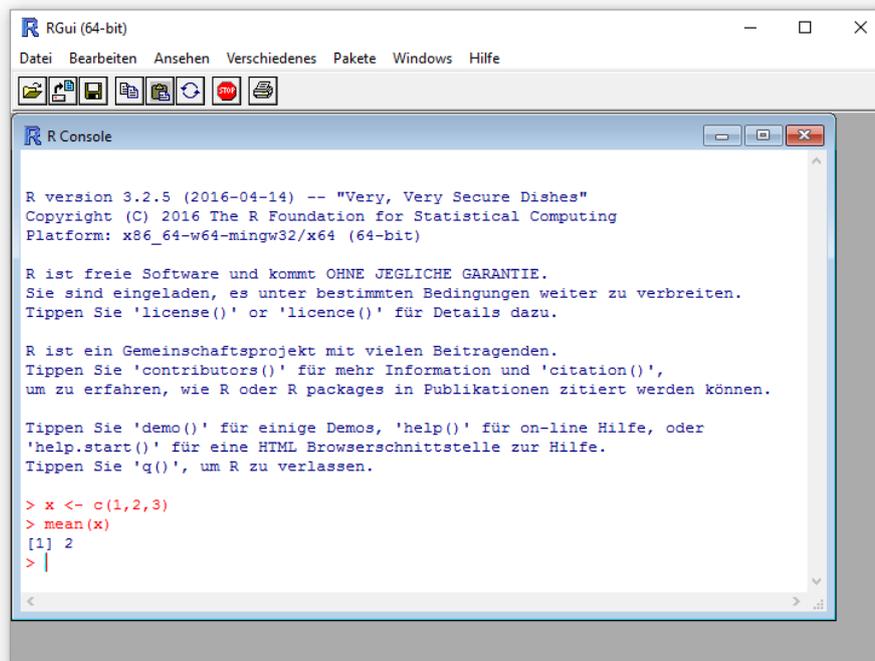
5.1 RGui zur direkten Interaktion mit R

Im aktuellen Abschnitt 5 werden die **R**-Funktionen der Kürze halber in der Regel *ohne* Einrahmung durch die SPSS-Kommandos

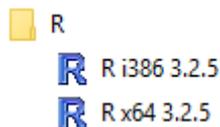
```
BEGIN PROGRAM R.
. . .
END PROGRAM.
```

gezeigt.

Zum Erlernen der **R**-Syntax ist die direkte Interaktion mit dem **R**-Interpreter über die grafische **R**-Bedienoberfläche (RGui) ohnehin sinnvoller:¹



Bei einer Windows-Version mit 64-Bit - Architektur steht das RGui als 32- und 64-Bit - Anwendung bereit. Mittlerweile ist auch die 64-Bit - Version ausgereift und sollte in der Regel bevorzugt werden. Hier sind die Startmenüeinträge für **R** 3.2.5 zu sehen:



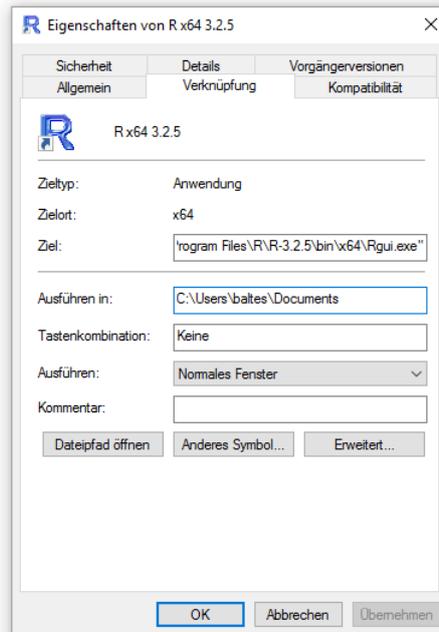
Die im RGui üblichen Farben für Eingaben (rot) und Ausgaben (blau) werden auch im Manuskript verwendet.

¹ Das populäre und mächtige **R**-Studio als Bedienoberfläche zu verwenden, wäre für unsere Zwecke übertrieben.

Weil im Konsolenfenster der RGui-Bedienoberfläche sämtliche Eingaben und **R**-Antworten protokolliert werden, ist gelegentlich sinnvoll, per Kontextmenü oder mit dem Tastenbefehl **Strg+L** für eine leere **R**-Konsole zu sorgen.

5.1.1 Arbeitsverzeichnis

In einer R-Sitzung dient das **Arbeitsverzeichnis** (engl.: *working directory*) als Voreinstellung beim Lesen und Schreiben von Dateien. Initial ist es identisch mit dem aktuellen Verzeichnis beim Start von **R** (Startverzeichnis). Startet man das RGui über seine Verknüpfung im Startmenü, hängt das Startverzeichnis vom Feld **Ausführen in** des Eigenschaftsdialogs zur Verknüpfung ab, z.B.:



Man kann in der **R**-Konsole das aktuelle Arbeitsverzeichnis mit der Funktion **getwd()** ermitteln, z.B.:

```
> getwd()
[1] "C:/Users/baltes/Documents"
```

Um das Arbeitsverzeichnis zu wechseln, kann man den RGui-Menübefehl

Datei > Verzeichnis wechseln

oder die **R**-Funktion **setwd()** verwenden, z.B.:

```
> setwd("U:/Eigene Dateien/R")
```

Weil in **R** (wie in vielen anderen Programmiersprachen) der Rückwärtsschrägstrich als Einleitungszeichen für Escape-Sequenzen reserviert ist (z.B. **\n** für den Zeilenwechsel), muss bei Windows-Pfadangaben ersatzweise ein doppelter Rückwärts-Schrägstrich oder ein Vorwärts-Schrägstrich verwendet werden.

Mit der Funktion **dir()** fordert man eine Liste der Dateisystemobjekte im Arbeitsverzeichnis an:

```
> dir()
```

5.1.2 Workspace und Anweisungsgedächtnis

Die in einer Sitzung erzeugten **R**-Objekte (z.B. Vektoren und Datentabellen) landen im so genannten **Workspace**. Mit der Funktion `ls()` kann man sich die Objekte im Workspace auflisten lassen, z.B.:

```
> ls()
[1] "eimer" "x"
```

Der etwas längere Aliasname für diese Funktion lautet `objects()`.

Über die `apropos()` - Funktion lassen sich alle Objekte ermitteln, die einen bestimmten Namensbestandteil besitzen, wobei aber nicht nur der Workspace inspiziert wird, sondern der gesamte Suchpfad der **R**-Sitzung (siehe Abschnitt 5.3.3.2). Die folgenden Kommandos

```
> eimer <- c(1, 2, 3)
> apropos("eim")
```

führen bei **R** 3.2.5 zur Ausgabe:

```
[1] ".mergeImportMethods" "eimer" "getNamespaceImports"
[4] "namespaceImport" "namespaceImportClasses" "namespaceImportFrom"
[7] "namespaceImportMethods"
```

Man kann die Funktion `apropos()` im RGui auch per Menübefehl nutzen:

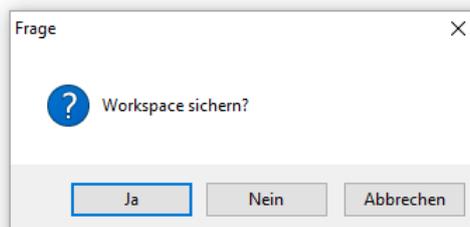
Hilfe > Apropos

Neben den Datenobjekten merkt sich das RGui auch die abgeschickten Anweisungen. Über die vertikalen Pfeiltasten kann man früher verwendete Anweisungen zurückholen.

Das RGui verwendet die folgenden Dateien im Arbeitsverzeichnis, um den Workspace und das Kommandogedächtnis zwischen den Sitzungen aufzubewahren:

- **.RData**
Hier speichert das RGui den Workspace.
- **.Rhistory**
Hier speichert das RGui das Kommandogedächtnis.

Beim Starten liest **R** den Workspace und das Kommandogedächtnis aus diesen Dateien, und beim Beenden einer Sitzung werden Datenobjekte sowie gespeicherte Anweisungen der Sitzung nach einer zustimmenden Antwort auf die folgende Frage



in die beiden Dateien gesichert.

Um alle Datenobjekte der aktuellen Sitzung in eine *wählbare* Datei mit der Namenserweiterung **.RData** zu sichern bzw. von dort zu laden, kann man die folgenden Menübefehle

- **Datei > Sichere Workspace**
- **Datei > Lade Workspace**

verwenden, oder mit Kommandos arbeiten, z.B.:

- `> save.image("ws.RData")`
- `> load("ws.RData")`

Vor dem Speichern aller Objekte kann es sinnvoll sein, mit der Funktion `rm()` überflüssige Objekt aus dem Workspace zu entfernen, z.B.:

```
> rm(v1, v2, matz)
```

Um das Kommandogedächtnis der aktuellen Sitzung in eine wählbare Datei zu sichern bzw. von dort zu laden, kann man die Menübefehle

- **Datei > Speichere History**
- **Datei > Lade History**

verwenden, oder mit Kommandos arbeiten, z.B.:

- `> savehistory("komm.Rhistory")`
- `> loadhistory("komm.Rhistory")`

Alle per `save.image()` oder `savehistory()` ohne Pfadangabe geschriebenen Dateien landen im Arbeitsverzeichnis. Wie man es ermittelt oder verändert, wurde im Abschnitt 5.1.1 erläutert.

5.1.3 Sichern und Laden einzelner Datenobjekte im Binärformat von R

Eine **RData**-Datei kann in **R** als binärformatige Datendatei analog zu einer **SAV**-Datei in SPSS genutzt werden. Oft ist es sinnvoll, ein wichtiges Datenobjekt (z.B. eine Datentabelle mit den Variablen einer Studie) in einer eigenen **RData**-Datei zu speichern. Dazu eignet sich die Funktion `save()`, z.B.:

```
> save(daten, file="daten.RData")
```

Man kann auch *mehrere* Datenobjekte in eine Datei befördern, z.B.:

```
> save(tabelle1, tabelle2, file="daten.RData")
```

Mit der Funktion `load()` befördert man die in einer Datei befindlichen Datenobjekte in den Workspace der aktuellen Sitzung:

```
> load("daten.RData")
```

Eine bequeme Möglichkeit zum Öffnen einer **RData**-Datei besteht darin, die Datei per Drag & Drop von einem Fenster des Windows-Explorers auf das Fenster der **R**-Konsole zu bewegen.

5.1.4 Konfigurationsoptionen

Diverse Verhaltensmerkmale von **R** lassen sich über die Funktion `options()` beeinflussen. Ein Funktionsaufruf ohne Argumente hat eine (längliche) Auflistung sämtlicher Optionen zur Folge. Übergibt man eine Option als Zeichenfolgen-Argument, erfährt man ihren aktuellen Wert, z.B.:¹

```
> options("digits")
$digits
[1] 7
```

¹ Man erhält eine Liste der Länge 1. Vom eigentlich zu bevorzugenden Funktionsaufruf `getOption()` erhält man einen Einzelwert, z.B.:

```
> getOption("digits")
[1] 7
```

Um die meisten Optionen muss man sich wegen sinnvoller Voreinstellungen nicht kümmern. Anschließend werden Optionen vorgestellt, bei denen sich eine Änderung lohnen könnte. Soll eine Einstellungsänderung bei jedem **R**-Start gültig sein, fügt man den zugehörigen **options()** - Aufruf in eine Initialisierungsdatei ein (siehe Abschnitt 5.1.5).

Genauigkeit der Ergebnisausgabe

Per Voreinstellung werden Berechnungsergebnisse mit 7 Stellen Genauigkeit angezeigt, z.B.:

```
> log(2)
[1] 0.6931472
```

Über das Argument **digits** der Funktion **options()** lässt sich eine alternative Anzeigegenauigkeit einstellen, z.B.:

```
> options(digits=15)
> log(2)
[1] 0.693147180559945
```

Automatisch geladene Pakete

Alle **R**-Funktionen befinden sich in Paketen und können erst nach dem Laden ihres Pakets genutzt werden (siehe Abschnitt 5.2). Über die Option **defaultPackages** legt man fest, welche Pakete automatisch geladen werden sollen. Hier ist die Voreinstellung für diese Option zu sehen:

```
> options("defaultPackages")
$defaultPackages
[1] "datasets" "utils" "grDevices" "graphics" "stats" "methods"
```

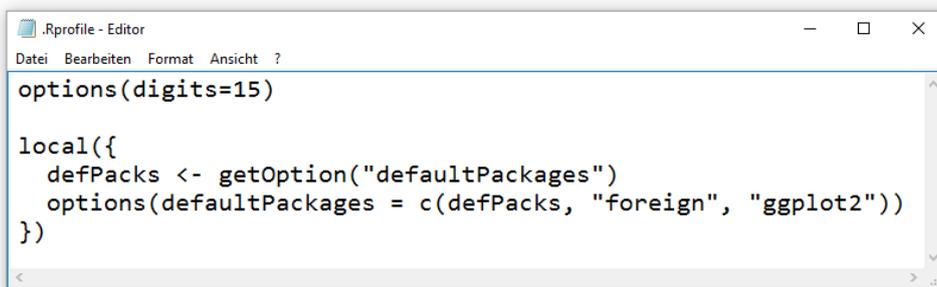
Eine Änderung dieser Option macht nur im Rahmen einer beim **R**-Start automatisch ausgeführten Initialisierungsdatei Sinn (siehe nächsten Abschnitt).

5.1.5 Initialisierungsdateien

Als Initialisierungsdateien mit Anweisungen, die beim Programmstart ausgeführt werden sollen, fungieren:

- **Rprofile.site** im **etc**-Unterordner des **R**-Installationsordners (z.B. **C:\Program Files\R\R-3.2.5\etc**)
- **.Rprofile** im Startverzeichnis (mit einem Punkt als erstem Zeichen im Dateinamen)

Mit der folgenden Datei **.Rprofile** im Startverzeichnis wird eine Ausgabegenauigkeit von 15 Dezimalstellen angeordnet und die Liste der automatisch zu ladenden Pakete erweitert:¹



```
.Rprofile - Editor
Datei Bearbeiten Format Ansicht ?
options(digits=15)

local({
  defPacks <- getOption("defaultPackages")
  options(defaultPackages = c(defPacks, "foreign", "ggplot2"))
})
```

¹ Damit sich die Änderung der Option **defaultPackages** auf den **R**-Start auswirkt, wird sie im Rahmen der **local**-Funktion vorgenommen (Muenchen 2011, S. 659). Zur Ermittlung des Ausgangszustands wird die Funktion **getOption()** verwendet, die im Unterschied zu der in Abschnitt 5.1.4 verwendeten Funktion **options()** keine Liste liefert, sondern einen Vektor mit Elementen vom Typ **character**.

5.2 Pakete

R-Funktionen befinden sich in Paketen, die geladen sein müssen, um die darin enthaltenen Funktionen nutzen zu können. Welche Pakete aktuell geladen sind, erfährt man durch einen Aufruf der Funktion `.packages()`, z.B.:

```
> (.packages())
[1] "ggplot2"   "foreign"   "stats"     "graphics"  "grDevices" "utils"
[7] "datasets" "methods"   "base"
```

In der etwas gewöhnungsbedürftigen Syntax des Beispiels sorgen die runden Klammern um den Funktionsaufruf dafür, dass seine Rückgabe (ein Vektor mit Elementen vom Typ **character**) ausgegeben wird. Einige Pakete (darunter stets das Paket **base**) werden automatisch geladen. Welche Pakete neben **base** automatisch geladen werden, kann man über die Option **defaultPackages** erfahren und festlegen (vgl. Abschnitt 5.1.4).

Alternativ zu `.packages()` kann man die Funktion `search()` dazu benutzen, die geladenen Pakete aufzulisten, z.B.:

```
> search()
[1] ".GlobalEnv"           "package:ggplot2"     "package:foreign"
[4] "package:stats"        "package:graphics"    "package:grDevices"
[7] "package:utils"        "package:datasets"    "package:methods"
[10] "Autoloads"            "package:base"
```

Die eigentliche Aufgabe der Funktion `search()` besteht darin, den Suchpfad für **R**-Objekte aufzulisten (siehe Abschnitt 5.3.3.2), der auch die geladenen Pakete enthält. Im RGui unter Windows lässt sich der `search()` - Aufruf auch mit dem folgenden Menübefehl auslösen:

Verschiedenes > Liste Suchpfad auf

5.2.1 Pakete laden

Zum Laden eines Pakets taugt die Funktion `library()`, z.B.:

```
> library(MASS)
```

Ein Funktionsaufruf ohne Argumente

```
> library()
```

führt zu einer Liste mit allen *installierten* Paketen ...

- in der installations-allgemeinen Bibliothek, z.B.: `C:\Program Files\R\R-3.2.5\library`
- und in der persönlichen Bibliothek, z.B. `C:\Users\baltes\Documents\R\win-library\3.2`

Über das Argument **lib.loc** ist es möglich, ein Paket aus einem beliebigen Ordner zu laden, z.B.:

```
> library(mice, lib.loc="E:\\Daten\\R\\library\\3.2 C:\Users\baltes\Documents\R\win-library\3.2")
```

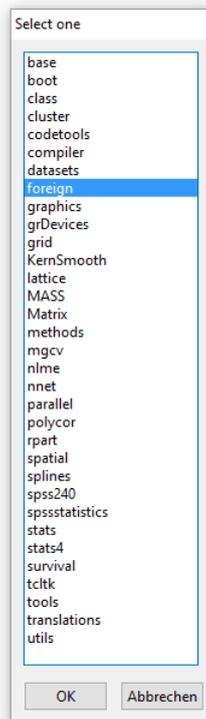
Beinhalten mehrere geladene Pakete namensgleiche Funktionen, gewinnt das zuletzt geladene Paket. Über den Operator `::` ist es aber möglich, die Funktion aus einem bestimmten Paket anzusprechen, z.B.

```
> mypack::func()
```

Um ein Paket *per Dialogbox* zu laden, kann man im RGui den Menübefehl

Pakete > Lade Paket

verwenden, aus der Liste mit allen installierten Paketen ein Exemplar wählen und mit **OK** quittieren:



In den meisten Fällen ist jedoch das Laden von Paketen per **library()** - Funktion sinnvoller (z.B. im Rahmen von Skripten).

Zum Laden von Paketen per Syntax taugt auch die Funktion **require()**, die im Vergleich zu **library()** folgende Vorteile besitzt:

- Scheitert das Laden, produziert **library()** eine Fehlermeldung, **require()** hingegen eine Warnung. Beim Aufruf innerhalb einer Funktion führt eine Fehlermeldung zum Abbruch der Funktion, eine Warnung hingegen nicht.
- **require()** liefert eine boolesche Rückgabe zum Erfolg des Aufrufs und kann daher gut in eine bedingte Anweisung integriert werden (vgl. Abschnitt 5.3.8).

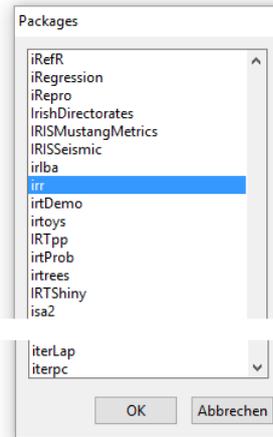
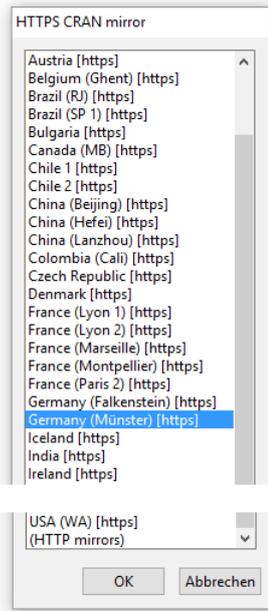
5.2.2 Pakete installieren

Bei der Arbeit mit **R** ist es häufig erforderlich, Zusatzpakete zu installieren, die bestimmte Auswertungsverfahren implementieren. Wir wollen in Abschnitt **Fehler! Verweisquelle konnte nicht gefunden werden.** das (in SPSS nicht verfügbare) **Fleiss-Kappa** zur Beurteilung der Inter-Rater-Übereinstimmung bei mehr als 2 Diagnostikern berechnen lassen und benötigen dazu das **R**-Paket **irr** (*Inter Rater Reliability*), das nicht zum Standardumfang einer **R**-Installation gehört.

Bei bestehender Internet-Verbindung kann die Ergänzungsinstallation bequem im RGui über den folgenden Menübefehl gestartet werden:

Pakete > Installiere Paket(e)

Es erscheint ein Dialog zur Wahl eines Spiegel-Servers zum *Comprehensive R Archive Network* (CRAN) (links). Danach wählt man das gewünschte Paket (rechts):



Wer lieber mit Kommandos arbeitet, wählt zur Installation eines Paketes die Funktion **install.packages()**, z.B.:

```
> install.packages("irr")
```

Auch bei diesem Einstieg fragt das RGui nach dem bevorzugten Spiegel-Server. Um diese Anfrage zu vermeiden, kann man im Aufruf der Funktion **install.packages()** über das **repos**-Argument einen Server angeben, z.B.:

```
> install.packages("irr", repos='http://cran.us.r-project.org')
```

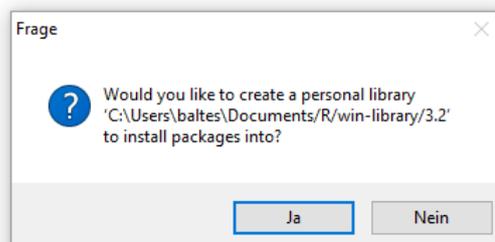
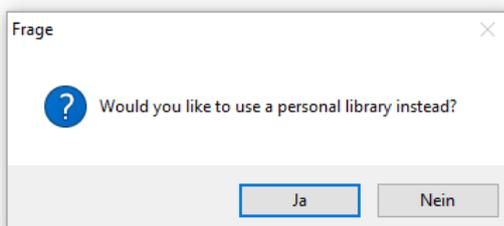
Ist auf dem Zielrechner keine Internet-Verbindung vorhanden, kann man die benötigten Pakete auf einem lokalen Datenträger als ZIP-Dateien bereithalten und über den folgenden Menübefehl installieren:

Pakete > Installiere Paket(e) aus lokalen Zip-Dateien

Sofern entsprechende Schreibrechte bestehen, werden neue Pakete im **library** - Unterordner des **R**-Installationsordners abgelegt, und die Pakete stehen allen Benutzern zur Verfügung. Bei einer Installation von **R** 3.2.5 unter Windows landen die Pakete per Voreinstellung hier:

C:\Program Files\R\R-3.2.5\library

Bestehen bei der Installation eines Zusatzpakets keine Schreibrechte für diesen Ordner, wird die Einrichtung einer **persönlichen Bibliothek** vorgeschlagen:



Vor der geplanten Installation von Paketen, die in der allgemeinen Bibliothek (im Programmordner) landen und damit für alle Benutzer verfügbar sein sollen, muss **R** mit administrativen Rechten gestartet werden, damit ein Schreibzugriff auf den Programmordner möglich ist. Sobald ein persönlicher Bibliotheksordner vorhanden ist, wird dieser jedoch vorgezogen, wie die Ausgabe der Funktion **.libPaths()** zeigt:

```
> .libPaths()
[1] "C:/Users/baltes/Documents/R/win-library/3.2"
[2] "C:/Program Files/R/R-3.2.5/library"
```

In dieser Lage kann man zur Installation die Funktion **install.packages()** benutzen und über das Argument **lib** den Installationsort bestimmen, z.B.:

```
> install.packages("irr", lib="C:/Program Files/R/R-3.2.5/library")
```

Eine weitere Möglichkeit, die Paketinstallation in der allgemeinen Bibliothek zu erzwingen, besteht darin, den persönlichen Bibliotheksordner vorübergehend (vor dem Start von **R**) durch Umbenennen zu deaktivieren.

Bei der Installation werden auch Abhängigkeiten von anderen Paketen berücksichtigt, wie die Konsolen-Protokollausgabe zum Beispiel **irr** zeigt:

```
--- Bitte einen CRAN Spiegel für diese Sitzung auswählen ---
Warnung in install.packages(NULL, .libPaths()[1L], dependencies = NA, type = type)
'lib = "C:/Program Files/R/R-3.2.5/library" ist nicht schreibbar
installiere auch Abhängigkeit 'lpSolve'

versuche URL 'https://cran.uni-muenster.de/bin/windows/contrib/3.2/lpSolve_5.6.13.zip'
Content type 'application/zip' length 676267 bytes (660 KB)
downloaded 660 KB

versuche URL 'https://cran.uni-muenster.de/bin/windows/contrib/3.2/irr_0.84.zip'
Content type 'application/zip' length 95558 bytes (93 KB)
downloaded 93 KB

Paket 'lpSolve' erfolgreich ausgepackt und MD5 Summen abgeglichen
Paket 'irr' erfolgreich ausgepackt und MD5 Summen abgeglichen
```

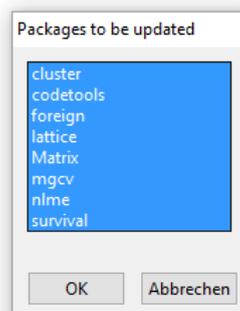
Weil **R** bei fehlenden Schreibrechten im Programmordner geschickt auf den persönlichen Bibliotheksordner ausweicht, können Benutzer ohne administrative Rechte benötigte Pakete problemlos nachrüsten (z.B. auf einem Pool-PC). Das gilt natürlich auch für Pakete, die von SPSS aus genutzt werden sollen, wobei die Installation aber über das RGui geschehen muss.

5.2.3 Installierte Pakete aktualisieren

Vor einer geplanten Aktualisierung der installierten Pakete muss **R** mit administrativen Rechten gestartet werden, damit ein Schreibzugriff auf das Programmverzeichnis möglich ist. Um die Aktualität der installierten Pakete zu prüfen und nötigenfalls Updates zu installieren, bietet das RGui den Menübefehl

Pakete > Aktualisiere Pakete

Nachdem die Liste der betroffenen Pakete mit **OK** quittiert worden ist, läuft die Aktualisierung automatisch ab:



Die Aktualisierung aller Pakete lässt sich auch über die Funktion **update.packages()** anfordern:

```
> update.packages()
```

5.2.4 Task Views

Um das Installieren und Aktualisieren von **R**-Paketen zu erleichtern, wurden sogenannte **Task Views** definiert, die aus einer mehr oder weniger großen Zahl von Paketen bestehen und komplett mit

```
> install.views("name")
```

installiert sowie mit

```
> update.views("name")
```

aktualisiert werden können. Eine Beschreibung der verfügbaren Paketbündel findet sich hier:

<http://cran.r-project.org/web/views/>

Um Task Views nutzen zu können, muss man zunächst das **R**-Paket **ctv** (*CRAN Task Views*) installieren.

5.2.5 Pakete entladen

Ein geladenes Paket kann per **detach()** wieder entladen werden, z.B.:

```
> detach(package:MASS)
```

5.2.6 Pakete zitieren

Wenn eine Veröffentlichung auf **R**-Paketen basiert, sollte der Urheber aus Respekt vor seiner geistigen Leistung und zur Orientierung des Lesers genannt werden. Mit der Funktion **citation()** erfährt man, wie ein **R**-Paket zitiert werden sollte, z.B.:

```
> citation("MASS")
```

Im Beispiel stellen sich mit W.N. Venables und B. D. Ripley zwei herausragende Förderer der **R**-Entwicklung als Autoren heraus:

To cite the MASS package in publications use:

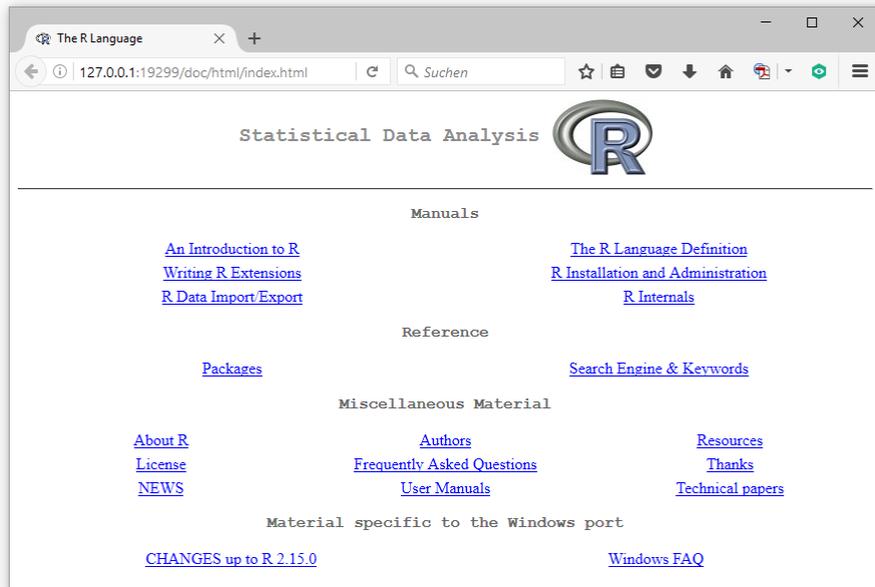
Venables, W. N. & Ripley, B. D. (2002) Modern Applied Statistics with S. Fourth Edition. Springer, New York. ISBN 0-387-95457-0

5.3 Elementare Eigenschaften der Programmiersprache R

5.3.1 Hilfe und Dokumentation

5.3.1.1 Hilfe aufrufen

Die HTML-Startseite der **R**-Hilfe



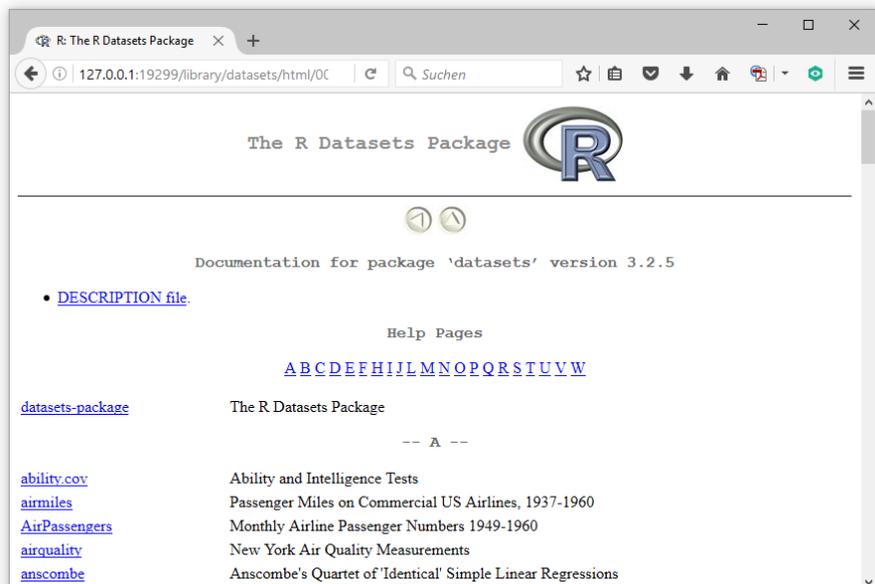
ist erreichbar über das Kommando

```
> help.start()
```

oder den RGui-Menübefehl:

Hilfe > HTML-Hilfe

Über den Link **Packages** gelangt man zu einer Liste mit allen installierten Paketen und über einen Mausklick auf ein Paket erreicht man die zugehörige Hilfeseite, z.B. beim Paket **datasets**:



Informationen zu einem *Schlüsselwort* der Programmiersprache **R** oder zu einer **R**-Funktion aus einem *geladenen* Paket erhält man über die Funktion **help()**, z.B.:

```
> help("mean")
```

Kurzform:

```
> ?"mean"
```

Bei vielen Suchbegriffen kann die Umrahmung durch Anführungszeichen entfallen, z.B.:

```
> help(mean)
> ?mean
```

Bei Schlüsselwörtern der Programmiersprache **R** ist sie jedoch erforderlich, z.B.:

```
> ?"if"
```

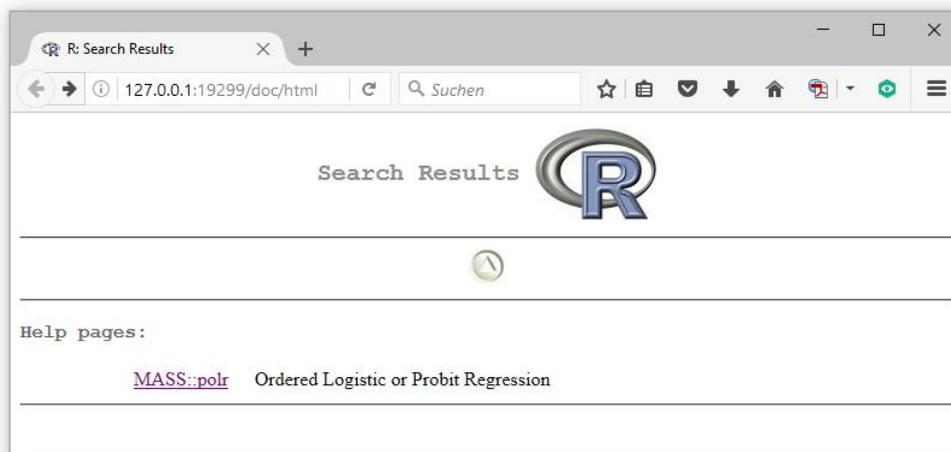
Um die Hilfedateien (auch von nicht geladenen Paketen) *fehlertolerant* nach einem Begriff zu durchsuchen, eignet sich die die Funktion **help.search()**, z.B.:

```
> help.search(probitt)
```

Kurzform (mit demselben Tippfehler):

```
> ??probitt
```

Im Beispiel stellt sich heraus, dass die Dokumentation zur Funktion **polr** im Paket **MASS** einen ähnlich geschriebenen Begriff enthält:



5.3.1.2 Beispiele aus den Hilfetexten ausführen lassen

In den Hilfetexten zu den **R**-Funktionen finden sich regelmäßig Beispiele, etwa zur **mean**-Funktion:

Examples

```
x <- c(0:10, 50)
xm <- mean(x)
c(xm, mean(x, trim = 0.10))
```

Über die Funktion **example()** kann man diese Beispiele in Aktion erleben. Durch den folgenden Aufruf

```
> example("mean")
```

erhält man das Ergebnis:

```
mean> x <- c(0:10, 50)

mean> xm <- mean(x)

mean> c(xm, mean(x, trim = 0.10))
[1] 8.75 5.50
```

5.3.1.3 Elektronische Handbücher

Das Hilfe-Menü im RGui bietet über

Hilfe > Handbücher (PDF)

zahlreiche Handbücher im PDF-Format an, z.B.:

- *An Introduction to R* (Venables et al. 2016).
- *R Language Definition* (R Development Core Team 2016)

Auf der CRAN-Webseite (*Comprehensive R Archive Network*)

<http://cran.r-project.org/>

finden sich unter der Überschrift *Documentation* diverse Handbücher zu **R**, die vom Kern-Team und aus anderen Quellen stammen.

5.3.2 Bezeichner und Kommentare

Für die Namen von Variablen und Funktionen in **R** gelten folgende Regeln:

- Erlaubte Zeichen:
 - Buchstaben (inkl. dt. Umlaute) und Ziffern
 - Punkt (.) und Unterstrich (_)
- Das erste Zeichen muss ein Buchstabe oder Punkt sein, wobei Bezeichner mit führendem Punkt für spezielle Zwecke reserviert bleiben sollten.
- Die Länge ist beliebig.
- Die **Groß-/Kleinschreibung ist relevant**.
- Schlüsselwörter der Programmiersprache **R** (z.B. **if**, **TRUE**) scheiden als Bezeichner aus. Eine Liste der reservierten Wörter erhält man mit:
`> ?Reserved`
- Die Namen müssen in ihrer Umgebung (z.B. in ihrem Paket) eindeutig sein.

In der empirischen Forschungspraxis lohnt es sich, einigen Aufwand an Zeit und Phantasie in die Bildung von Variablennamen zu stecken, die einerseits kurz und andererseits informativ sein sollten. Auf der Webseite

<https://google.github.io/styleguide/Rguide.xml>

werden die folgenden Konventionen vorgeschlagen:

- Es sollten ausschließlich Kleinbuchstaben verwendet werden. Dieses Kriterium zu verletzen und Variablennamen mit einem Großbuchstaben beginnen zu lassen, hat allerdings bei graphischen Darstellungen die angenehme Konsequenz, dass spontan perfekte Beschriftungen erscheinen (z.B. der Achsentitel *Gewicht* statt *gewicht*).
- Bei zusammengesetzten Bezeichnungen sollten die Wörter durch einen Punkt getrennt werden, z.B.: `buecher.gelesen`.
- Alternativ ist bei zusammengesetzten Bezeichnungen das sogenannte *Camel Casing* erlaubt, z.B.: `buecherGelesen`, womit an speziellen Positionen doch Großbuchstaben toleriert werden.
- Unterstriche zur Trennung von Wörtern sind unerwünscht, z.B.: `buecher_gelesen`.

Für numerische Literale (Zahlen als Bestandteile von **R**-Anweisungen) ist als **Dezimaltrennzeichen** der **Punkt** zu verwenden, z.B.:

```
> 2*3.1415926  
[1] 6.283185
```

Speziell bei gespeicherten Sequenzen von **R**-Anweisungen, die als *Skripte* bezeichnet werden (vgl. Abschnitt 5.4), sind **Kommentare** hilfreich für der Verwendung durch andere Personen und für die spätere

Nutzung durch den Autor. Mit dem Doppelkreuz wird der Zeilenrest als Kommentar gekennzeichnet, z.B.:

```
> sq <- sum(a^2) # Summe der quadrierten Elemente im Vektor a
```

5.3.3 Funktionen

R besitzt eine große Zahl von eingebauten Funktionen, die sich durch Zusatzpakete oder eigene Definitionen noch beliebig erweitern lässt. Damit kann man ...

- mathematische und insbesondere auch statistische Operationen anfordern
- Diagramme erstellen
- Datenverwaltungsaufgaben erledigen

R betrachtet Funktionen als modifizierbare Datenstrukturen. Eine **R**-Funktion kann andere Funktionen verändern oder auch neu erstellen. Details finden sich in der **R**-Sprachbeschreibung (R Development Core Team 2016).

5.3.3.1 Regeln für den Aufruf von Funktionen

Funktionen erhalten die zu verarbeitenden Daten oder Informationen zur gewünschten Arbeitsweise über Argumente. Diese Argumente werden als Kopien übergeben, so dass funktionsintern vorgenommene Veränderungen ohne Effekt auf Variablen bleiben, die der Aufrufer als Argumente übergeben hat.

Beim Funktionsaufruf sind **Positions- und Namensargumente** möglich:

- Wird die Argumentenreihenfolge der Funktionsdefinition eingehalten, kann man die Werte ohne Namen hintereinander durch Kommata getrennt schreiben.
- Paare aus einem Argumentnamen, einem Gleichheitszeichen und dem zugehörigen Wert dürfen in beliebiger Reihenfolge auftreten. Man darf die Argumentnamen abkürzen, solange Eindeutigkeit besteht, sollte aber dabei auf die Lesbarkeit achten.
- Namentlich bedachte Argumente gelten als versorgt und werden nicht mehr berücksichtigt, wenn die unbenannten Werte des Funktionsaufrufs in der Reihenfolge ihres Auftretens den Argumenten aus der Funktionsdefinition zugeordnet werden.

Besitzt ein Argument eine Voreinstellung, muss beim Aufruf kein Wert angegeben werden.

Eine Besonderheit von **R** ist das **Dreipunktargument** (...), das für zwei Aufgaben verwendet wird:

- Serienargumente
Im Definitionskopf der bereits mehrfach erwähnten und benutzten Verknüpfungsfunktion `c()` vertritt das Dreipunktargument eine beliebige Anzahl von Elementen, aus denen die Funktion einen Vektor erstellt:

```
c(..., recursive = FALSE)
```

Es sind z.B. die folgenden Aufrufe erlaubt:

```
> c(1,2)
[1] 1 2
> c(1,2,3)
[1] 1 2 3
```

- Transferargumente
Eine Funktion kann per Dreipunktargument signalisieren, dass an dieser Stelle angegebene Namensargumente an intern genutzte Funktionen weitergeleitet werden. Während die Gestaltungsflexibilität der intern aufgerufenen Funktionen erhalten bleibt, müssen deren formale Argumente

nicht in der Hüllenfunktionsdefinition wiederholt werden. Als Beispiel betrachten wir die **R**-Funktion zum Erstellen von Histogrammen, die einen sehr einfachen Definitionskopf besitzt:

```
hist(x, ...)
```

Trotzdem stehen beim Aufruf diverse Argumente zur Verfügung, z.B.:

```
> hist(x, ylab="MX")
```

Ist ein angegebener Transferparametername in keiner aufgerufenen Funktion bekannt, ignoriert **R** diesen Parameter und gibt eine Warnung aus.

Ob eine Funktion mit den drei Punkten in ihrem Definitionskopf Serien- oder Transferargumente unterstützt, ist der Dokumentation zu entnehmen.

Besitzt eine Funktion keine Argumente, ist beim Aufruf trotzdem eine leere Parameterliste (über ein Paar runder Klammern) an den Namen anzuhängen.

5.3.3.2 Elementare Funktionen

Die anschließend vorgestellten **R**-Funktionen werden häufig benötigt und sind generell (ohne vorheriges Laden von Paketen) verfügbar.

Verkettungsfunktion **c()**

Die Verkettungsfunktion **c()** (*combine, concatenate*) erzeugt einen Vektor (vgl. Abschnitt 5.3.4.2) aus den Argumenten, die beim Aufruf durch Kommata getrennt anzugeben sind, z.B.:

```
> x <- c(1, 2, 3)
```

Im Beispiel wird das Ergebnis der Variablen **x** zugewiesen, wobei der aus den beiden Zeichen „<-“ bestehende **Zuweisungsoperator** zum Einsatz kommt, den wir vor seiner „offiziellen“ Behandlung (siehe Abschnitt 5.3.7.6) in vielen Beispielen benötigen werden.

Als Argumente der Verkettungsfunktion sind auch Vektoren erlaubt, so dass man z.B. bequem einen Vektor verlängern kann:

```
> x <- c(x, 4, 5)
```

Mit Hilfe des Zuweisungsoperators lässt sich auch ein einzelner Wert, der in **R** als einelementiger Vektor aufgefasst wird, in eine Variable schreiben.

```
> k <- 13
```

Wertausgabe mit **print()**

Mit der Funktion **print()** gibt man den Inhalt einer Variablen aus, z.B.:

```
> print(x)
[1] 1 2 3 4 5
```

Weil dies sehr oft erforderlich ist, kann man die **print()** - Funktion beim interaktiven Arbeiten implizit aufrufen, indem man einen Variablennamen per Kommandozeile abschickt, z.B.:

```
> x
[1] 1 2 3 4 5
```

Der implizite **print()** - Aufruf klappt auch bei der **R**-Nutzung über ein SPSS-Syntaxfenster. Er klappt hingegen nicht ...

- in der eingebetteten Anweisung einer **for**-Schleife (vgl. Abschnitt 5.3.8.4)
- in einem per **source()** ausgeführten Skript (vgl. Abschnitt 5.4)
- in einer selbst definierten **R**-Funktion (vgl. Abschnitt 5.6)

Am Zeilenanfang einer **print()** - Ausgabe erscheint zwischen eckigen Klammern die Indexnummer des ersten Ergebniswerts in der jeweiligen Zeile, was nur bei einer mehrzeiligen Ergebnisausgabe von Bedeutung ist. Zwecks Demonstration verwenden wir im Vorgriff den Sequenzoperator (vgl. Abschnitt 5.3.7.4):

```
> y <- 1:30
> y
[1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
[26] 26 27 28 29 30
```

Um den Wert eines Ausdrucks (z.B. einer Zuweisung) ohne expliziten **print()** - Aufruf ausgeben zu lassen, setzt man runde Klammern um den Ausdruck, z.B.:

```
> (x <- c(1, 2, 3, 4, 5))
[1] 1 2 3 4 5
```

length(x)

Die Funktion **length()** liefert die Anzahl der Elemente in einem Objekt (z.B. in einem Vektor oder in einer Liste), z.B.:

```
> length(x)
[1] 5
```

Man kann der Länge eine positive Ganzzahl zuweisen und so am Ende des Vektors Elemente ergänzen oder löschen. Auf diese Weise angehängte Elemente haben den Wert **NA** (*Not Available*), z.B.:

```
> length(x) <- 7
> x
[1] 1 2 3 4 5 NA NA
> length(x) <- 3
> x
[1] 1 2 3
```

ls()

Diese Funktion listet die Objekte im Workspace auf, z.B.:

```
> ls()
[1] "eimer" "x"
```

Im RGui unter Windows kann man den **ls()** - Aufruf auch mit dem folgenden Menübefehl auslösen:

Verschiedenes > Liste Objekte auf

rm()

Mit der Funktion **rm()** kann man ein Objekt aus dem Workspace entfernen, z.B.:

```
> x
[1] 1 2 3 4 5
> rm(x)
> x
Fehler: Objekt 'x' nicht gefunden
```

Mit dem folgenden Funktionsaufruf, der die **ls()** - Ausgabe als Wert für das Argument **list** verwendet, werden alle Objekte im Workspace gelöscht (ohne Warnung):

```
> rm(list = ls())
```

Im RGui unter Windows kann man diesen Aufruf auch mit dem folgenden Menübefehl auslösen:

Verschiedenes > Entferne alle Objekte

Durch einen Aufruf der Funktion `gc()` (*Garbage Collector*, dt.: *Müllsammler*) kann man **R** auffordern, den zuvor von gelöschten Objekten belegten Speicherplatz sofort an das Betriebssystem zurückzugeben. Ohne explizite Aufforderung erledigt **R** diese Arbeit irgendwann automatisch.

Wieviel Speicher die **R**-Objekte aktuell belegen, verrät die Funktion `memory.size()` in der Einheit MB, z.B.:

```
> memory.size()
[1] 29.73
```

`search()`

Dieser Funktion listet den Suchpfad für Objekte auf, z.B.:

```
> search()
[1] ".GlobalEnv"      "package:stats"    "package:graphics"
[4] "package:grDevices" "package:utils"    "package:datasets"
[7] "package:methods" "Autoloads"        "package:base"
```

Er startet mit der globalen Umgebung, welche die Workspace-Objekte enthält, und endet mit dem Paket `base`.

Im RGui unter Windows kann man diesen Aufruf auch mit dem folgenden Menübefehl auslösen:

Verschiedenes > Liste Suchpfad auf

`quit()` bzw. `q()`

Mit dieser Funktion wird **R** beendet, z.B.:

```
> q()
```

Das RGui unter Windows kann auch mit den üblichen Methoden des Betriebssystems beendet werden, z.B. mit dem Menübefehl:

Datei > Beenden

5.3.4 Datentypen

Während SPSS mit der aus Variablen unterschiedlichen Typs (zur Aufnahme von Zahlen, Zeichenfolgen oder Datumsangaben) bestehenden Datenmatrix als der einzigen Organisationsform für Daten auskommt, bietet **R** *mehrere* Datentypen an (Vektor, Faktor, Matrix, Array, Liste, Datentabelle). Statt von *Datentypen* spricht man auch von *Klassen*. Dementsprechend werden Variablen oft als *Objekte* bezeichnet. Zusammen mit einem Datentyp werden anschließend auch die bei seiner Verwendung häufig benötigten Funktionen vorgestellt.

5.3.4.1 Datentypbezogene Funktionen

Mit der `class()` - Funktion befragt, nennt ein Objekt seine Klasse, z.B.:

```
> mat <- matrix(c(1,2,3,4),2)
> class(mat)
[1] "matrix"
> v <- c(1,2)
> class(v)
[1] "numeric"
```

Befragte Vektor-Objekte nennen den Typ ihrer Elemente (siehe zweites Beispiel). Bei numerischen Vektoren erscheint statt **numeric** in Abhängigkeit von der Erstellungsmethode eventuell eine genauere Typangabe. Im nächsten Beispiel wird der Vektor `v` per Sequenzoperator (vgl. Abschnitt 5.3.7.4) erstellt:

```
> v <- 1:2
> class(v)
[1] "integer"
```

Bei Objekten mit typidentischen Elementen (Vektor, Faktor, Matrix, Array), meldet die Funktion **mode()** den Elementtyp, z.B.:

```
> mode(mat)
[1] "numeric"
```

Über die Funktion **is.typ()** mit booleschem Rückgabewert stellt man für ein Objekt fest, ob es den per Argument angefragten Datentyp besitzt, z.B.:

```
> is.matrix(mat)
[1] TRUE
> is.vector(mat)
[1] FALSE
```

Über die Funktion **as.typ()** lässt sich eine Typumwandlung erzwingen, z.B.:

```
> v <- c(1,2)
> mode(v)
[1] "numeric"
> v <- as.character(v)
> v
[1] "1" "2"
> mode(v)
[1] "character"
```

R verwendet eine **implizite und dynamische Typisierung** von Variablen. Der Datentyp einer Variablen wird also implizit festgelegt und kann später wieder geändert werden, z.B.:

```
> y <- 3
> mode(y)
[1] "numeric"
> y <- 'a'
> mode(y)
[1] "character"
```

5.3.4.2 Vektor

Ein Vektor ist ein Objekt, das eine geordnete Anzahl von Elementen desselben Typs enthält. Man kann den Vektor als den elementarsten und wichtigsten Datentyp in **R** betrachten, weil fast alle anderen Datentypen intern als Vektoren mit speziellen Eigenschaften realisiert sind (Ligges 2007, S. 33).

R kennt keine Skalare und behandelt z.B. eine einzelne Zahl als einen numerischen Vektor der Länge 1.

5.3.4.2.1 Vektoren erstellen

Wir haben bereits in zahlreichen Beispielen die Verkettungsfunktion **c()** dazu verwendet, einen Vektor zu erstellen:

```
> x <- c(1, 2, 5)
```

Über eine Funktion mit dem Namen des Elementtyps und einem Argument zur Längenbestimmung erhält man einen Vektor mit Nullinitialisierung der Elemente. Im folgenden Beispiel entsteht ein numerischer Vektor mit 8 Nullen:

```
> n8 <- numeric(8)
> n8
[1] 0 0 0 0 0 0 0 0
```

Mit der **character()** - Funktion erhält man einen Vektor mit leeren Zeichenfolgen als Elementen, z.B.:

```
> c3 <- character(3)
> c3
[1] "" "" ""
```

Mit dem durch einen Doppelpunkt bezeichneten **Sequenzoperator** lässt sich ein Vektor aus Zahlen produzieren, beginnend mit dem linken Operanden und dann im Einserabstand wachsend bis zur letzten Zahl, die den rechten Operanden nicht übertrifft, z.B.:

```
> (x <- 1:5)
[1] 1 2 3 4 5
```

Um eine Sequenz mit alternativer (auch negativer) Schrittweite zu erzeugen, verwendet man die **seq()** - Funktion, z.B.:

```
> (x <- seq(0.3, 0.36, by=0.01))
[1] 0.30 0.31 0.32 0.33 0.34 0.35 0.36
```

5.3.4.2.2 Vektorelemente ansprechen

Ein Element eines Vektors lässt sich per `[]` - Operator über einen (1-basierten) Indexwert ansprechen, z.B.:

```
> x <- c(1, 2, 5)
> x[2]
[1] 2
```

Ein versuchter Zugriff auf eine nicht vorhandene Indexposition liefert den Ersatzwert **NA** (*Not Available*), z.B.:

```
x[7]
[1] NA
```

In Abschnitt 5.3.6 werden weitere Optionen für den Indexzugriff behandelt.

5.3.4.2.3 Automatische Typanpassung

Werden Elemente mit *verschiedenen* Typen in einen Vektor eingefügt, findet eine Anpassung zum allgemeinsten Typ statt, z.B.:

```
> weck <- 1:3
> mode(weck)
[1] "numeric"
> weck <- c(weck, "a")
> mode(weck)
[1] "character"
> weck
[1] "1" "2" "3" "a"
```

5.3.4.2.4 Vektoren mit Elementtyp character oder logical

Als Elemente vom Typ **character** sind einzelne Zeichen und Zeichenfolgen erlaubt, wobei zur Begrenzung alternativ einfache oder doppelte Anführungszeichen verwendet werden dürfen, z.B.:

```
> s <- c("Otto's Welt", 'ist simpel')
> s
[1] "Otto's Welt" "ist simple"
```

Hinter den Werten vom Typ **logical** (**TRUE**, **FALSE**) stecken intern die Zahlen 1 und 0, z.B.:

```
> cond <- TRUE
> mode(cond)
[1] "logical"
> cond+4
[1] 5
```

Es ist erlaubt, aber *nicht* empfehlenswert, die Wahrheitswerte durch ihre Anfangsbuchstaben abzukürzen, weil diese Symbole (**T**, **F**) undefiniert worden sein könnten, z.B.:

```
> cond <- T
> mode(cond)
[1] "logical"
> T <- 4
> cond <- T
> mode(cond)
[1] "numeric"
```

5.3.4.2.5 Vektoren mit Datumsangaben

Um Vektoren mit Datumsangaben zu erzeugen, die z.B. eine Differenzberechnung zulassen, kann man so vorgehen:

- Per Verkettungsoperator einen Zeichenkettenvektor mit Werten im Format „jjjj-mm-tt“ erstellen
- und die Funktion **as.Date()** auf diesen Vektor anwenden.

Beispiel:

```
> geboren <- as.Date(c("1978-12-31", "1988-07-24"))
> geboren[2] - geboren[1]
Time difference of 3493 days
```

5.3.4.2.6 Elemente benennen

Beim Erzeugen eines Vektors über die Verkettungsfunktion kann man die Elemente benennen, z.B.:

```
> b <- c(one=1, too=2)
> b
one too
 1  2
```

Über die Funktion **names()** kann man die Namen der Elemente ermitteln,

```
> names(b)
[1] "one" "too"
```

und auch ändern, z.B.:

```
> names(b) <- c("alpha", "beta")
> b
alpha beta
  1     2
```

Die Namen lassen sich an Stelle der zugehörigen Indexwerte für den Elementzugriff verwenden, z.B.:

```
> b["alpha"]
alpha
  1
```

5.3.4.2.7 Sortieren und Ränge

Die Funktion **sort()** liefert die sortierte Variante eines Argumentvektors, z.B.:

```
> v <- c(40,55,23,11,77)
> (sv = sort(v))
[1] 11 23 40 55 77
```

Mit dem Argument **decreasing** lässt sich eine *absteigende* Sortierung veranlassen, z.B.:

```
> (sv = sort(v, decreasing=TRUE))
[1] 77 55 40 23 11
```

Die Funktion **order()** liefert einen Ergebnisvektor, der zu den Elementen des sortierten Vektors ihre Indexpositionen im Argumentvektor angibt, z.B.:

```
> (ov = order(v))
[1] 4 3 1 2 5
```

Das *i*-te Element von **ov** gibt an, welche Indexposition das Element mit dem Rang *i* (das *i*-te Element von **sv**) im Originalvektor **v** besitzt. Verwendet man das **order()** - Ergebnis als Indexvektor (vgl. Abschnitt 5.3.6.3) auf den ursprünglichen Vektor an, dann resultiert das **sort()** - Ergebnis, z.B.:

```
> v[ov]
[1] 11 23 40 55 77
```

Auch die Funktion **order()** kennt das Argument **decreasing**.

Über die Funktion **rank()** erhält man einen Ergebnisvektor mit den Rängen der Elemente des Argumentvektors, z.B.:

```
> rank(v)
[1] 3 4 2 1 5
```

Zur Veranschaulichung der drei Funktionen **sort()**, **order()** und **rank()** betrachten wir jeweils das erste Element des Ergebnisvektors (bei aufsteigender Sortierung in **sort()** und **order()**):

- **sort(v)[1]** enthält den Wert des kleinsten Elements von **v**

```
> sort(v)[1]
[1] 11
```
- **order(v)[1]** enthält die Indexposition des kleinsten Elements von **v**

```
> order(v)[1]
[1] 4
```
- **rank(v)[1]** enthält den Rangplatz des ersten Elements von **v**

```
> rank(v)[1]
[1] 3
```

5.3.4.2.8 Elemente replizieren

Die **rep()** - Funktion leistet eine Replikation der Elemente von Vektoren (oder auch von Faktoren und Listen, siehe unten). Ihre wichtigsten Argumente sind:

- **x**
Objekt mit den zu wiederholenden Elementen
- **times**
Ist das **times**-Argument ein Vektor mit Länge 1, legt es eine identische Anzahl von Wiederholungen für alle **x**-Elemente fest. Ist das **times**-Argument ein Vektor mit der Länge von **x**, legt es für jedes **x**-Element einen individuellen Wiederholungsfaktor fest.

Beispiel:

```
> v <- c(1,2,3)
> rep(v,2)
[1] 1 2 3 1 2 3
> rep(v,c(1,2,3))
[1] 1 2 2 3 3 3
```

5.3.4.3 Faktor

Ein Faktor in **R** ist das Analogon zu einer kategorialen (nominalen oder ordinalen) Variablen in SPSS. Als Designvariable in einer Analysefunktion (z.B. **lm()**) wird ein Faktor automatisch korrekt behandelt und durch Kodiervariablen repräsentiert.

Für den Indexzugriff auf einzelne Elemente gelten bei Faktoren dieselben Regeln wie bei Vektoren (siehe Abschnitt 5.3.4.2.2).

Ein Faktor wird aus einem Vektor mit Modus **numeric** oder **character** erstellt und hat selbst stets den Modus **numeric**.

5.3.4.3.1 Nominale Faktoren

Einen **nominalen Faktor** erstellt man aus einem Vektor mit der Funktion **factor()**. Das folgende Beispiel verwendet einen Vektor mit Modus **character**:

```
> nf <- factor(c("A", "C", "B", "A", "A", "C", "B"))
> nf
[1] A C B A A C B
Levels: A B C
```

Die Funktion **levels()** liefert die Kategorienbeschriftungen als Vektor mit Elementen vom Typ **character**

```
> levels(nf)
[1] "A" "B" "C"
```

und erlaubt auch eine Änderung der Beschriftungen, z.B.:

```
> levels(nf) <- c("X", "Y", "Z")
```

Mit der **str()** - Funktion stellt man fest, dass die Kategorien eines Faktors intern durch natürliche Zahlen kodiert werden:

```
> str(nf)
Factor w/ 3 levels "A","B","C": 1 3 2 1 1 3 2
```

Auch aus einem *numerischen* Vektor lässt sich ein Faktor erstellen. Im folgenden Beispiel wird das optionale Argument **labels** der Funktion **factor()** zum Etikettieren der Faktorstufen verwendet (vergleichbar mit den Wertelabels in SPSS):

```
> gruppe <- factor(c(1,2,2,1,1,3,3), labels=c("KG","EG1","EG2"))
> gruppe
[1] KG EG1 EG1 KG EG2 EG2
Levels: KG EG1 EG2
```

Fehlt das Argument **labels**, dienen die Werte auch als Etiketten und werden dazu in Zeichenfolgen gewandelt.

Wenn der Faktor in einem Modell als unabhängige Variable mit Dummy- bzw. Indikatorkodierung zum Einsatz kommt, wird die *erste* Kategorie als Referenz verwendet.

Über das optionale Argument **levels** der Funktion **factor()** lassen sich folgende Effekte erzielen:

- Man kann die Faktorwerte rekodieren, also insbesondere die Referenzkategorie für die Dummy- bzw. Indikatorkodierung neu festlegen. Im Beispiel

```
> gruppe <- factor(c(1,2,2,1,1,3,3), levels=c(3,1,2), labels=c("KG","EG1","EG2"))
> str(gruppe)
Factor w/ 3 levels "KG","EG1","EG2": 2 3 3 2 2 1 1
```

findet die folgende Rekodierung statt:

```
1 → 2
2 → 3
3 → 1
```

Die Labels werden für die *rekodierten* Werte vergeben, so dass im Beispiel das Etikett *KG* zum neuen Wert 1 (hervorgegangen aus dem alten Wert 3) gehört:

```
> gruppe
[1] EG1 EG2 EG2 EG1 EG1 KG KG
```

- Man kann einzelne Werte des zu Grunde liegenden Datenvektors ausschließen (als fehlende Werte behandeln), in dem man sie im **levels**-Argument weglässt. Im folgenden Beispiel werden Fälle mit dem Wert 3 ausgeschlossen:

```
> gruppe2 <- factor(c(1,2,2,1,1,3,3), levels=c(1,2), labels=c("KG","EG"))
> str(gruppe2)
Factor w/ 2 levels "KG","EG": 1 2 2 1 1 NA NA
```

5.3.4.3.2 Ordinale Faktoren

Einen **ordinalen Faktor** erstellt man mit der Funktion **ordered()**, z.B.:

```
> of <- ordered(c("A", "C", "B", "A", "A", "C", "B"))
> of
[1] A C B A A C B
Levels: A < B < C
```

5.3.4.4 Matrix

Eine Matrix ist zur Aufnahme empirischer Daten geeignet, wenn alle Werte vom selben Datentyp sind (z.B. **numeric**). Im Vergleich zur später vorzustellenden Datentabelle (siehe Abschnitt 5.3.4.7) ist die Matrix ...

- weniger flexibel,
- aber Speicherplatz schonender.

Damit eignet sich die Matrix speziell für Anwendungen aus dem Bereich der linearen Algebra, für den zahlreiche Matrixfunktionen vorhanden sind (siehe Abschnitt 10.2).

Sollen numerische Variablen in einer Datentabelle durch eine **R**-Funktion verarbeitet werden, die nur Matrizen unterstützt, dann hilft die Funktion **as.matrix()** weiter.

5.3.4.4.1 Matrix aus einem Vektor erstellen

Um aus einem Vektor eine Matrix zu erstellen, verwendet man in der Regel die Funktion **matrix()** mit den folgenden Argumenten:

- **data**
Hier ist ein Vektor anzugeben.
- **nrow**
Anzahl der Zeilen
- **ncol**
Anzahl der Spalten
- **byrow**
Zeilen- statt Spaltendominanz beim Verteilen der Vektorelemente auf die Matrixzellen (siehe unten)

Man muss nur die Zeilen- *oder* die Spaltenzahl angeben. Um z.B. auf die Spaltenangabe zu verzichten, lässt man den Parameter **ncol** weg, z.B.:

```
> matz <- matrix(c(1,2,3,4,5,6), 2)
> matz
      [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
```

Soll nur die Spaltenzahl genannt werden, ist **ncol** als **Namensargument** (statt als Positionsargument) zu verwenden, z.B.:

```
> matz <- matrix(c(1,2,3,4,5,6), ncol=2)
> matz
      [,1] [,2]
[1,]    1    4
[2,]    2    5
[3,]    3    6
```

Wie die Beispiele zeigen, verwendet **R** per Voreinstellung ein Indizierungsschema mit **Spaltendominanz**, d.h. die verfügbaren Vektorelemente füllen zunächst die Spalte 1 von oben nach unten, dann die Spalte 2 usw. Sollen statt dessen die Zeilen nacheinander befüllt werden, ist für das Argument **byrow** der Wert **TRUE** anzugeben, z.B.:

```
> matz <- matrix(c(1,2,3,4,5,6), 2, byrow=TRUE)
> matz
      [,1] [,2] [,3]
[1,]    1    2    3
[2,]    4    5    6
```

5.3.4.4.2 Vektoren und/oder Matrizen zu einer neuen Matrix koppeln

Mit den Funktionen **rbind()** und **cbind()** lassen sich dimensional passende Matrizen und/oder Vektoren hintereinander bzw. nebeneinander koppeln, z.B.:

```
> matz2 <- matz + 3
> matz2
      [,1] [,2] [,3]
[1,]    4    5    6
[2,]    7    8    9
> fertig <- rbind(cbind(matz,matz2),c(1,1,1,1,1,1))
```

```
> fertig
      [,1] [,2] [,3] [,4] [,5] [,6]
[1,]    1    2    3    4    5    6
[2,]    4    5    6    7    8    9
[3,]    1    1    1    1    1    1
```

Werden *benannte* Vektoren gekoppelt, resultiert eine Matrix mit benannten Spalten bzw. Zeilen, z.B.:

```
> v1 <- c(1,2,3)
> v2 <- c(4,5,6)
> m <- cbind(v1,v2)
> m
      v1 v2
[1,]  1  4
[2,]  2  5
[3,]  3  6
```

Im Abschnitt 5.3.4.4.3 ist zu erfahren, wie sich Zeilen- bzw. Spaltennamen einer Matrix über die Funktionen **rownames()** bzw. **colnames()** setzen lassen.

5.3.4.4.3 Eigenschaften einer Matrix

Über die Struktur einer Matrix informiert die Funktion **str()**, z.B.:

```
> str(matz)
 num [1:2, 1:3] 1 4 2 5 3 6
```

Wir erfahren im Beispiel, dass die Matrix 2 Zeilen und 3 Spalten hat. Ihre Elemente werden spaltendominant aufgelistet.

Die Funktion **dim()** liefert die Maximalindizes zu den beiden Dimensionen, z.B.

```
> dim(matz)
[1] 2 3
```

Über die Funktionen **nrow()** bzw. **ncol()** kann man die Anzahl der Zeilen bzw. Spalten auch separat ermitteln, z.B.:

```
> nrow(matz)
[1] 2
> ncol(matz)
[1] 3
```

Über die Funktionen **rownames()** bzw. **colnames()** lassen sich die Zeilen- bzw. Spaltennamen einer Matrix setzen oder ermitteln, z.B.:

```
> rownames(matz) <- c('A', 'B', 'C')
> colnames(matz) <- c('I', 'II')
> matz
  I II
A 1  4
B 2  5
C 3  6
> rownames(matz)
[1] "A" "B" "C"
```

Wie bei einem Vektor gilt bei einer Matrix für den Datentyp der Elemente:

- Alle Elemente müssen vom selben Datentyp sein.
- Neben **numeric** sind auch alternative Datentypen möglich (z.B. **character**, **logical**).

5.3.4.4.4 Matrizen transponieren

Zum **Transponieren** einer Matrix steht die Funktion **t()** bereit, z.B.:

```
> matz <- matrix(c(1,2,3,4,5,6), 2, byrow=TRUE)
> matz
      [,1] [,2] [,3]
[1,]    1    2    3
[2,]    4    5    6
> tmatz <- t(matz)
> tmatz
      [,1] [,2]
[1,]    1    4
[2,]    2    5
[3,]    3    6
> str(tmatz)
num [1:3, 1:2] 1 2 3 4 5 6
```

5.3.4.4.5 Funktionen auf Zeilen oder Spalten anwenden

Häufig soll eine Funktion (z.B. **max()** zur Bestimmung des Maximums) auf alle Zeilen oder Spalten einer Matrix angewendet werden. In dieser Situation bietet die Funktion **apply()** eine bequeme Lösung. Ihre Argumente sind:

- **X**
Hier ist die zu analysierende Matrix anzugeben.
- **MARGIN**
Hier wählt man zwischen den Werten
 - 1 für Zeilenbezug und
 - 2 für Spaltenbezug
- **FUN**
Hier ist der Name einer Funktion, die Vektoren verarbeiten kann, ohne Parameterliste anzugeben.

Im folgenden Beispiel werden die maximalen Spaltenwerte ermittelt:

```
> daten <- matrix(c(4,5,7,4,3,7,8,9, 3,4,6,2,3,7,3,4), ncol=2)
> daten
      [,1] [,2]
[1,]    4    3
[2,]    5    4
[3,]    7    6
[4,]    4    2
[5,]    3    3
[6,]    7    7
[7,]    8    3
[8,]    9    4
> apply(daten, 2, max)
[1] 9 7
```

5.3.4.5 Array

Die Klasse Array erweitert die Klasse Matrix, indem statt zwei beliebig viele Dimensionen erlaubt sind. Wie bei der Matrix müssen auch beim Array alle Elemente denselben elementaren Datentyp besitzen (z.B. **numeric**, **logical**, **character**).

Man erstellt einen Array über die Funktion **array()** mit den folgenden Parametern:

- **data**
Vektor mit den Elementen
- **dim**
Hier ist ein Vektor anzugeben, der durch seine Länge die Anzahl der Array-Dimensionen und durch seine Elemente den Maximalindex pro Dimension angibt.

Beispiel:

```
> arri <- array(c(1,2,3,4,5,6,7,8), c(2,2,2))
> arri
, , 1
      [,1] [,2]
[1,]    1    3
[2,]    2    4

, , 2
      [,1] [,2]
[1,]    5    7
[2,]    6    8
```

5.3.4.6 Liste

Eine Liste besteht aus einer Anzahl von Objekten, die jeweils einen beliebigen **R**-Datentyp haben dürfen, insbesondere auch den Typ Liste. Weil die Liste somit eine enorme Flexibilität bietet, wird sie von vielen **R**-Funktionen zur Organisation ihrer Ausgaben verwendet. Die von der Funktion **length()** gelieferte Länge einer Liste wird durch die Anzahl ihrer Objekte auf oberster Ebene festgelegt. Um die strukturelle Flexibilität bei den Bestandteilen einer Liste zu betonen, spricht man von auch von den *Komponenten* einer Liste (Muenchen 2011, S. 83).

5.3.4.6.1 Liste erstellen

Man erstellt eine Liste über die Funktion **list()**, z.B.:

```
> lis <- list(c(1,2,3), matrix(c(1,2,3,4), 2), list("Jetzt schlägt es", 13))
> lis
[[1]]
[1] 1 2 3

[[2]]
      [,1] [,2]
[1,]    1    3
[2,]    2    4

[[3]]
[[3]][[1]]
[1] "Jetzt schlägt es"

[[3]][[2]]
[1] 13
```

Die Liste **lis** enthält die folgenden drei Komponenten

- einen Vektor
- eine Matrix
- eine Liste

und hat damit die Länge 3:

```
> length(lis)
[1] 3
```

Beim Erzeugen einer Liste kann man ihre Elemente benennen, z.B.:

```
> lis <- list(vektor=c(1,2,3), matrix=matrix(c(1,2,3,4), 2), liste=list("Jetzt schlägt es", 13))
```

Wie bei einem Vektor kann man auch bei einer Liste über die Funktion **names()** die Komponenten benennen bzw. vorhandene Labels ermitteln, z.B.:

```
> names(lis) <- c("vektor", "matrix", "liste")
```

Sind Namen vorhanden, erscheinen diese in der Ausgabe anstelle von Indexnummern, z.B.:

```
> lis
$vektor
[1] 1 2 3

$matrix
  [,1] [,2]
[1,]  1   3
[2,]  2   4

$liste
$liste[[1]]
[1] "Jetzt schlägt es"

$liste[[2]]
[1] 13
```

5.3.4.6.2 Zugriff auf Bestandteile einer Liste

Für den Zugriff auf die Komponenten einer Liste ist ein Indexausdruck zwischen eckigen Doppelklammern zu verwenden, z.B.:

```
> lis[[2]]
  [,1] [,2]
[1,]  1   3
[2,]  2   4
```

Sind Namen für die Komponenten vorhanden, können diese für den Zugriff auf die Komponenten verwendet werden, entweder über die `[[]]` - Syntax

```
> lis[["matrix"]]
  [,1] [,2]
[1,]  1   3
[2,]  2   4
```

oder per `$`-Präfix:

```
> lis$matrix
  [,1] [,2]
[1,]  1   3
[2,]  2   4
```

Man muss schon genau hinsehen, um den Wirkungsunterschied zwischen dem Komponentenauswahl-Operator `[[]]` und dem allgemeinen Indizierungs-Operator `[]` bei einer Liste zu erkennen:

```
> lis[[1]]
[1] 1 2 3
```

```
> lis[1]
$vektor
[1] 1 2 3
```

`lis[[1]]` ist die erste Komponente in der Liste (ein Vektor mit numerischen Elementen). `lis[1]` ist hingegen eine Liste, die allerdings nur noch die erste Komponente von `lis` enthält. Mit der `class()`-Funktion befragt, nennen die beiden Objekte ihren Typ:

```
> class(lis[[1]])
[1] "numeric"
> class(lis[1])
[1] "list"
```

Das Verhalten des allgemeinen Indizierungsoperators `[]` in **R** ist durchaus konsistent. Angewandt auf einen Vektor liefert er einen neuen Vektor mit Länge 1, der genau das vom Indexausdruck angesprochene Element enthält, und angewandt auf eine Liste liefert er eine neue Liste der Länge 1, welche genau die vom Indexausdruck angesprochene Komponente enthält. Verwendet man einen Indexvektor (siehe Abschnitt 5.3.6.3) in Kombination mit dem allgemeinen Indizierungs-Operator lässt sich ein Teilvektor oder eine Teilliste mit einer Länge *größer* 1 gewinnen.

Um ein bestimmtes Element aus einer Listenkomponente anzusprechen, lässt man die beiden Indexoperatoren aufeinander folgen, z.B.:

```
> lis[[1]][3]
[1] 3
```

Es ist aber auch folgende Schreibweise erlaubt:

```
> lis[[c(1,3)]]
[1] 3
```

Natürlich kann auch eine per Namen angesprochene Komponente indiziert werden, z.B.:

```
> lis$vektor[3]
[1] 3
```

5.3.4.7 Datentabelle

Ein Data Frame (in Anlehnung an Ligges 2007 übersetzt mit: *Datentabelle*) in **R** entspricht einem Datenblatt (mit einer (Fälle \times Variablen)-Datenmatrix) in SPSS. Es ist eine *Liste* von Vektoren, Faktoren und Matrizen, die *alle gleich lang* sind.

5.3.4.7.1 Datentabelle erzeugen

Man kann eine Datentabelle mit der Funktion `data.frame()` erstellen, z.B.:

```
> alter <- c(45, 32, 58)
> geschlecht <- factor(c(1, 1, 2))
> dt <- data.frame(alter, geschlecht)
```

```
> dt
  alter geschlecht
1    45          1
2    32          1
3    58          2
```

Ein Vektor mit Elementen vom Typ **character** wird bei Aufnahme in eine Datentabelle per Voreinstellung in einen Faktor umgewandelt, z.B.:

```
> augenfarbe <- c("blau","braun","grün")
> dt <- data.frame(alter, geschlecht, augenfarbe)
> str(dt)
'data.frame':  3 obs. of  3 variables:
 $ alter      : num  45 32 58
 $ geschlecht: Factor w/ 2 levels "1","2": 1 1 2
 $ augenfarbe: Factor w/ 3 levels "blau","braun",...: 1 2 3
```

Soll stattdessen ein **character**-Vektor in die Datentabelle aufgenommen werden, kann die Wandlung folgendermaßen verhindert werden:

- für einen einzelnen **character**-Vektor durch Schachteln in einen Aufruf der Funktion **I()**, z.B.:


```
> values <- c(1, 2, 3)
> names <- c("A", "B", "C")
> tcv <- data.frame(values, I(names))
```
- für den gesamten **data.frame()** - Aufruf durch den Wert **FALSE** für das Argument **stringsAsFactors**, z.B.:


```
> tcv <- data.frame(values, names, stringsAsFactors=FALSE)
```

Mit dem Typ der in eine Datentabelle aufgenommenen Vektoren legt man das **Messniveau** der repräsentierten Merkmale fest:

- Für quantitative (kontinuierliche, intervallskalierte) Merkmale verwendet man numerische Vektoren.
- Für rangskalierte Merkmale verwendet man ordinale Faktoren (vgl. Abschnitt 5.3.4.3.2).
- Für nominalskalierte Merkmale verwendet man nominale Faktoren (vgl. Abschnitt 5.3.4.3.1).
- Ein Zeichenkettenvektor (z.B. mit den Namen der Fälle) wird von **R** als nominalskaliert behandelt.

5.3.4.7.2 Eigenschaften einer Datentabelle ermitteln und ändern

Über die aktuelle Struktur einer Datentabelle informiert die Funktion **str()**:

```
> str(dt)
'data.frame':  3 obs. of  2 variables:
 $ alter      : num  45 32 58
 $ geschlecht: Factor w/ 2 levels "1","2": 1 1 2
```

Wie bei einer Matrix kann man über die Funktionen **dim()**, **nrow()** bzw. **ncol()** die Anzahl der Fälle bzw. Variablen ermitteln, z.B.:

```
> dim(dt)
[1] 3 2
> nrow(dt)
[1] 3
```

```
> ncol(dt)
[1] 2
```

Die Namen der Spalten (Variablen) können über die Funktion **names()** abgefragt und geändert werden, z.B.:

```
> names(dt)
[1] "alter"      "geschlecht"
> names(dt) <- c("age", "gender")
```

Matrixspalten werden per Voreinstellung durch den Matrixnamen mit angehängter Spaltennummer benannt, z.B.:

```
> dtm <- data.frame(mat = matrix(c(1,2,3,4,5,6),ncol=2), vec = factor(c(1,1,2)))
> names(dtm)
[1] "mat.1" "mat.2" "vec"
> names(dtm) <- c("v1", "v2", "v3")
```

Die Namen der Zeilen (Fälle) können über die Funktion **row.names()** abgefragt und geändert werden, z.B.:

```
> row.names(dt)
[1] "1" "2" "3"
> row.names(dt) <- ("alpha", "beta", "gamma")
```

Per Voreinstellung werden laufende Nummern verwendet.

5.3.4.7.3 Bestandteile einer Datentabelle ansprechen

Auf eine einzelne Variable greift man über den **\$**-Operator zu, z.B.:¹

```
> mean(dt$alter)
[1] 45
```

Befindet sich eine Matrix als Komponente in einer Datentabelle, werden ihre Spalten wie Vektoren behandelt, anzusprechen über den Matrixnamen mit angehängter Spaltennummer, z.B.:

```
> dtm <- data.frame(mat = matrix(c(1,2,3,4,5,6),ncol=2), vec = factor(c(1,1,2)))
> dtm
  mat.1 mat.2 vec
1     1     4  1
2     2     5  1
3     3     6  2
> dtm$mat.2
[1] 4 5 6
```

Zur Ansprache einzelner Werte kann bei Datentabellen (abweichend von gewöhnlichen Listen) die Notation der Matrizen verwendet werden, z.B.:

```
> dt <- data.frame(alter = c(45, 32, 58), geschlecht = factor(c(1, 1, 2)))
> dt[3,1]
[1] 58
```

Weitere Möglichkeiten zur Auswahl von Bestandteilen einer Datentabelle werden in Abschnitt 5.3.6 vorgestellt.

¹ Mit der offiziell erst in Abschnitt 9.1.1.2 behandelten Funktion **mean()** wird das arithmetische Mittel eines Vektors berechnet.

5.3.4.7.4 Variablen ergänzen oder entfernen

Um eine zusätzliche Variable in eine Datentabelle

```
> dt <- data.frame(alter=c(45, 32, 58), geschlecht=factor(c(1, 1, 2)))
```

aufzunehmen, kann man die Funktion `data.frame()` erneut aufrufen. Man verwendet die vorhandene Datentabelle als erstes Argument, ergänzt die neue Variable als zweites Argument und weist der Datentabelle die Funktionsrückgabe zu, z.B.:

```
> dt <- data.frame(dt, bildung = factor(c(3,4,2)))
```

Eine alternative Möglichkeit besteht darin, den Namen der neuen Variablen per `$`-Syntax an den Datentabellennamen anhängen und per Zuweisungsoperator die Werte folgen zu lassen, z.B.:

```
> dt$bildung <- factor(c(3,4,2))
```

```
> dt
```

```
  alter geschlecht bildung
1    45           1       3
2    32           1       4
3    58           2       2
```

Eine einfache Möglichkeit, eine Variable aus einer Datentabelle zu entfernen, besteht darin, dem Variablennamen den Wert `NULL` zuzuweisen, z.B.:

```
> dt$bildung <- NULL
```

```
> dt
```

```
  alter geschlecht
1    45           1
2    32           1
3    58           2
```

5.3.4.7.5 Datentabelle in den Suchpfad der R-Sitzung aufnehmen

Um die Angabe des Datentabellennamens einzusparen, kann man einen Data Frame per `attach()` in den Suchpfad der **R**-Sitzung aufnehmen, z.B.:

```
> attach(dt)
> mean(alter)
[1] 45
```

Ein Aufruf der Funktion `search()` zeigt, dass sich die Datentabelle nun an 2. Position im Suchpfad befindet, direkt hinter der globalen Umgebung der **R**-Sitzung (mit den Workspace-Objekten):

```
> search()
[1] ".GlobalEnv"      "dt"                "package:stats"     "package:graphics"
[5] "package:grDevices" "package:utils"     "package:datasets"  "package:methods"
[9] "Autoloads"       "package:base"
```

Bei Namensgleichheit gewinnt also das Objekt in `dt`, sofern sich der Konkurrent nicht in der globalen Umgebung befindet.

Mit der Funktion `detach()` lässt sich eine Datentabelle wieder aus dem Suchpfad entfernen, z.B.:

```
> detach(dt)
```

Es ist zu beachten, dass `attach()` von den Variablen in der Datentabelle *Kopien* erstellt und im Workspace ablegt, so dass sich Schreibzugriffe *nicht* auf die Originale auswirken, z.B.:

```
> dt <- data.frame(alter=c(45, 32, 58), geschlecht=factor(c(1, 1, 2)))
> attach(dt)
> alter
[1] 45 32 58
> alter[1] <- 99
> alter
[1] 99 32 58
> dt$alter
[1] 45 32 58
```

Die Kopien persistieren auch nach der **detach()** - Anweisung, z.B.:

```
> detach(dt)
> alter
[1] 99 32 58
```

Wegen des großen Fehlerrisikos warnen viele **R**-Kenner vor der **attach()** - Funktion (z.B. Bliese 2013, S. 11; Muenchen 2011, S. 422ff). Auf keinen Fall sollte man die Variablen einer „eingehängten“ Datentabelle ändern.

Wenn in *einem* Funktionsaufruf mehrere Variablen aus einer Datentabelle angesprochen werden müssen, kann die **with()** - Funktion eine Vereinfachung bewirken. Im folgenden Aufruf der **hist()** - Funktion zur Erstellung eines Histogramms für eine Teilstichprobe (vgl. Abschnitt 8.2.5.5) stammen die Variablen **x** und **treat** aus der Datentabelle **dtlongname**:

```
> with(dtlongname, hist(x[treat==0]))
```

Im Vergleich zur alternativen Schreibweise

```
> hist(dtlongname$x[dtlongname$treat==0])
```

ist allerdings nur ein kleiner Einspareffekt festzustellen. Der **with()** - Rahmen muss für jeden Funktionsaufruf wiederholt werden und trägt durch das zusätzliche Klammernpaar nicht zur Übersichtlichkeit bei. Daher wird die syntaktisch weit bequemere **attach()** - Funktion trotz der damit verbundenen Risiken vielfach doch verwendet.

Bei Modellierungsfunktionen (z.B. **lm()**) kann man im **data**-Argument eine Datentabelle angeben und deren Variablen in anderen Argumenten mit einfachen Namen ansprechen, z.B.:

```
> lm(formula = y ~ x, data = casedata)
```

5.3.4.7.6 Funktionen auf Variablen oder Fälle anwenden

Wie bei Matrix-Objekten (siehe Abschnitt 5.3.4.4.5) lässt sich auch bei Datentabellen per **apply()** eine Funktion auf alle Spalten oder alle Zeilen anwenden. Im folgenden Beispiel werden die maximalen Werte der Variablen in einer Datentabelle ermittelt:

```
> dt <- data.frame(alter = c(45,32,58), groesse=c(166,167,178))
> apply(dt, 2, max)
alter groesse
58      178
```

Mit dem zweiten Parameter legt man fest, ob die Funktion auf die Zeilen (Wert = 1) oder auf die Spalten (Wert = 2) wirken soll.

5.3.5 Fehlende Werte

Bei beliebigen Datentypen dient **NA** (*Not Available*) als Ersatz für fehlende Werte, z.B. bei einer Datentabelle mit zwei numerischen Vektoren

```
> dt <- data.frame(alter = c(NA, 32, 58), groesse=c(166, 167, 178))
> dt
  alter groesse
1    NA     166
2    32     167
3    58     178
```

oder bei einem Vektor mit Modus **character**

```
> vc <- c("a", NA, "b", "a")
> vc
[1] "a" NA  "b" "a"
```

Es ist zu beachten, dass **NA** auch bei Variablen mit **character**-Modus *ohne* Anführungszeichen geschrieben wird.

Über die Funktion **is.na()** lässt sich überprüfen, ob der Ersatzwert **NA** vorliegt, z.B.:

```
> is.na(dt$alter[1])
[1] TRUE
```

Bei einem mehrelementigen Objekt (z.B. Vektor, Matrix, Datentabelle) als Argument erhält man ein strukturgleiches Ergebnisobjekt, z.B.:

```
> is.na(dt)
      alter groesse
[1,]  TRUE  FALSE
[2,] FALSE  FALSE
[3,] FALSE  FALSE
```

Über die Funktion **any()** (vgl. Abschnitt 5.3.7.3) kann man feststellen, ob überhaupt Werte fehlen, z.B.:

```
> any(is.na(dt))
[1] TRUE
```

Die Anzahl fehlender Werte lässt sich mit der Funktion **sum()** ermitteln, weil die logischen Werte **TRUE** und **FALSE** intern als 1 und 0 gespeichert werden, z.B.:

```
> sum(is.na(dt))
[1] 1
```

Misslingt eine Berechnung (z.B. 0/0), resultiert der Ersatzwert **NaN** (*Not a Number*), z.B.

```
> x <- 0/0
> x
[1] NaN
```

Über die Funktion **is.nan()** lässt sich überprüfen, ob der Ersatzwert **NaN** vorliegt, z.B.:

```
> is.nan(x)
[1] TRUE
```

In der Regel benimmt sich ein **NaN** wie **NA**, und entsprechend liefert die Funktion **is.na()** auch beim Ersatzwert **NaN** das Ergebnis **TRUE**:

```
> is.na(x)
[1] TRUE
```

Oft ist es sinnvoll, sich auf die Fälle mit einem vollständigen Datensatz zu beschränken. Von der Funktion **na.omit()** erhält man die um Fälle mit fehlenden Werten (**NA** oder **NaN**) bereinigte Variante einer Datentabelle, z.B.:

```
> dt <- data.frame(alter = c(45,32,NA), groesse=c(NaN,167,178))
> dt <- na.omit(dt)
> dt
  alter groesse
2    32    167
```

Enthält ein Vektor mindestens einen fehlenden Wert (**NA** oder **NaN**) liefern statistische Funktionen wie **sum()** oder **mean()** (vgl. Abschnitt 9.1.1.2) das Ergebnis **NA** oder **NaN**, z.B.:

```
> a <- c(1, NA, 3)
> mean(a)
[1] NA
```

Soll stattdessen aus den vorhandenen Argumenten ein Ergebnis ermittelt werden, ist das Argument **na.rm** auf den Wert **TRUE** zu setzen, z.B.:

```
> mean(a, na.rm=TRUE)
[1] 2
```

5.3.6 Indexzugriff

5.3.6.1 Zugriff auf einzelne Elemente

Zugriff auf **einzelne Elemente** (Indexstart mit 1):

- Bei einem Vektor: $[i]$

Beispiel:

```
> v <- c(1,2,3)
> i <- 3
> v[i]
[1] 3
```

- Bei einer Matrix: $[i, j]$

Beispiel:

```
> m <- matrix(c(1,2,3,4,5,6),3)
> m
      [,1] [,2]
[1,]    1    4
[2,]    2    5
[3,]    3    6
> m[2,2]
[1] 5
```

- Bei einer Liste: $[[i]]$,

Bei einer Liste mit benannten Elementen: $[["name"]]$ oder $liste$name$

Beispiel:

```
> lst <- list(name="Brgl", vorname="Thea", anzKinder=4, alterKinder=c(12,10,8,2))
```

Ansprache per Index

```
> lst[[2]]
[1] "Thea"
```

Ansprache per Elementname:

```
> lst[["Vorname"]]
[1] "Thea"
> lst$alterKinder
[1] 12 10 8 2
```

5.3.6.2 Zugriff auf einen Zeilen- oder Spaltenvektor aus einer Matrix oder Datentabelle

Um die k -te Zeile oder Spalte einer Matrix auszuwählen, verwendet man den Indexausdruck $[k,]$ oder $[,k]$, z.B.

```
> m <- matrix(c(1,2,3,4,5,6), 2)
> m
      [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
> m[2,]
[1] 2 4 6
> m[,1]
[1] 1 2
```

Analog kann man einen Fall bzw. eine Variable aus einer Datentabelle wählen, z.B.:

```
> daten <- data.frame(alter=c(45,55,NA,58), geschlecht=factor(c(1,2,1,2)))
> daten[1,]
  alter geschlecht
1    45           1
```

Die bei Datentabellen erlaubte $\$$ -Notation zur Auswahl von Spalten (Variablen)

```
> daten$alter
[1] 45 55 NA 58
```

klappt bei Matrizen mit benannten Spalten *nicht*, z.B.:

```
> v1 <- c(1,2,3)
> v2 <- c(4,5,6)
> m <- cbind(v1,v2)
> m
      v1 v2
[1,]  1  4
[2,]  2  5
[3,]  3  6
> m$v1 # Fehler!
Fehler in m$v1 : $ operator is invalid for atomic vectors
```

Auf folgende Weise lassen sich vorhandene Spaltennamen bei Matrizen aber doch zur Spaltenauswahl nutzen:

```
> m[, "v1"]
[1] 1 2 3
```

5.3.6.3 Indexvektoren

Über ein Vektor-Argument für den allgemeinen Index-Operator $[]$ wählt man eine Teilmenge der Elemente ...

- in einem Vektor
- in den Zeilen oder Spalten einer Matrix
- in einer Datentabelle
- in einer Liste

Als Indexvektor sind u.a. erlaubt:

- Numerischer Indexvektor

Durch einen Vektor mit ausschließlich positiven Einträgen werden die Elemente mit den entsprechenden Indexwerten ausgewählt, z.B.:

```
> v <- c(11,7,13,54,76)
> v[c(1,3,4)]
[1] 11 13 54
```

Im nächsten Beispiel wird per Sequenzoperator ein Indexvektor für die Auswahl von Matrixspalten gebildet:

```
> m <- matrix(c(1,2,3,4,5,6), 2)
> m[,1:2]
      [,1] [,2]
[1,]    1    3
[2,]    2    4
```

Durch *negative* Einträge im Indexvektor werden die Elemente an den entsprechenden Positionen ausgeschlossen, z.B.:

```
> v[c(-1,-3,-4)]
[1] 7 76
```

- Logischer Indexvektor

Durch einen Indexvektor mit dem Modus **logical** werden alle Indexelemente mit den Wert **TRUE** ausgewählt, z.B.:

```
> t <- c(TRUE, TRUE, FALSE, FALSE, FALSE)
> v[t]
[1] 11 7
```

Im folgenden Beispiel werden mit Hilfe des Modulo-Operators (vgl. Abschnitt 5.3.7.1) die Elemente mit geradzahligem Wert gewählt:

```
> v[v %% 2 == 0]
[1] 54 76
```

Über eine Wertzuweisung an den Indexvektorausdruck lassen sich die ausgewählten Elemente ändern. Im folgenden Beispiel werden die negativen Werte eines Vektors auf 0 gesetzt:

```
> v <- c(-1,3,5,-7,2)
> v[v<0] <- 0
> v
[1] 0 3 5 0 2
```

Während der numerische Indexvektor eine Teilmenge der Indexpositionen aus dem Startobjekt enthalten kann, benötigt der logische Indexvektor dieselbe Länge wie das Startobjekt und wird nötigenfalls nach der Recycling-Regel aufgefüllt (vgl. Abschnitt 5.3.7.5).

Enthält ein logischer Indexvektor **NA**-Werte, sollte er durch die Funktion **which()** (siehe Abschnitt 5.3.7.3) in einen numerischen Indexvektor überführt werden (Wollschläger, 2010, S. 38). **which()** liefert einen Vektor mit den Indexpositionen der **TRUE**-Werte im Argumentvektor, z.B.:

```
> daten <- data.frame(alter=c(45,55,NA,58), geschlecht=factor(c(1,2,1,2)))
> (indLog <- Daten$alter > 50)
[1] FALSE TRUE NA TRUE
> (indNum <- which(indLog))
[1] 2 4
```

Über den logischen Indexvektor gelingt es im Beispiel nicht, eine Datentabelle mit den Positivfällen zu extrahieren:

```
> Daten[indLog,]
  alter geschlecht
2    55           2
NA   NA         <NA>
4    58           2
```

Mit der **which()**-Rückgabe wird dieses Ziel hingegen erreicht:

```
> Daten[indNum,]
  alter geschlecht
2    55           2
4    58           2
```

Weitere Details zu Indexvektoren erläutert das **R** Development Core Team (2016, Abschnitt 3.4.1).

5.3.6.4 Indexmatrizen

Einen Vektor mit einer Auswahl der Elemente einer Matrix gewinnt man mit einer Indexmatrix, die 2 Spalten aufweist, so dass jede Zeile die Indexpositionen eines ausgewählten Elements der Ausgangsmatrix enthält, z.B.:

```
> m <- matrix(c(1,2,3,4), 2)
> m
      [,1] [,2]
[1,]    1    3
[2,]    2    4
> im <- matrix(c(1,1,2,2), ncol=2, byrow=TRUE)
> im
      [,1] [,2]
[1,]    1    1
[2,]    2    2
> m[im]
[1] 1 4
```

5.3.7 Operatoren

In **R** arbeiten die Operatoren meist **elementweise**.

5.3.7.1 Arithmetische Operatoren

Symbole für die arithmetischen Operationen:

+	Addition
-	Subtraktion
*	Multiplikation
/	Division
^ oder **	Potenzieren
%%	Modulo (Divisionsrest)

Diese Operatoren werden auf Vektoren und Matrizen *elementweise* angewendet, wie das folgende Beispiel zeigt:

```
> m <- matrix(c(1,2,3,4), 2)
> m
      [,1] [,2]
[1,]    1    3
[2,]    2    4
```

```

> im <- solve(m)
> im
      [,1] [,2]
[1,]   -2  1.5
[2,]    1 -0.5
> m*im
      [,1] [,2]
[1,]   -2  4.5
[2,]    2 -2.0

```

Obiger Aufruf der Funktion `solve()` liefert die inverse Matrix `im` zum Argument `m` (vgl. Abschnitt 10.2). Weil der Produktoperator `*` elementweise arbeitet, ergibt das Produkt `m*im` *nicht* die Einheitsmatrix. Dazu muss mit dem `%*%` - Operator das *Matrixprodukt* berechnet werden:

```

> m %*% im
      [,1] [,2]
[1,]    1  0
[2,]    0  1

```

5.3.7.2 Vergleichsoperatoren

Bei den Symbolen für die Vergleichsoperatoren in **R** ist vor allem zu beachten, dass der Identitätsoperator (wie in den Programmiersprachen C und Java) durch zwei "=" - Zeichen ausgedrückt wird:

==	gleich
!=	verschieden
<=	kleiner oder gleich
<	kleiner
>	größer
>=	größer oder gleich

Beispiel:

```

> a <- c(1,2,3)
> b <- c(1,4,5)
> a<b
[1] FALSE TRUE TRUE

```

Beim Identitätstest für Gleitkommazahlen müssen technisch bedingte Abweichungen von der reinen Mathematik berücksichtigt werden, z.B.:

```

> 0.1 == (10.0 - 9.9)
[1] FALSE

```

Hier hilft die Funktion `all.equal()`:

```

> all.equal(0.1, (10.0 - 9.9))
[1] TRUE

```

5.3.7.3 Logische Operatoren

Die folgenden logischen Operatoren wirken elementweise:

!	nicht
&	und
	oder

Beispiel:

```
> a <- c(1, 3, 5)
> a > 1 & a < 5
[1] FALSE TRUE FALSE
```

Bei elementweisen logischen Operationen (mit einem Vektor von Wahrheitswerten als Ergebnis) kann man ...

- über die Funktion **all()** prüfen, ob *alle* Ergebniselemente gleich **TRUE** sind, z.B.:


```
> all(a > 1 & a < 5)
[1] FALSE
```
- über die Funktion **any()** prüfen, ob mindestens ein Ergebniselement gleich **TRUE** ist, z.B.:


```
> any(a > 1 & a < 5)
[1] TRUE
```

Ist man an den Indexpositionen der **TRUE**-Werte interessiert, hilft die Funktion **which()** weiter, z.B.:

```
> which(a > 1 & a < 5)
[1] 2
```

Im Unterschied zu den elementweisen logischen Operatoren berücksichtigen die folgenden Operatoren nur das jeweils *erste* Element:

&&	und
 	oder

Beispiele:

```
> t1 <- c(FALSE, FALSE, TRUE)
> t2 <- c(FALSE, TRUE, FALSE)
> t1 || t2
[1] FALSE
```

5.3.7.4 Sequenzoperator

Das Zeichen **:** steht in **R** für den Sequenzoperator, der einen Vektor aus Zahlen produziert, beginnend mit dem linken Operanden und dann im Einserabstand wachsend bis zur letzten Zahl, die den rechten Operanden nicht übertrifft, z.B.:

```
> 1:5
[1] 1 2 3 4 5
```

Mit der **seq()** - Funktion lässt sich eine alternative Schrittweite einstellen (auch eine negative), z.B.:

```
> seq(0.3, 0.36, by=0.01)
[1] 0.30 0.31 0.32 0.33 0.34 0.35 0.36
```

Der Sequenzoperator und die **seq()** - Funktion wurden schon im Zusammenhang mit dem Erstellen von Vektoren vorgestellt (in Abschnitt 5.3.4.2.1).

5.3.7.5 Recycling-Regel

Sind zwei **R**-Objekte elementweise zu verarbeiten, werden bei ungleicher Länge Elemente des kürzeren Objekts wiederverwendet, bis Längengleichheit besteht, z.B.:

```
> a <- 5
> b <- c(1,2,3)
> a+b
[1] 6 7 8
```

Ist die Länge des größeren Objekts *kein* Vielfaches der kleineren Länge, vermutet **R** ein Problem und warnt, z.B.:

```
> a <- c(5,6)
> b <- c(1,2,3)
> a+b
[1] 6 8 8
Warnmeldung:
In a + b : Länge des längeren Objektes
           ist kein Vielfaches der Länge des kürzeren Objektes
```

5.3.7.6 Zuweisungsoperatoren

Wie Sie mittlerweile aus zahlreichen Beispielen wissen, wird in **R** bevorzugt das Zeichenduo <- bei Wertzuweisungen verwendet. In fast allen Situationen ist das "=" - Zeichen äquivalent.¹

Beim Zweizeichen-Zuweisungsoperator kann es durch ein versehentlich eingefügtes Leerzeichen zu einem Fehler kommen, z.B.:

```
> oh
[1] 4
> oh < - 3
[1] FALSE
```

Es ist erlaubt, aber nicht üblich, den rechtsorientierten Zuweisungsoperator -> zu verwenden, z.B.:

```
> 13 -> k
```

5.3.8 Anweisungen

Eine **R**-Anweisung kann durch ein Semikolon oder eine Zeilentrennung abgeschlossen werden. Über das Semikolon ist es möglich, mehrere Anweisungen in einer Zeile unterzubringen, z.B.:

```
> g<-1:5; mean(g)
[1] 3
```

Bei einer unfertigen Anweisung erwartet **R** nach einem Zeilenwechsel die Fortsetzung der Anweisung. In dieser Situation präsentiert die graphische **R**-Bedienoberfläche (RGui) einen Fortsetzungs-Prompt, z.B.:

```
> g <-
+
```

¹ Exotische Ausnahme: Wer bei einem Funktionsaufruf in einem Aktualparameterausdruck eine Variable definieren möchte, muss den offiziellen Zuweisungsoperator verwenden, weil an dieser Stelle das "=" - Zeichen für Namensparameter reserviert ist, z.B.:

```
> mean(g=1:5)
Fehler .....
> mean(g<-1:5)
[1] 3
> g
[1] 1 2 3 4 5
```

5.3.8.1 Blockanweisung

Speziell bei bedingten Anweisungen oder Wiederholungsanweisungen (siehe unten) ist es oft nützlich, aus mehreren Einzelanweisungen eine Blockanweisung zu erstellen. Der gesamte Block ist durch ein Paar geschweifeter Klammern zu begrenzen, z.B.:

```
> if(a>0) {b<-log(a); c <- b+a}
```

5.3.8.2 if - Anweisungen

Durch die **if**-Anweisung kann man die Ausführung einer Anweisung von einer Bedingung abhängig machen. In der folgenden Syntaxbeschreibung sind die kursiv gesetzten Platzhalter durch zulässige Konkretisierungen zu ersetzen:

```
if (Logischer Ausdruck)  
  Anweisung
```

Die Anweisung wird nur dann ausgeführt, wenn der logische Ausdruck den Wert **TRUE** besitzt.

Beispiel:

```
> if(a >= 1) b <- log(a)
```

5.3.8.3 if-else - Anweisung

Soll auch dann etwas passieren, wenn der steuernde logische Ausdruck den Wert **FALSE** besitzt, erweitert man die **if**-Anweisung um eine **else**-Klausel:

```
if (Logischer Ausdruck) {  
  Anweisungsblock 1  
} else {  
  Anweisungsblock 2  
}
```

Beispiel:

```
> if(a >= 1) {  
+   b <- log(a)  
+ } else {  
+   b <- 0  
+ }
```

Mit dem vorgeschlagenen Aufbau der **if-else** - Anweisung wird verhindert, dass **R** ein **if** mit folgender Anweisung für abgeschlossen hält und sich über ein unerwartetes **else** beschwert. Wenn die gesamte Anweisung in eine einzige Zeile geschrieben wird, und die Blöcke jeweils nur eine einzige Anweisung enthalten, dann sind die geschweiften Klammern überflüssig:

```
> if(a >= 1) b <- log(a) else b <- 0
```

5.3.8.4 Wiederholungsanweisungen

Durch eine **for**-Schleife sorgt man dafür, dass eine Anweisung wiederholt ausgeführt wird, wobei in der Anweisung eine Schleifenvariable *Var* auftritt, die beim *i*-ten Schleifendurchgang das *i*-te Element eines Objekts als Wert annimmt:

```
for (Var in Objekt)  
  Anweisung
```

Beispiel:

```
> sq <- 0; a <- c(4,7,8)
> for (i in a) sq <- sq + i^2
> sq
[1] 129
```

In **R** kann und sollte man Schleifen weitgehend vermeiden, um eine elegante und performante Programmierung zu erhalten. Auch die **for**-Schleife im letzten Beispiel lässt sich leicht ersetzen:

```
> sq <- sum(a^2)
> sq
[1] 129
```

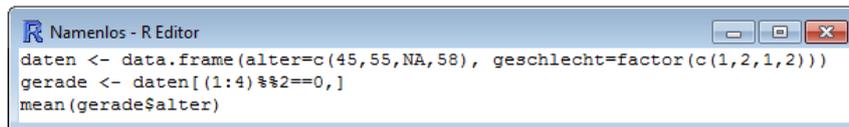
Die Funktion **sum()** addiert die Elemente eines Vektors (siehe Abschnitt 9.1.1.2).

5.4 Mit Skripten arbeiten

Statt mehrere zusammengehörige Anweisungen Zeile für Zeile in der **R**-Bedienoberfläche RGui abzuschicken, erstellt man besser ein **R-Skript**, was im RGui nach dem Menübefehl

Datei > Neues Skript

über einen integrierten Editor unterstützt wird:



```
Namenlos - R Editor
daten <- data.frame(alter=c(45,55,NA,58), geschlecht=factor(c(1,2,1,2)))
gerade <- daten[(1:4)%%2==0,]
mean(gerade$alter)
```

Um das Skript *komplett* ausführen zu lassen, kann man den Menübefehl

Bearbeiten > Alles ausführen

verwenden oder ...

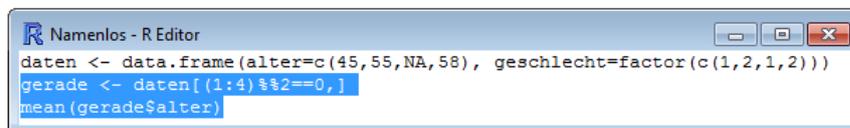
- den Text vollständig markieren (z.B. über die Tastenkombination **Strg+A**)
- und dann den Schalter  betätigen oder die Tastenkombination **Strg+R** verwenden.

Im Beispiel erscheint in der R-Konsole das Ergebnis:

```
[1] 56.5
```

Um ein Skript *partiell* ausführen zu lassen, ...

- markiert man die gewünschten Kommandos wie im folgenden Beispiel



```
Namenlos - R Editor
daten <- data.frame(alter=c(45,55,NA,58), geschlecht=factor(c(1,2,1,2)))
gerade <- daten[(1:4)%%2==0,]
mean(gerade$alter)
```

- und veranlasst die Ausführung mit dem Schalter  oder der Tastenkombination **Strg+R**.

Um die Zeile, in der sich gerade die Einfügemarke befindet, mit dem Schalter  oder der Tastenkombination **Strg+R** auszuführen, ist keine Markierung erforderlich.

Mit dem Skripteditor entstehen wiederverwendbare Anweisungsfolgen, die über den Menübefehl

Datei > Speichern unter

in eine Datei mit der Namenserweiterung **.R** gespeichert werden können.

Um ein vorhandenes Skript zu nutzen, kann man es über den Menübefehl

Datei > Öffne Skript

öffnen und dann wie oben beschrieben ausführen lassen.

Manchmal ist es sinnvoller, ein Skript über die `source()` - Funktion auszuführen, wobei im Skript definierte Daten- und Funktionsobjekte angelegt werden. Wenn sich die Skript-Datei nicht im Arbeitsverzeichnis befindet, ist der komplette Pfadname anzugeben, wobei die Pfadb Bestandteile unter Windows durch einen einfachen Vorwärtsschrägstrich oder einen doppelten Rückwärts-Schrägstrich zu trennen sind, z.B.:

```
> source("u:/eigene dateien/r/ds.r")
```

Wenn in **ds.R** das obige Skript gespeichert ist, bleibt der `source()` -Aufruf allerdings ohne Ausgabe, weil in dieser Situation implizite `print()` - Aufrufe nicht klappen. Wer eine Ausgabe sehen möchte, muss entweder im Skript den impliziten `print()` - Aufruf durch einen expliziten ersetzen (vgl. Abschnitt 5.3.3.2), oder im `source()` - Aufruf ein Echo anfordern, z.B.:

```
> source("u:/eigene dateien/r/ds.R", echo=TRUE)
```

Trotz der Empfehlung, das direkte Abschicken von Anweisungen an der Eingabeaufforderung der **R**-Konsole eher zu vermeiden, werden auch im weiteren Verlauf des Manuskripts der besseren Unterscheidbarkeit halber die Anweisungen im Stil von Direkteingaben präsentiert (mit Prompt und in roter Farbe).

5.5 Generische Funktionen und Ausgabenverwaltung

Viele **R**-Funktionen können mit Daten unterschiedlichen Typs arbeiten und dabei das jeweils passende Verhalten zeigen. Man spricht hier von *generischen Funktionen*. Ein Beispiel ist die Funktion `summary()`, der man das Ausgabeobjekt einer Statistikprozedur übergibt, um über den meist spärlichen Standardausgabeumfang hinaus weitere Details zu erfahren. Im folgenden Beispiel wird für 3 künstliche Fälle mit der Funktion `lm()` eine Regression der Variablen y auf die Variable x angefordert, wobei die Ergebnisse im Objekt `lmod` landen. Die mit einem impliziten `print()` - Aufruf angeforderte Ausgabe beschränkt sich auf die Modellformel und die Regressionskoeffizienten (ohne Signifikanztests):

```
> x <- c(1,2,3)
> y <- x + rnorm(3,0,1)
> lmod <- lm(y ~ x)
> lmod
```

```
Call:
lm(formula = y ~ x)
```

```
Coefficients:
(Intercept)          x
    1.2394         0.6929
```

Im Beispiel stecken die Ergebnisse in einem Objekt mit der Klasse **lm**,

```
> class(lmod)
[1] "lm"
```

und die `print()` - Funktion zeigt eine für **R** typische Zurückhaltung bei der Ergebnisausgabe. Um eine ausführlichere Ausgabe zu erhalten, wendet man die `summary()` - Funktion auf das **lm**-Objekt an:

```
> summary(lmod)

Call:
lm(formula = y ~ x)

Residuals:
    1      2      3 
-0.01326  0.02651 -0.01326

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.23940    0.04960   24.99  0.0255 *
x            0.69291    0.02296   30.18  0.0211 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.03247 on 1 degrees of freedom
Multiple R-squared:  0.9989,    Adjusted R-squared:  0.9978
F-statistic: 910.8 on 1 and 1 DF,  p-value: 0.02109
```

Auf den Vektor `x` angewendet, zeigt die Funktion `summary()` ein anderes Verhalten:

```
> summary(x)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  1.0     1.5     2.0     2.0     2.5     3.0
```

Im der unterschiedlichen Arbeitsweise in Abhängigkeit vom Typ des Arguments zeigt sich die Generizität der `summary()` - Funktion. Im Grunde dient die `summary()` - Funktion nur dazu, aufgrund des Argumenttyps die im Hintergrund tatsächlich auszuführende Funktion zu wählen.

5.6 Eigene Funktionen

Als einfaches Beispiel für die Entwicklung von statistischen Algorithmen mit **R** soll eine Funktion namens `pCor()` erstellt werden, die für zwei numerische Vektoren `x` und `y` die folgendermaßen definierte Pearson-Korrelation berechnet:

$$r(x, y) := \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}$$

Bei der Definition kommt man dank der nützlichen Operatoren und Funktionen in **R** mit drei Anweisungen aus (siehe Abschnitt 9.1.1.2 zu den Funktionen `mean()` und `sum()`):

```
pCor <- function(x, y) {
  xz <- x - mean(x)
  yz <- y - mean(y)
  c(cor = sum(xz*yz) / (sum(xz^2) * sum(yz^2))^0.5)
}
```

Dabei wird allerdings die bequeme und leider oft unrealistische Annahme gemacht, dass alle Werte vorhanden sind.

Auf das Schlüsselwort **function** folgt zwischen runden Klammern eine Liste mit den formalen Argumenten. Im Rumpf einer Funktion verwendet man die formalen Argumente wie lokale Variablen, die beim Aufruf durch die übergebenen Argumente initialisiert worden sind. **R** verwendet generell Wertargumente, d.h. funktionsintern vorgenommene Änderungen haben keinen Effekt auf die aufrufende Umgebung.

In der letzten Anweisung einer Funktion legt man ihre Rückgabe fest. Im Beispiel wird (in Anlehnung an Muenchen 2011, S. 107) mit der `c()` - Funktion ein Vektor mit einem benannten Element erstellt (vgl. Abschnitt 5.3.4.2.6), um für eine informative Ausgabe der Funktionsrückgabe zu sorgen (siehe unten).

Zur Funktionsdefinition kann man das in Abschnitt 5.4 beschriebene Skriptfenster verwenden, z.B.:

```
# Funktion zur Berechnung der Pearson-Korrelation
# Zum Testen zwei Vektoren mit r(x,y) = 0,8894574
x <- c(1, 2, 3, 4, 5, 6, 7, 6, 5, 4, 3, 2, 1)
y <- c(2, 2, 3, 4, 4, 5, 6, 4, 3, 3, 2, 1, 2)

pCor <- function(x, y) {
  xz <- x - mean(x)
  yz <- y - mean(y)
  c(cor = sum(xz*yz) / (sum(xz^2) * sum(yz^2))^0.5)
}
```

Nachdem das Skript ausgeführt worden ist (z.B. über den Menübefehl **Bearbeiten > Alles ausführen**), hat **R** die Funktionsdefinition zur Kenntnis genommen, und es ist ein Objekt von Typ **function** vorhanden:

```
> class(pCor)
[1] "function"
```

Zum Testen verwenden wir die folgenden Vektoren mit der Korrelation $r = 0,8894574$:

```
x <- c(1, 2, 3, 4, 5, 6, 7, 6, 5, 4, 3, 2, 1)
y <- c(2, 2, 3, 4, 4, 5, 6, 4, 3, 3, 2, 1, 2)
```

Die Funktion arbeitet offenbar korrekt:

```
> pCor(x, y)
      cor
0.8894574
```

Wenn eine Funktion ein Objekt ausgeben soll, muss die **print()** - Funktion explizit ausgerufen werden, weil implizite **print()** - Aufrufe in Funktionen nicht klappen (vgl. Abschnitt 5.3.3.2).

Abgespeicherte Funktionen lassen sich nach dem Öffnen der Skriptdatei wiederverwenden (vgl. Abschnitt 5.4).

Manchmal ist es sinnvoller, ein Skript über die **source()** - Funktion einzulesen, wobei im Skript definierte Funktionsobjekte angelegt werden. Wenn sich die Skript-Datei nicht im Arbeitsverzeichnis befindet, ist der komplette Pfadname anzugeben, wobei die Pfadb Bestandteile unter Windows durch einen einfachen Vorwärtsschrägstrich oder einen doppelten Rückwärtsschrägstrich zu trennen sind, z.B.:

```
> source("u:/eigene dateien/r/pcor.r")
```

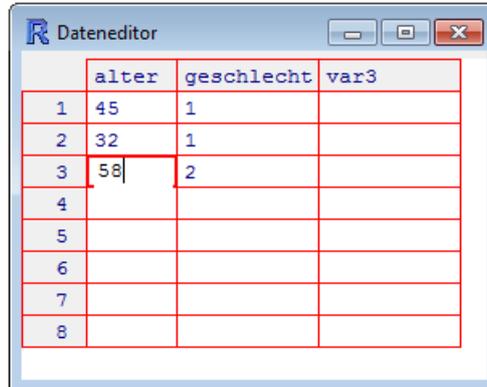
Selbstverständlich lässt sich eine selbst definierte und in einer Datei gespeicherte **R**-Funktion auch von SPSS aus nutzen. Im folgenden Beispiel wird für zwei Variablen der SPSS-Arbeitsdatei die Pearson-Korrelation mit der Funktion `pCor()` berechnet:

```
BEGIN PROGRAM R.
  casedata <- spssdata.GetDataFromSPSS()
  source("u:/eigene dateien/r/pCor.r")
  pCor(casedata$x, casedata$y)
END PROGRAM.
```

6 Bedienungserleichterungen für R

6.1 Dateneditor

Die **R**-Benutzeroberfläche hält einen einfachen Dateneditor bereit, mit dem sich Datentabellen und Matrizen anzeigen und ändern lassen:



	alter	geschlecht	var3
1	45	1	
2	32	1	
3	58	2	
4			
5			
6			
7			
8			

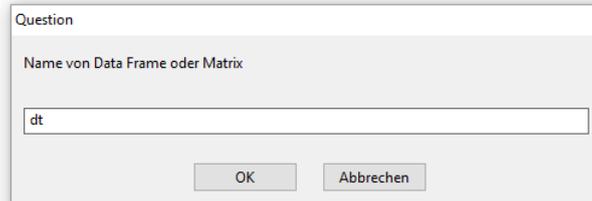
Um die Bearbeitung eines vorhandenen Datenobjekts zu starten, verwendet man entweder die Funktion **fix()**,

```
> fix(dt)
```

oder man wählt bei aktiver Konsole den Menübefehl

Bearbeiten > Dateneditor

und nennt anschließend das zu bearbeitende Datenobjekt:



Soll das Bearbeitungsergebnis in einem *anderen* Datenobjekt landen, wählt man die Funktion **edit()** mit Angabe des Rückgabeziels:

```
> dt2 <- edit(dt)
```

Der **fix()** - Funktion entspricht ein Aufruf der **edit()** - Funktion mit einem Ausgabeziel, das mit dem Argument übereinstimmt, z.B.:

```
> dt <- edit(dt)
```

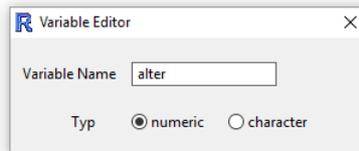
Achtung: Nach dem folgenden Einstieg

```
> edit(dt)
```

wird das Bearbeitungsergebnis nach Verlassen des Editors *nicht* in **dt** gespeichert. Um die Änderungen zu retten, bleibt dann noch der Zugriff auf das zuletzt abgeschickte Objekt, das in **R** über den Namen **.Last.value** angesprochen werden kann (Muenchen 2011, S. 117):

```
> dt <- .Last.value
```

Mit dem Editor kann man nicht nur Daten anzeigen und editieren, sondern auch Variablen deklarieren. Den folgenden Dialog mit dem Namen und dem Datentyp zu einer Variablen erreicht man per Mausklick auf ihre Spaltenbeschriftung:



Hinweise zum Editieren:

- Aktivieren Sie nötigenfalls die Zelle zur ersten Variablen des ersten Falles, und tippen Sie den zugehörigen Wert ein.
- Drücken Sie die Tabulatortaste  oder die Taste mit dem Rechtspfeil , 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).
- Verwenden Sie den Punkt als Dezimaltrennzeichen.
- Wenn ein Wert fehlt, lassen Sie die betroffene Zelle einfach leer. Dort erscheint später die Anzeige **NA** (Not Available, vgl. Abschnitt 5.3.5).
- Um einen fehlerhaften Wert zu ersetzen, tragen Sie nach dem Markieren der Zelle den korrekten Wert ein.
- Um eine Eintragung zu verändern, starten Sie nach einem Doppelklick auf die betroffene Zelle das Editieren.
- Es ist *nicht* möglich, eine Änderung rückgängig zu machen.
- Eine Zeile (einen Fall) zu löschen, ist mir per Dateneditor nicht gelungen. Per Syntax gelingt es z.B. folgendermaßen, die Zeile 7 aus der Datentabelle `dt` zu löschen (vgl. Abschnitt 5.3.6):

```
> dt <- dt[-7,]
```

Zur Anpassung der Spaltenbreiten können Sie in der Kopfzeile mit den Spaltenbeschriftungen die rechten Spaltengrenzen per Maus packen und verschieben.¹

Um den Editor zu beenden, können Sie ...

- den Menübefehl **Datei > Schließe** benutzen
- oder auf das Schließkreuz in der Titelzeile des Fensters klicken.

Es gibt kein Bedienelement zum Speichern. Wenn Sie das Fenster des Dateneditors schließen, wird das Ergebnis Ihrer Arbeit in das beim Editorstart vereinbarte Zielobjekt (im flüchtigen Hauptspeicher!) übertragen. Damit ist der Dateneditor für umfangreiche manuelle Dateneingaben bzw. -modifikationen wenig geeignet, denn zum Zwischenspeichern an einen sicheren Ort muss man ...

- den Dateneditor schließen,
- das Datenobjekt auf einen nichtflüchtigen Speicher (z.B. eine Festplatte) sichern, z.B. mit der `save()` - Funktion (siehe Abschnitt 5.1.2),
- den Dateneditor erneut öffnen.

Als Alternativen bieten sich an:

¹ Sobald der Editor zum horizontalen Rollen gezwungen ist und mindestens eine Variable nach links aus dem sichtbaren Fensterbereich verschwunden ist, wird eine geänderte Breite auf eine falsche Spalte bezogen, sodass die Funktion nicht mehr nutzbar ist (beobachtet u.a. in **R** 3.2.5).

- bei sehr großen Datenmengen ein Datenerfassungsspezialist wie Data Collection Data Entry aus der SPSS-Familie
- der SPSS-Dateneditor
- ein Tabellenkalkulationsprogramm wie z.B. Excel aus dem Office-Paket von Microsoft oder Calc aus dem freien Paket Open- bzw. LibreOffice.

Wer zur Datenerfassung ein Tabellenkalkulationsprogramm verwendet, sollte nominalskalierte Merkmale (z.B. Geschlecht) *numerisch* kodieren (z.B. 1 = Frau, 2 = Mann) und später in **R** die explizite Wandlung in einen Faktor vornehmen (Field 2012, S. 97).

6.2 R Commander

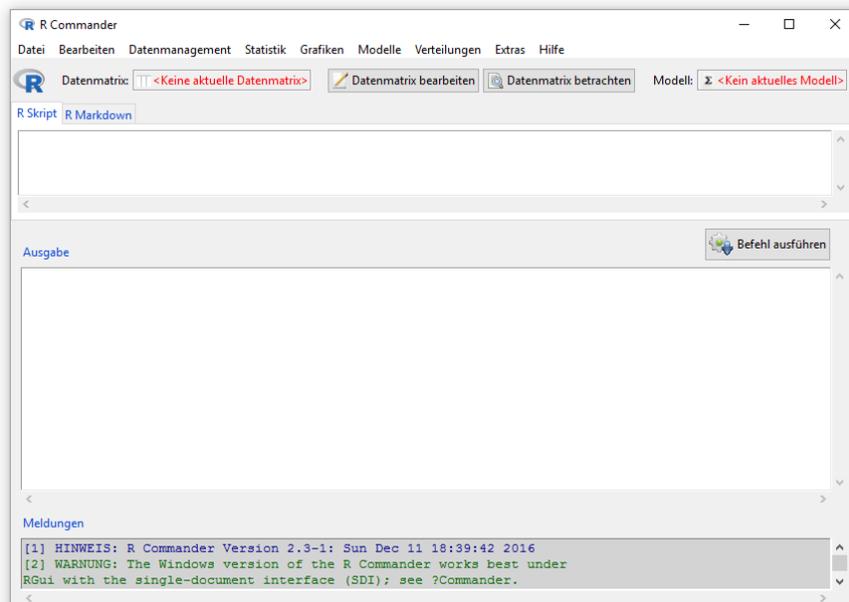
Der von John Fox (2005) entwickelte R Commander realisiert in **R** partiell die von SPSS gewohnte Bedienungsbequemlichkeit. Um Ihn nutzen zu können, muss zunächst das Paket **Rcmdr** installiert werden. Während bei **R**-Paketinstallationen die vorausgesetzten Pakete in der Regel automatisch mitinstalliert werden, ist bei **Rcmdr** aufgrund der Vielzahl von benötigten Paketen die explizite Aufforderung zur Auflösung der Abhängigkeiten durch den Wert **TRUE** für das Argument **dependencies** nach wie vor erforderlich:

```
> install.packages("Rcmdr", dependencies=TRUE)
```

Zum Starten von **Rcmdr** lädt man das Paket wie gewohnt mit der **library()** - Funktion:

```
> library(Rcmdr)
```

Es erscheint ein separates Fenster mit viel versprechenden Menüitems:



In diesem Abschnitt werden wir den R Commander zu einfachen Datenverwaltungsarbeiten verwenden. Im weiteren Verlauf des Manuskripts kommt das Programm immer wieder bei Datentransformationen und -auswertungen zum Einsatz. Weiterführende Erläuterungen zum R Commander bietet z.B. ein im Internet kostenlos verfügbares Manuskript von Fox & Bouchet-Valat (2015).

6.2.1 Datentabelle anlegen, definieren und füllen

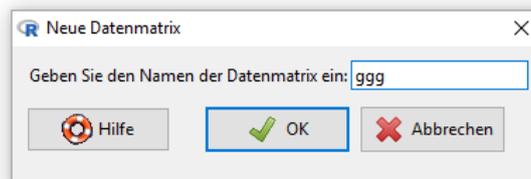
Wir werden anschließend die folgenden Daten zum Ernährungsverhalten einer studentischen Stichprobe in eine Datentabelle eintragen:

Geschlecht	Größe (in cm)	Gewicht (in kg)
2	186	82
2	178	72
2	182	75,5
1	160	65
1	168	66
1	unbekannt	76
1	165	55
2	179	76,5
1	158	50,5
2	175	80
1	176	62
2	176	unbekannt

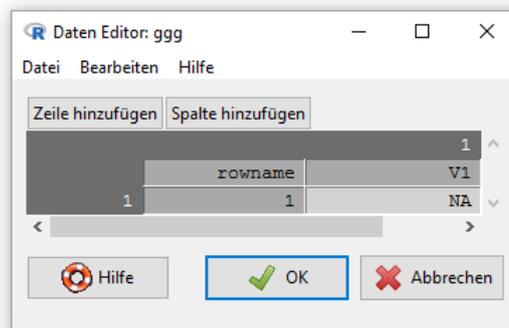
Starten Sie im Commander mit dem Menübefehl

Datenmanagement > Neue Datenmatrix

die Definition einer neuen Datentabelle, und tragen Sie im folgenden Dialog den gewünschten Namen ein, z.B.:¹

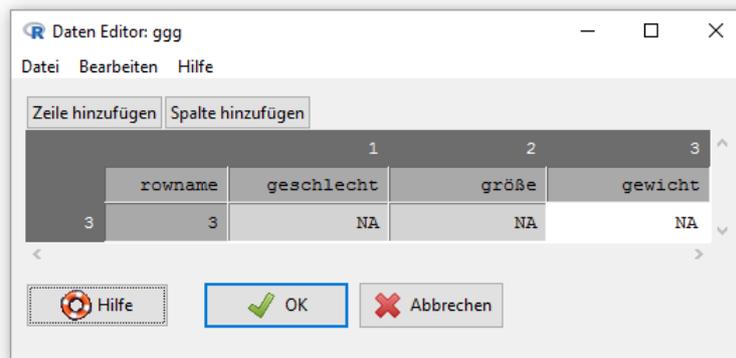


Nach dem Abschicken dieses Dialogs präsentiert der R Commander einen eigenen Dateneditor:



Legen Sie für die drei Variablen passende Namen fest, indem Sie jeweils nach einem Mausklick auf die Spaltenüberschrift den gewünschten Namen eintragen und nach den beiden ersten Variablen noch eine **Spalte hinzufügen**:

¹ Im Manuskript wird der Name **ggg** verwendet, weil die drei Merkmale (Geschlecht, Größe, Gewicht) zufälligerweise alle mit dem Buchstaben **G** beginnen:

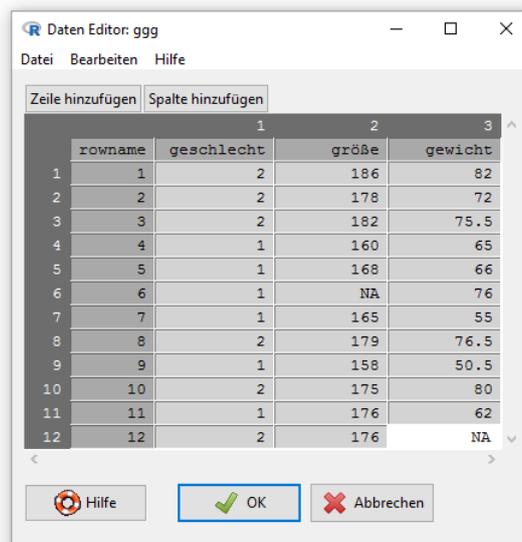


Zur Anpassung der Spaltenbreiten können Sie in der Kopfzeile mit den Spaltenbeschriftungen die rechten Spaltengrenzen per Maus packen und verschieben.

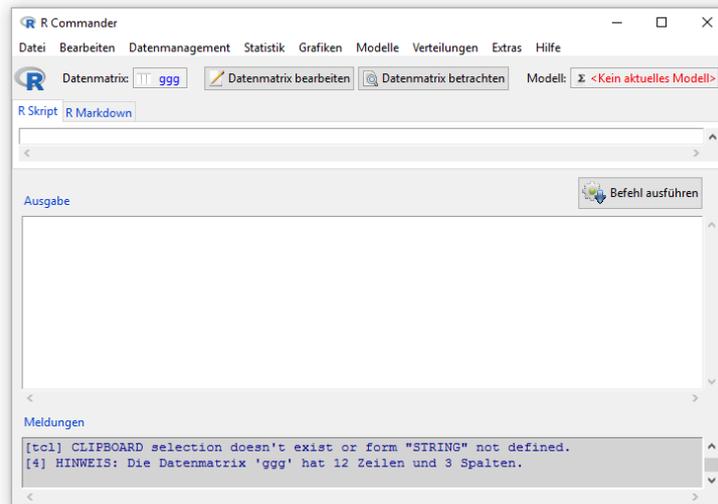
Geben Sie nun die Daten unter Beachtung der folgenden Hinweise ein:

- Drücken Sie die Rechtspfeil \rightarrow , um den eingetippten Wert zu quittieren und die Zellenmarkierung um eine Spalte nach *rechts* zu verschieben (zur nächsten Variablen).
- Klicken Sie nach der vollständigen Erfassung eines Falles auf den Schalter **Zeile hinzufügen**.
- Verwenden Sie den Punkt als Dezimaltrennzeichen.
- Wenn ein Wert fehlt, belassen Sie in der betroffenen Zelle den Initialwert **NA**.

So wollte das Ergebnis aussehen:



Wenn Sie die Dateneingabe mit **OK** beenden, ist im R Commander die **Datenmatrix ggg** eingestellt:



Man kann sie im folgenden Fenster **betrachten**

	geschlecht	größe	gewicht
1	2	186	82.0
2	2	178	72.0
3	2	182	75.5
4	1	160	65.0
5	1	168	66.0
6	1	NA	76.0
7	1	165	55.0
8	2	179	76.5
9	1	158	50.5
10	2	175	80.0
11	1	176	62.0
12	2	176	NA

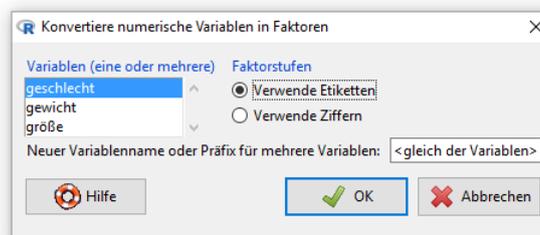
oder erneut **bearbeiten**.

6.2.2 Datenverwaltung

Als Beispiel für die Datenverwaltung mit dem R Commander wandeln wir den Vektor `geschlecht` in einen Faktor (vgl. Abschnitt 5.3.4.7.1). Nach dem Menübefehl

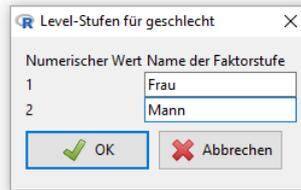
**Datenmanagement > Variablen bearbeiten >
Konvertiere numerische Variablen in Faktoren**

markieren wir im folgenden Dialog

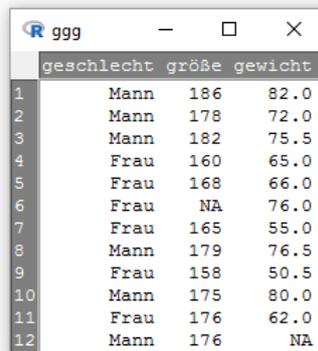


die Variable `geschlecht` und geben *keinen* **neuen Variablennamen** an, so dass keine neue Variable entsteht, sondern die vorhandene einen neuen Typ erhält.

Um Wertbeschriftungen zu ermöglichen, wählen wir die **Faktorstufen**-Option **Verwende Etiketten**, so dass nach dem Quittieren mit **OK** der folgende Dialog erscheint und die Definition von Etiketten erlaubt:



Wenn wir die Datentabelle erneut **betrachten**, ist das Ergebnis der Datentyp-Konvertierung zu besichtigen:



	geschlecht	größe	gewicht
1	Mann	186	82.0
2	Mann	178	72.0
3	Mann	182	75.5
4	Frau	160	65.0
5	Frau	168	66.0
6	Frau	NA	76.0
7	Frau	165	55.0
8	Mann	179	76.5
9	Frau	158	50.5
10	Mann	175	80.0
11	Frau	176	62.0
12	Mann	176	NA

Wir beenden den Commander und verzichten darauf, das Skript und weitere Commander-Produktionen zu speichern.

Im Konsolenfenster der RGui-Bedienoberfläche speichern wir die neue Datentabelle in eine **RData**-Datei:

```
> setwd("U:/Eigene Dateien/R")  
> save(ggg, file="ggg.RData")
```

7 Datenverwaltung und -transformation mit R

Mit der Datenverwaltung und -transformation in **R** (siehe z.B. Hain 2011, Kap. 3 und 4; Muenchen 2011, Kap. 10) werden wir uns nur knapp beschäftigen. Einfache und häufig anfallende Datentransformationen (z.B. Variablen berechnen oder rekodieren) sind bequemer in SPSS (per Syntax oder mit Dialogboxen) zu erledigen. Bei komplexen Operationen, die mit der regulären SPSS-Syntax kaum zu realisieren sind (z.B. Implementierung von Algorithmen), haben SPSS-Anwender die Wahl zwischen zahlreichen Programmieroptionen:

- Die in SPSS integrierte Programmiersprache **MATRIX** zu nutzen, hat den Vorteil, dass die erarbeiteten Lösungen von anderen SPSS-Anwendern direkt (ohne Installation von Zusatzkomponenten) übernommen werden können.
- Bei Verwendung von **R** besteht eine größere Wahrscheinlichkeit, dass man auf vorhandene Teill- oder Fertiglösungen zurückgreifen kann.
- Als weitere, im Kurs nicht behandelte Programmieroptionen unterstützt SPSS noch Python, Java und die im Windows-Bereich populären .NET-Sprachen.

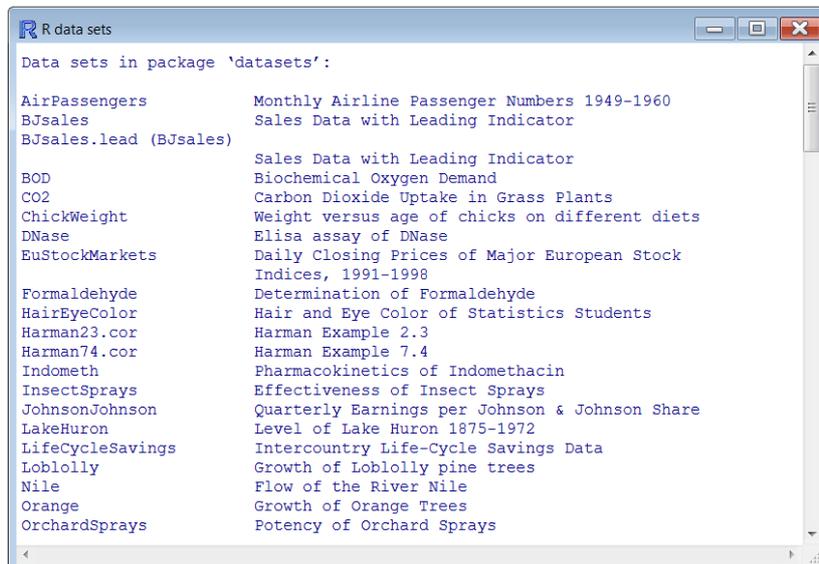
Wird ein **R**-Skript erstellt, sind dort natürlich auch einfache Datentransformationen zu erledigen. Daher werden in diesem Abschnitt nach einer Erläuterung des lesenden und schreibenden Zugriffs auf Datendateien auch elementare Datentransformationen mit **R** behandelt (z.B. Berechnung neuer Variablen, Verwendung von Zufallszahlen, Fallauswahl).

7.1 Beispieldaten in R-Paketen nutzen

R-Pakete enthalten oft illustrierende Beispieldaten, und diese sind sehr einfach zu nutzen. Über den Funktionsaufruf

```
> data()
```

erhält man eine Liste mit allen Datensätzen, die in aktuell geladenen Paketen verfügbar sind. In der Ausgabe sind die Pakete und deren Datensätze jeweils alphabetisch sortiert, z.B.:



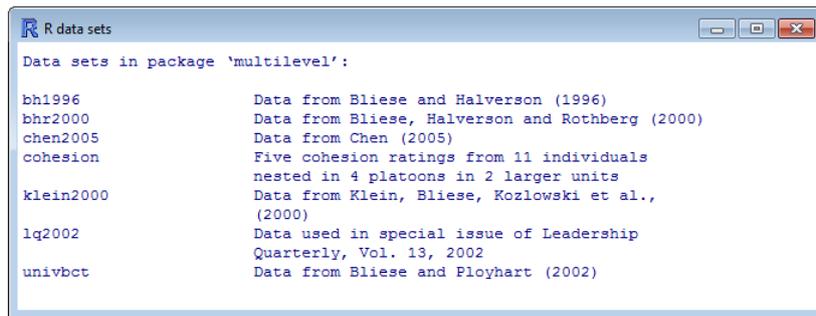
```
R data sets
Data sets in package 'datasets':

AirPassengers      Monthly Airline Passenger Numbers 1949-1960
BJsales            Sales Data with Leading Indicator
BJSales.lead (BJsales) Sales Data with Leading Indicator
BOD                Biochemical Oxygen Demand
CO2                Carbon Dioxide Uptake in Grass Plants
ChickWeight        Weight versus age of chicks on different diets
DNase              Elisa assay of DNase
EuStockMarkets     Daily Closing Prices of Major European Stock
Indices, 1991-1998
Formaldehyde        Determination of Formaldehyde
HairEyeColor        Hair and Eye Color of Statistics Students
Harman23.cor        Harman Example 2.3
Harman74.cor        Harman Example 7.4
Indometh            Pharmacokinetics of Indomethacin
InsectSprays        Effectiveness of Insect Sprays
JohnsonJohnson     Quarterly Earnings per Johnson & Johnson Share
LakeHuron           Level of Lake Huron 1875-1972
LifecycleSavings    Intercountry Life-Cycle Savings Data
Loblolly            Growth of Loblolly pine trees
Nile                Flow of the River Nile
Orange              Growth of Orange Trees
OrchardSprays       Potency of Orchard Sprays
```

Interessiert man sich für die Datensätze in einem bestimmten Paket, kann man wie im folgenden Beispiel vorgehen:

```
> library(multilevel)
> data(package="multilevel")
```

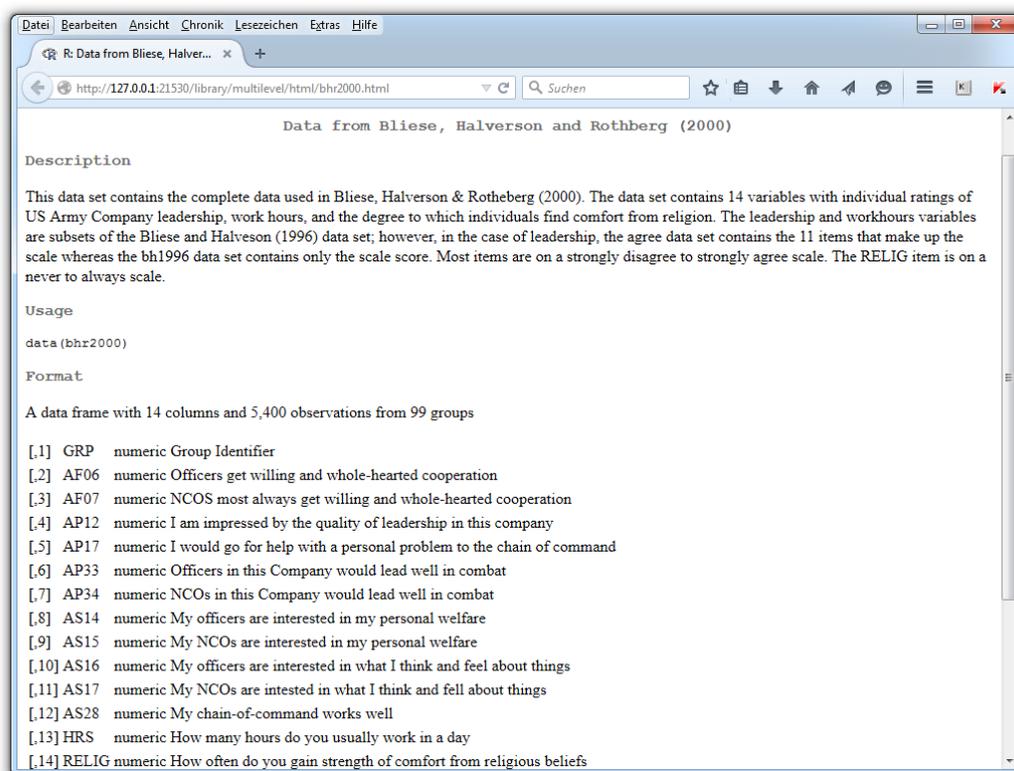
Im Beispiel resultiert die folgende Ausgabe:



An die zur Nutzung eines Datensatzes erforderlichen Informationen kommt man mit der **help()** - Funktion heran, z.B.:

```
> help(bhr2000)
```

Man erhält u.a. technische und inhaltliche Informationen über die Variablen:



Um ein Datenobjekt analysieren zu können, befördert man es mit einem **data()** - Aufruf in den Workspace, z.B.:

```
> data(bhr2000)
```

7.2 Daten in Fremdformaten lesen und schreiben

7.2.1 Textdateien mit separierten Daten

Eine Textdatei mit separierten Daten enthält pro Fall eine Zeile, in der alle Werte in einer festen Reihenfolge stehen, wobei sich zwischen zwei Werten ein Trennzeichen befindet. In Abhängigkeit vom Separatortzeichen unterscheidet man die Typen:

- Tabulator-separierte Textdatei
Typische Namensweiterungen sind **.dat** (beim Erstellen durch SPSS) und **.txt** (beim Erstellen durch Excel).
- Komma-separierte Textdatei
Trotz der Typbezeichnung verwenden SPSS und Excel beim Erstellen einer solchen Datei (zumindest in unseren Ländern) ein Semikolon als Trennzeichen, weil das Komma als Dezimaltrennzeichen benötigt wird. Als Namensweiterung wird einheitlich **.csv** verwendet.

Eine Datei mit separierten Daten enthält oft in der ersten Zeile die Variablennamen, so dass z.B. die in Abschnitt 6.2.1 vorgestellten Daten in einer Tabulator-getrennten Textdatei so aussehen:

geschlecht	größe	gewicht
2	186	82
2	178	72
2	182	75,5
1	160	65
1	168	66
1		76
1	165	55
2	179	76,5
1	158	50,5
2	175	80
1	176	62
2	176	

Fehlt bei einem Fall eine Variablenausprägung, bleibt die betroffene Zelle leer.

7.2.1.1 Lesen

Um die eben vorgestellten Daten in eine Datentabelle einzulesen, eignet sich in **R** die Funktion **read.delim2()**, die im Unterschied zur Funktion **read.delim()** das *Komma* als Dezimaltrennzeichen interpretiert. Im folgenden Beispiel wird mit dem **header**-Argument die Anwesenheit einer einleitenden Zeile mit Variablennamen bekannt gegeben:

```
> ggg <- read.delim2("U:/Eigene Dateien/R/ggg.dat", header = TRUE)
> ggg
```

Im Ergebnis

```
i..geschlecht grÃ.Ãye gewicht
1          2      186      82.0
2          2      178      72.0
3          2      182      75.5
4          1      160      65.0
5          1      168      66.0
6          1       NA      76.0
7          1      165      55.0
8          2      179      76.5
9          1      158      50.5
10         2      175      80.0
11         1      176      62.0
12         2      176       NA
```

zeigen sich zwei Macken:

- Am Dateianfang sind unerwartete Zeichen erkannt und dem ersten Variablennamen zugeschlagen worden (mit dem Ergebnis: **i..geschlecht**).
- Die Kodierung der Zeichen *ö* und *ß* im Variablennamen **größe** macht Probleme.

Die Fehler resultieren daraus, dass die Eingabedatei keine ANSI-Kodierung verwendet, sondern die modernere UTF-8 - Kodierung. Im folgenden Funktionsaufruf

```
> ggg <- read.delim2("U:/Eigene Dateien/R/ggg.dat", header = TRUE, encoding="UTF-8")
> ggg
```

wird **R** per **encoding**-Argument über die Kodierung der Eingabedatei informiert, was zumindest das Umlaute-Problem behebt:

```
  X.U.FEFF.geschlecht  gröÙe  gewicht
1                    2   186    82.0
2                    2   178    72.0
.                    .     .     .
```

Der verunstaltete Name der ersten Variablen resultiert aus der falsch verstandenen BOM-Sequenz (*Byte Order Mark*), die viele unter Windows erstellte UTF-8 - Dateien am Dateianfang enthalten. Mit der folgenden, ab **R** 3.0 unterstützten, Variante

```
> ggg <- read.delim2("U:/Eigene Dateien/R/ggg.dat", header = TRUE, fileEncoding="UTF-8-BOM")
> ggg
```

gelingt ein korrekter Datenimport:

```
  geschlecht  gröÙe  gewicht
1           2   186    82.0
2           2   178    72.0
.           .     .     .
```

Die eingelesenen Fälle erhalten automatisch fortlaufende Nummern (im Beispiel von 1 bis 12). Gelegentlich sind Datendateien einzulesen, die eine Fallidentifikation enthalten, die nicht unbedingt eine mit 1 beginnende, lückenlose Nummerierung enthält, In der folgenden Variante

```
nr      geschlecht  gröÙe  gewicht
1       2          186    82
2       2          178    72
.       .           .     .
13      2          176    .
```

der obigen Beispieldatei ist eine einleitende Fallidentifikationsvariable namens **nr** vorhanden, wobei der letzte Fall (vermutlich aus gutem Grund) den Wert 13 besitzt. In einer solchen Situation sollte man den **R**-Textdatenimport dazu überreden, die vorhandene Fallidentifikation zu verwenden und auf eine automatische Nummerierung zu verzichten:

- Wenn die Fallidentifikationsvariable am Dateianfang steht, genügt es, ihren Namen aus der Kopfzeile zu löschen.
- Um das Ziel ohne Änderung der Eingabedatei und bei beliebiger Position der Identifikationsvariablen zu erreichen, informiert man die Importfunktion mit dem Argument **row.names** darüber, in welcher Variablen sich die Fallidentifikation befindet.

Im folgenden Beispiel wird das Argument **row.names** verwendet:

```
> setwd("U:/Eigene Dateien/R")
> ggg <- read.delim2("gggNr.dat", header = TRUE, fileEncoding="UTF-8-BOM", row.names="nr")
> ggg
```

Ergebnis:

```
  geschlecht  gröÙe  gewicht
1           2   186    82.0
2           2   178    72.0
.           .     .     .
13          2   176     NA
```

Beim Lesen der CSV-Variante der Eingabedaten (inkl. Fallidentifikationsvariable **nr**)

```
nr;geschlecht;größe;gewicht
1;2;186;82
2;2;178;72
3;2;182;75,5
4;1;160;65
5;1;168;66
6;1; ;76
7;1;165;55
8;2;179;76,5
9;1;158;50,5
10;2;175;80
11;1;176;62
13;2;176;
```

ersparen wir uns den UTF-8 - Ärger, verwenden also eine Textdatei mit ANSI-Kodierung. Als **R**-Funktion verwenden wir **read.csv2()**, wobei die 2 am Ende des Funktionsnamens signalisiert, dass ein Komma als Dezimaltrennzeichen und ein Semikolon als Separatorzeichen erwartet werden. Im folgenden Kommando verwenden wir unser Wissen über die Zeilenbezeichnungsvariable und kommen gleich im ersten Versuch

```
> ggg <- read.csv2("U:/Eigene Dateien/R/gggNr.csv", header = TRUE, row.names="nr")
> ggg
```

zum gewünschten Ergebnis:

```
  geschlecht  gröÙe  gewicht
1           2    186    82.0
2           2    178    72.0
.           .      .      .
13          2    176     NA
```

Die im aktuellen Abschnitt vorgestellten Funktionen zum Lesen von Textdateien stellen Aufrufvereinfachungen für die Funktion **read.table()** dar, die jeweils im Hintergrund verwendet wird. Der folgenden Tabelle nach Hain (2011, S. 67) ist zu entnehmen, welche Einstellungen bzgl.

- der Einführungszeile mit Variablennamen (Argument **header**)
- des Separatorzeichens (Argument **sep**)
- und des Dezimaltrennzeichens (Argument **dec**)

mit den verschiedenen Funktionen verbunden sind ("**\t**" steht für das Tabulatorzeichen):

Funktion	header	sep	dec
read.table()	FALSE	""	"."
read.csv()	TRUE	","	"."
read.csv2()	TRUE	";"	";"
read.delim()	TRUE	"\t"	"."
read.delim2()	TRUE	"\t"	";"

Weil das Merkmal Geschlecht in den Beispieldateien numerisch kodiert war, ist beim Einlesen ein numerischer *Vektor* entstanden, der durch die folgende Anweisung in einen *Faktor* gewandelt werden sollte:

```
> ggg$geschlecht <- factor(ggg$geschlecht, labels=c('Frau','Mann'))
> ggg
```

Das Ergebnis:

```
  geschlecht  gröÙe  gewicht
1        Mann    186    82.0
2        Mann    178    72.0
.         .      .      .
13       Mann    176     NA
```

7.2.1.2 Schreiben

Zum Zweck der Kooperation mit anderen Programmen ist es oft erforderlich, die Variablen einer **R**-Datentabelle in eine Textdatei mit separierten Daten zu befördern, weil dieses Dateiformat von praktisch jeder Statistik-Software gelesen werden kann. Zuständig ist die **R**-Funktion `write.table()`, die als erstes Argument den Namen der Datentabelle und als zweites Argument den Namen der Zielformat erwartet. Häufig werden außerdem die Argumente mit den folgenden Namen benötigt:

- **sep**
Zwischen Anführungszeichen kann das gewünschte Trennzeichen angegeben werden, wobei insbesondere in Frage kommen:
 - `\t`
Tabulatorzeichen
 - `,` oder `;`
Um eine CSV-Datei zu erzielen, verwendet man in Abhängigkeit vom Dezimaltrennzeichen zum Separieren das Komma oder das Semikolon.
- **dec**
Per Voreinstellung verwendet **R** in der Ausgabedatei den Punkt als Dezimaltrennzeichen. Ist das Komma gewünscht, setzt man den Parameter **dec** auf den Wert `","`.
- **na**
Bei fehlenden Werten schreibt **R** per Voreinstellung **NA** in die Exportdatei. Über das Argument **na** lässt sich ein alternativer Ersatzwert vereinbaren (z.B. die leere Zeichenfolge durch `"`).
- **row.names**
Per Voreinstellung gibt **R** die Zeilenbeschriftungen aus, verzichtet aber in der einleitenden Zeile auf einen zugehörigen Variablennamen. Beim Einlesen durch Fremdprogramme kommt es daher oft zu einer falschen Zuordnung der Variablennamen oder zu einem Fehler. Um dieses Problem zu vermeiden, kann man auf die Ausgabe der Zeilenbeschriftungen verzichten und das Argument **row.names** auf den Wert **FALSE** setzen.
- **quote**
Per Voreinstellung begrenzt **R** die Variablennamen und die Werte von Faktoren durch Anführungszeichen. Dies lässt sich mit dem Wert **FALSE** für das Argument **quote** verhindern.

Um die Datentabelle `ggg`

```
> ggg
  geschlecht grösse gewicht
1      Mann   186    82.0
2      Mann   178    72.0
.
13     Mann   176     NA
```

in eine ...

- Textdatei im aktuellen Arbeitsverzeichnis
- mit Tabulator-getrennten Daten
- und Dezimalkomma
- ohne Zeilenbeschriftungen
- mit einer leeren Zeichenfolge statt fehlender Werte
- ohne Anführungszeichen um Variablennamen und Faktorwerte

zu schreiben, eignet sich das folgende Kommando:

```
> write.table(ggg, "ggg.dat", sep="\t", dec=",", row.names=FALSE, na="", quote=FALSE)
```

Das Ergebnis:

geschlecht	größe	gewicht
Mann	186	82
Mann	178	72
Mann	182	75,5
Frau	160	65
Frau	168	66
Frau		76

. . .

Unter Windows ist die Ausgabedatei ANSI-kodiert.

Bei der Ausgabe in eine ...

- Textdatei im aktuellen Arbeitsverzeichnis
- mit Semikolon-getrennten Daten und Dezimalkomma
- ohne Zeilenbeschriftungen
- mit der voreingestellten Ausgabe fehlender Werte
- ohne Anführungszeichen um Variablenamen und Faktorwerte

erspart die Funktion **write.csv2()** im Vergleich zu **write.table()** etwas Schreibarbeit, z.B.:

```
> write.csv2(ggg, "ggg.csv", row.names=FALSE, quote=FALSE)
```

Bei Faktoren schreibt **R** Platz verschwendend die Kategorietiketten (labels) in die Ausgabedatei. Von der Funktion **write.foreign()**, die zum Datenexport an andere Statistikprogramme (z.B. SPSS) dient (siehe Abschnitt 7.2.3), erhält man Textdateien mit den numerischen Werten an Stelle der Etiketten.

7.2.2 SPSS-Datendateien lesen

In einer aus SPSS gestarteten **R**-Sitzung kann man bequem auf die Variablen der SPSS-Arbeitsdatei zugreifen (siehe Abschnitt 4.1). **R** kommt aber auch im selbständigen Einsatz an SPSS-Variablen heran, sofern sich diese in einer SPSS-Datendatei befinden. Dazu muss zunächst das **R**-Paket **foreign** geladen werden:

```
> library(foreign)
```

Es gehört zu den so genannten *recommended packages*, dürfte also in praktisch jeder **R**-Installation vorhanden sein.

Eine SPSS-Datendatei (Namenserweiterung **.SAV**) kann mit der Funktion **read.spss()** gelesen werden, wobei für Dateien im Arbeitsverzeichnis (vgl. Abschnitt 5.1.1) keine Pfadangabe erforderlich ist, z.B.:

```
> ggg <- read.spss("ggg.sav", to.data.frame = TRUE, reencode="utf-8")
```

Als Rückgabe liefert die Funktion **read.spss()** per Voreinstellung eine *Liste* (vgl. Abschnitt 5.3.4.6), was man in der Regel durch den Wert **TRUE** für das Argument **to.data.frame** verhindert.

Wurde die SAV-Datei von einer SPSS-Version ab 18 erstellt, kommt es beim Öffnen zu Warnungen, z.B.:

Warnmeldung:

```
In read.spss("ggg.sav", to.data.frame = TRUE, reencode = "utf-8") :
  ggg.sav: Unrecognized record type 7, subtype 18 encountered in system file
```

Trotz der Warnungen scheint die resultierende Datentabelle intakt zu sein:

```

nr geschlecht grÖÖe gewicht
1 1 Mann 186 82.0
2 2 Mann 178 72.0
. . . . .

```

SAV-Dateien werden seit der SPSS-Version 21 bevorzugt im Unicode-Modus erstellt, wobei die UTF-8 -Kodierung zum Einsatz kommt. Um beim Lesen solcher Dateien z.B. eine falsche Interpretation von Umlauten in Variablennamen und Wertbeschriftungen

```

nr geschlecht grÄ.Äÿe gewicht
1 1 Mann 186 82.0
. . . . .

```

zu verhindern, muss in der Funktion `read.spss()` das Argument `reencode` den Wert `utf-8` erhalten. Bei SPSS-Dateien in traditioneller Kodierung wirkt der Wert `urt-8` kontraproduktiv. Die korrekte Einstellung muss man durch Ausprobieren ermitteln.

Über das Argument `use.value.labels` mit dem Voreinstellungswert `TRUE` entscheidet man darüber, ob numerische Variablen in Faktoren konvertiert werden sollen, wenn für mindestens einen Wert ein Etikett vorhanden ist. Im Beispiel passiert diese Konvertierung bei der Variablen `geschlecht`.

Im Unterschied zu den Funktionen zum Lesen von separierten Textdateien (siehe Abschnitt 7.2) besitzt `read.spss()` kein Argument, um für eine vorhandene Eingabevariable die Verwendung zur Zeilenbeschriftung zu veranlassen. Dieses Ziel lässt sich aber nach dem Import in zwei kurzen Anweisungen doch noch realisieren. Zunächst werden die Fallbeschriftungen aus der Variablen `nr` gelesen. Dann wird der Variablen `nr` der Wert `NULL` zugewiesen, um sie aus der Datentabelle zu entfernen:

```

> row.names(ggg) <- ggg$nr
> ggg$nr <- NULL

```

Das Ergebnis:

```

geschlecht grÖÖe gewicht
1 Mann 186 82.0
. . . . .

```

7.2.3 Datenexport an SPSS

Eine **R**-Datentabelle direkt in eine SPSS-Datendatei (mit Namensweiterung `.SAV`) zu schreiben, ist derzeit noch nicht möglich, doch kommt man mit der Funktion `write.foreign()` diesem Ziel recht nahe. Man erhält eine Textausdatei samt einer begleitenden SPSS-Syntaxdatei, die in SPSS zum Lesen der Textdaten genutzt werden kann. Für die aktuelle Beispieldatei kann der Export an SPSS so angefordert werden:

```

> write.foreign(ggg, datafile="Daten.dat", codefile="Syntax.sps", package="SPSS")

```

R warnt, es habe den Variablennamen `geschlecht` gekürzt, weil SPSS-Variablennamen in sehr grauer Vorzeit auf 8 Zeichen beschränkt waren. Die resultierende Textdatei hat das CSV-Format und enthält bei Faktoren platz sparend die Kategorienwerte an Stelle der von `write.table()` (vgl. Abschnitt 7.2.1.2) ausgegebenen Kategorienetiketten, was im Beispiel an der ersten Spalte mit dem Geschlecht zu sehen ist:

```

2,186,82,
2,178,72,
. . .
2,176,,

```

Die ergänzend gelieferte Syntax kann in SPSS ausgeführt werden, um die Textdaten in ein Datenblatt einzulesen:

```
DATA LIST FILE= "Data4SPSS.dat"  free (",")
/ gschlcht gröÙe gewicht  .

VARIABLE LABELS
geschlcht "geschlecht"
gröÙe "gröÙe"
gewicht "gewicht"
.

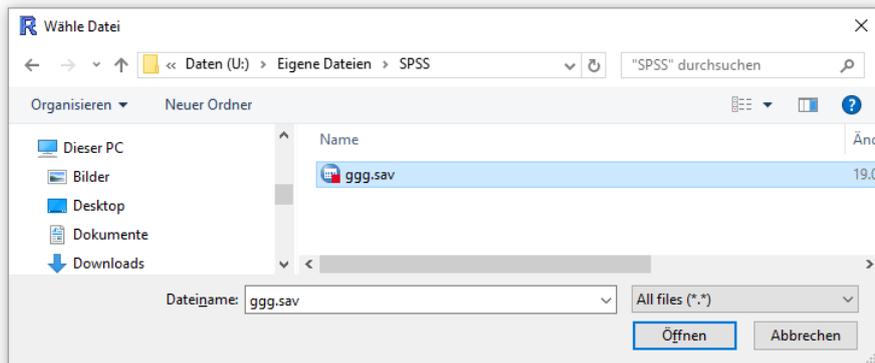
VALUE LABELS
/
geschlcht
1 "Frau"
2 "Mann"
.

EXECUTE.
```

Insgesamt sind für den Datenaustausch zwischen **R** und SPSS die in Abschnitt 4 beschriebenen Techniken auf der Basis der R-Essentials gegenüber dem Dateitransfer zu bevorzugen.

7.2.4 Dateiauswahl per Dialogbox

Um eine zu öffnende Datei per GUI-Dialog wählen zu können,



verwendet man den Funktionsaufruf **file.choose()**, der überall erlaubt ist, wo in der **R**-Syntax ein Dateiname erwartet wird, z.B.:

```
> ggg <- read.spss(file.choose(),to.data.frame = TRUE, reencode="utf-8")
```

7.3 Variablen erstellen oder modifizieren

7.3.1 Umkodieren

In der empirischen Forschungspraxis ist es oft erforderlich, die vorhandenen Werte einer Variablen auf neue abzubilden, z.B. um ein Item „umzupolen“:

5	→	1
4	→	2
3	→	3
2	→	4
1	→	5

7.3.1.1 Lösung per `recode()`

Diese Aufgabe lässt mit der Funktion `recode()` aus dem Paket `car` sehr bequem erledigen. Wir erzeugen zunächst einen numerischen Beispielvektor, und laden dann das Paket `car`:

```
> v <- c(2, NA, 4, 1, 5, 2, 4, 3)
> library(car)
```

Die Funktion `recode()` erwartet im ersten Argument den Namen der zu rekodierenden Variablen und im zweiten Argument eine Zeichenfolge mit Semikolon-separierten Ersetzungsvorschriften nach dem Muster „Alter Wert=Neuer Wert“, z.B.:

```
> recode(v, "5=1;4=2;2=4;1=5")
[1] 4 NA 2 5 1 4 2 3
```

Über flexible Möglichkeiten, eine Liste alter Werte anzugeben, die auf einen gemeinsamen Zielwert abgebildet werden sollen, informiert die Hilfe zur Funktion `recode()`.

7.3.1.2 Lösung mit logischen Indexvektoren

Mit Hilfe von logischen Indexvektoren (vgl. Abschnitt 5.3.6.3) lässt sich eine Rekodierung auch *ohne* nachzuladende Pakete erledigen. Es ist oft sinnvoll und beim gleich beschriebenen Verfahren nach einem von Wollschläger (2010, S. 39) auch von praktischem Vorteil, den Ausgangsvektor unverändert zu lassen und einen neuen Vektor mit den umkodierten Werten zu erstellen:

```
> v <- c(2, NA, 4, 1, 5, 2, 4, 3)
> w <- numeric(length(v))
> w[v==1] <- 5
> w[v==2] <- 4
> w[v==3] <- 3
> w[v==4] <- 2
> w[v==5] <- 1
> w[w==0] <- NA
> w
[1] 4 NA 2 5 1 4 2 3
```

Mit der Funktion `numeric()` (siehe Abschnitt 5.3.4.2.1) wird ein neuer numerischer Vektor passender Länge erstellt, dessen Elemente alle mit 0 initialisiert sind. Dann werden für jeden möglichen Wert des Ausgangsvektors über einen logischen Indexvektor die passenden Elemente des Zielvektors bestimmt und auf den gewünschten Wert gesetzt. `NA`-Werte des Ausgangsvektors werden bei dieser Prozedur auf die 0 abgebildet (= initialer Wert des Zielvektors). Durch die abschließende Zuweisung

```
> w[w==0] <- NA
```

erhalten diese Elemente den Wert `NA`.

7.3.1.3 Lösung mit `ifelse()`

Soll das Umkodieren zu einem Vektor mit lediglich zwei möglichen Werten führen, eignet sich die `ifelse()` - Funktion, z.B.:

```
> w <- ifelse(test=(v>3), yes=1, no=0)
> w
[1] 0 NA 1 0 1 0 1 0
```

Ein logischer Vektor fungiert als `test`-Argument, und die Funktion `ifelse()` liefert einen Vektor entsprechender Länge als Ergebnis. Die Elemente dieses Ergebnisvektors sind ...

- identisch mit dem **yes**-Argument, wenn das korrespondierende **test**-Element gleich **TRUE** ist,
- identisch mit dem **no**-Argument, wenn das korrespondierende **test**-Element gleich **FALSE** ist
- und gleich **NA**, wenn das **test**-Element gleich **NA** ist.

Als **yes**- bzw. **no**-Argument sind auch Vektoren erlaubt. In diesem Fall ist das *i*-te Element des Ergebnisvektors identisch mit dem *i*-ten Element des **yes**-Vektors, wenn der *i*-te Wahrheitswert im **test**-Vektor gleich **TRUE** ist, und gleich dem *i*-ten Element des **no**-Vektors, wenn der *i*-te Wahrheitswert gleich **FALSE** ist.

7.3.1.4 Metrische Variable kategorisieren

Mit der Funktion **cut()** lässt sich ein numerischer Vektor über eine Liste von Aufteilungspunkten in einen Faktor wandeln. Im folgenden Beispiel enthält der Ausgangsvektor 100 standardnormalverteilte Zufallszahlen, die von der Funktion **rnorm()** erstellt werden (siehe Abschnitt 7.4.1.1):

```
> numvec <- rnorm(100, 0, 1)
> fac <- cut(numvec, breaks=c(-Inf,-2,-1,0,1,2,Inf))
> str(fac)
Factor w/ 6 levels "(-Inf,-2]","(-2,-1]",...: 2 4 3 5 1 4 2 4 3 5 ...
> table(fac)
fac
(-Inf,-2]  (-2,-1]  (-1,0]  (0,1]  (1,2]  (2, Inf]
           2         14         34         40         8         2
```

Durch die im Argumentvektor **breaks** enthaltenen Aufteilungspunkte werden links-offene und rechts-abgeschlossene Intervalle festgelegt. Die Funktion **table()** (siehe Abschnitt 9.1.2) liefert die absoluten Häufigkeiten der Kategorien.

7.3.2 Berechnen

Einen neuen Vektor aus bereits vorhandenen Vektoren zu berechnen, ist in **R** eine leichte Übung. Als Beispiel soll aus zwei Vektoren mit der Körpergröße (in cm) bzw. mit dem Körpergewicht (in kg) von 12 Personen

```
> gröÙe <- c(186,178,182,160,168,NA,165,179,158,175,176,176)
> gewicht <- c(80,71,75.5,65,66,76,55,76.5,50.5,80,62,NA)
```

der *Body Mass Index* (BMI) nach folgender Formel berechnet werden:

$$\frac{\text{Gewicht (in kg)}}{\text{GröÙe}^2 \text{ (in m)}}$$

Hinweise zur Formulierung des numerischen Ausdrucks:

- Als Symbol für das Potenzieren sind zwei unmittelbar aufeinanderfolgende Sternchen (**) oder ein Dach (^) zu verwenden.
- Die Körpergröße muss von der Einheit Zentimeter in die Einheit Meter umgerechnet werden.
- Wegen der Auswertungsprioritäten der beteiligten Operatoren (in absteigender Reihenfolge: Potenzieren, Dividieren) kann z.B. der folgende numerische Ausdruck verwendet werden:

```
> bmi <- gewicht/(gröÙe/100)^2
```

Das Ergebnis:

```
[1] 23.12406 22.40879 22.79314 25.39062 23.38435      NA 20.20202 23.87566
[9] 20.22913 26.12245 20.01550      NA
```

Erwartungsgemäß haben die beiden Fälle mit einem fehlenden Wert bei einer Ausgangsvariablen als Berechnungsergebnis den Wert **NA** erhalten.

7.4 Zufallszahlen erzeugen

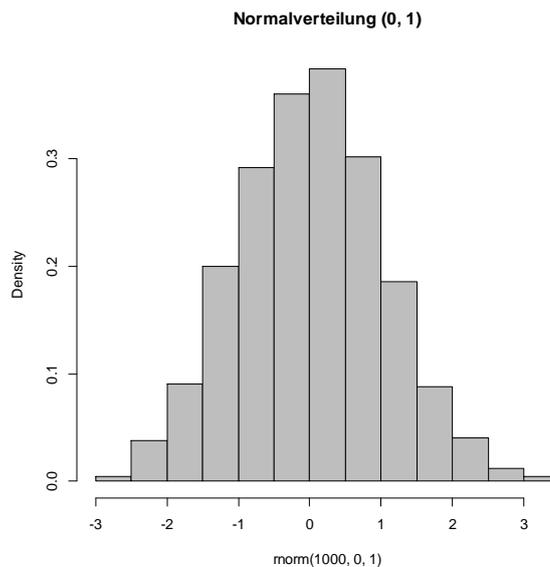
Für Simulationsstudien zu Phänomenen und Methoden der Statistik werden Stichproben aus einer Population mit bekannter Verteilung benötigt. So lässt sich etwa das Verhalten von statistischen Auswertungsverfahren bei bekannten Verteilungsverhältnissen untersuchen. **R** enthält diverse Funktionen, um Zufallsstichproben aus einer definierten Verteilung zu ziehen. Anschließend werden einige Vertreter vorgestellt.

7.4.1.1 Normalverteilte Zufallszahlen

Über die Funktion `rnorm()` erhält man n Zufallszahlen aus einer normalverteilten Population mit bestimmtem Erwartungswert und bestimmter Standardabweichung, z.B.:

```
> (sampnor <- rnorm(10, 0, 1))
[1] 0.5901100 0.5949126 0.7150877 1.1859644 0.8376390 -0.4352961
[7] -1.3281079 0.7989749 1.4540338 -1.5676006
```

Hier ist das Histogramm für eine erheblich größere Stichprobe aus einer Standardnormalverteilung zu sehen:



Es wurde durch den folgenden Aufruf der **R**-Funktion `hist()` erstellt (vgl. Abschnitt 8.2.5.4)

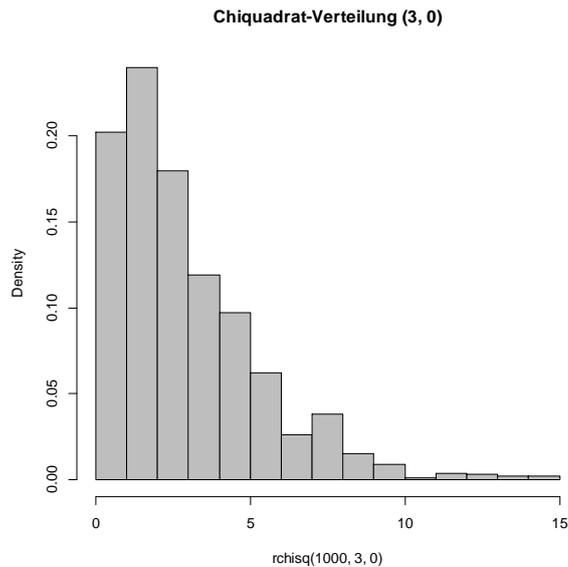
```
> hist(rnorm(1000, 0, 1), freq=FALSE, col="grey", main="Normalverteilung (0, 1)")
```

7.4.1.2 χ^2 -verteilte Zufallszahlen

Über die Funktion `rchisq()` erhält man n Zufallszahlen aus einer χ^2 -Verteilung, wobei die Anzahl der Freiheitsgrade und der Nonzentralitätsparameter einstellbar sind, z.B.:

```
> (sampchisq <- rchisq(10, 3, 0))
[1] 1.1991623 1.1538841 0.9253693 0.3756986 1.2608549 6.2716147 0.6392600
[8] 1.9859259 0.9282716 3.4587970
```

Hier ist das Histogramm für eine erheblich größere Stichprobe aus einer χ^2 -Verteilung mit 3 Freiheitsgraden und Nonzentralitätsparameter 0 zu sehen:



Es wurde durch den folgenden Aufruf der **R**-Funktion **hist()** erstellt (vgl. Abschnitt 8.2.5.4)

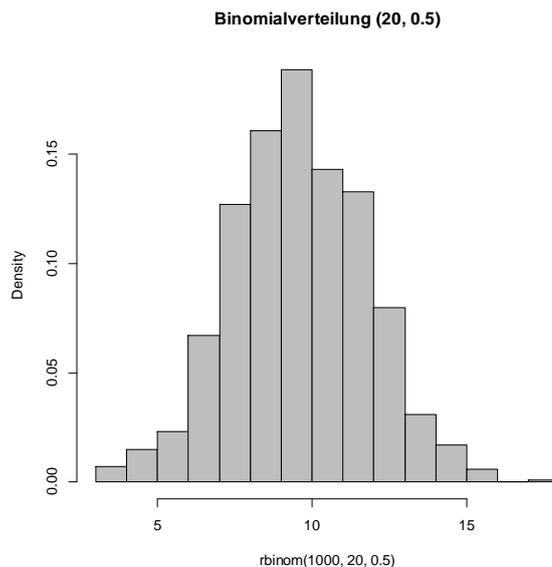
```
> hist(rchisq(1000, 3, 0), freq=FALSE, col="grey", main="Chiquadrat-Verteilung (3, 0)")
```

7.4.1.3 Binomialverteilte Zufallszahlen

Über die Funktion **rbinom()** erhält man n Zufallszahlen aus einer binomialverteilten Population mit bestimmten Parametern für das zugrunde liegende Bernoulli-Experiment (Anzahl der Wiederholungen, Wahrscheinlichkeit im Einzelexperiment), z.B.:

```
> sampbin <- rbinom(10, 20, 0.5); sampbin
[1] 13  8 12 11 10 12 12 10 10 10
```

Hier ist das Histogramm für eine erheblich größere Stichprobe aus derselben Binomialverteilung zu sehen:



Es wurde durch den folgenden Aufruf der **R**-Funktion **hist()** erstellt (vgl. Abschnitt 8.2.5.4)

```
> hist(rbinom(1000,20,0.5), freq=FALSE, col="grey", main="Binomialverteilung (20, 0.5)")
```

7.5 Auswahl von Fällen und/oder Variablen

7.5.1 Auswahl von Fällen

Oft soll bei einer Auswertung nur eine Teilmenge der Fälle in einer Datentabelle einbezogen werden. Mit der Funktion `subset()` lässt sich aus einer Datentabelle über einen logischen Ausdruck leicht eine Teilmenge mit allen Fällen ermitteln, die den logischen Ausdruck erfüllen, z.B.:

```
> dt <- data.frame(alter=c(45,32,58), ges=c('f','f','m'), mot=c(4,2,9))
> dt
  alter ges mot
1    45  f   4
2    32  f   2
3    58  m   9
> dtf <- subset(dt, ges=='f')
> dtf
  alter ges mot
1    45  f   4
2    32  f   2
```

Zu den im logischen Ausdruck verwendbaren relationalen und logischen Operatoren siehe Abschnitt 5.3.7.

Eine alternative Möglichkeit zur Fallauswahl bietet der in Abschnitt 5.3.6.3 beschriebene logische Indexvektor, z.B.:

```
> dtf <- dt[dt$ges=='f',]
```

Beim Indexzugriff auf eine Datentabelle muss hinter dem logischen Ausdruck für die Auswahl der Zeilen auf jeden Fall ein Komma stehen. Dann kann optional z.B. ein numerischer Vektor mit Spaltennummern zur einschränkenden Wahl der Variablen folgen (siehe Abschnitt 7.5.2).

7.5.2 Auswahl von Variablen

Bei einer großen Datentabelle kann es sich lohnen, die bei einer Auswertung tatsächlich benötigten Variablen in eine reduzierte Datentabelle zu extrahieren. Weil **R** alle geladenen Datenobjekte komplett im Hauptspeicher des Computers hält, ist bei einer großen Datentabelle der Verzicht auf irrelevante Variablen sinnvoll.

Die in Abschnitt 7.5.1 beschriebene Funktion `subset()` ermöglicht ergänzend zur bzw. anstatt einer Fallauswahl über das Argument `select` auch eine Variablenauswahl. Im ersten Beispiel findet ausschließlich eine Variablenauswahl statt, die über einen Indexvektor erfolgt:

```
> dt12 <- subset(dt, select=c(1,2))
```

Im zweiten Beispiel finden eine Fall- und eine Variablenauswahl statt. Zur Variablenauswahl dient ein logischer Indexvektor, wobei mit Hilfe der Funktion `grepl()` alle Variablen mit einer bestimmten Teilzeichenfolge im Namen ausgewählt werden:

```
> dtf <- subset(dt, ges=='f', select=grepl("e", names(dt)))
```

Auch beim Indexzugriff auf eine Datentabelle ist eine Variablenauswahl möglich, wobei die zweite Indextdimension (nach dem Komma) einen numerischen oder logischen Auswahlvektor erhält. Es folgen die beiden obigen Beispiele mit Indexsyntax:

```
> dt12 <- dt[, c(1,2)]
> dtf <- dt[dt$ges=='f', grepl("e", names(dt))]
```

7.6 Detail- und Master-Daten kombinieren

Wurden Daten mit einer hierarchischen Struktur erhoben (z.B. Clusterstichproben mit Beobachtungen zu den Individuen und den Gruppen, Längsschnittdaten mit Beobachtungen zu den Messzeitpunkten und den Subjekten), dann liegen die Daten zur Mikro- bzw. Makroebene oft in zwei getrennten Datentabellen vor. Im folgenden Beispiel sind die Daten der Mikroebene (Leistungsmessungen in den Fächern Mathematik und Geographie für Schüler in drei Klassen) in der Datentabelle `dfInd` zu finden und die Daten der Makroebene (Größe der Klasse) in der Datentabelle `dfGr`:

```
> dfInd <- data.frame(group=c(1,1,1,2,2,2,3,3,3), math=c(2,3,2,4,3,5,3,2,4),
+   geo=c(1,2,1,3,3,4,2,2,3))
> dfInd
  group math geo
1     1    2   1
2     1    3   2
3     1    2   1
4     2    4   3
5     2    3   3
6     2    5   4
7     3    3   2
8     3    2   2
9     3    4   3
> dfGr <- data.frame(group=c(1,2,3), size=c(21,40,32))
> dfGr
  group size
1     1   21
2     2   40
3     3   32
```

Um eine Mehrebenenanalyse zu ermöglichen, muss die Tabelle mit den Daten der Mikroebene um die Makrovariable `size` erweitert werden (siehe z.B. Bliese 2013), wobei die in beiden Datentabellen vorhandene Variable `group` für eine korrekte Zuordnung sorgt. Für das Zusammenführen der Daten eignet sich in **R** die Funktion `merge()`. Man übergibt die beiden Datenquellen als erstes bzw. zweites Argument und definiert anschließend die Zuordnung, was im Beispiel über eine einzige, in beiden Tabellen vorhandene und identisch benannte Variable geschehen kann:

```
> dfIndAug <- merge(dfInd, dfGr, by="group")
> dfIndAug
  group math geo size
1     1    2   1   21
2     1    3   2   21
3     1    2   1   21
4     2    4   3   40
5     2    3   3   40
6     2    5   4   40
7     3    3   2   32
8     3    2   2   32
9     3    4   3   32
```

7.7 Daten aggregieren

Oft sind für die Fälle in einer Datentabelle Gruppierungsvariablen vorhanden, und es wird eine neue Datentabelle benötigt, die für jede Ausprägung einer Gruppierungsvariablen oder für jede Ausprägungskombination von mehreren Gruppierungsvariablen genau einen Fall enthält, wobei die Werte dieses Falles bei den nicht zur Gruppierung verwendeten Variablen durch Aggregieren über die zugehörige Gruppe entstehen. Im folgenden Beispiel, das auch schon in Abschnitt 7.6 verwendet wurde, liegen Leistungsmessungen in den Fächern Mathematik und Geographie für Schüler in drei Gruppen (z.B. Klassen) vor:

```
> dfInd <- data.frame(group=c(1,1,1,2,2,2,3,3,3), math=c(2,3,2,4,3,5,3,2,4),
+   geo=c(1,2,1,3,3,4,2,2,3))
> dfInd
  group math geo
1     1    2   1
2     1    3   2
3     1    2   1
4     2    4   3
5     2    3   3
6     2    5   4
7     3    3   2
8     3    2   2
9     3    4   3
```

Mit der Funktion **aggregate()** lässt sich daraus eine Datentabelle mit den *Klassen* als Fällen erstellen, wobei die Ausprägungen bei den Variablen **math** und **geo** durch Mittelwertbildung in den Klassen entstehen:

```
> dfGr <- aggregate(dfInd[,2:3], list(dfInd$group), mean, na.rm=TRUE)
```

Im ersten Argument werden aus **dfInd** per Indexzugriff die Aggregierungsvariablen gewählt. Das zweite Argument legt die Liste mit den Gruppierungsvariablen fest, wobei das Beispiel mit *einer* Variablen auskommt. Im dritten Argument wird die Aggregierungsfunktion **mean** gewählt. Mit dem optionalen Argument **na.rm** wird schließlich dafür gesorgt, dass bei Gruppen mit fehlenden Werten *nicht* das Ergebnis **NA** resultiert, sondern das arithmetische Mittel der *vorhandenen* Werte. Im Beispiel entsteht die folgende Datentabelle:

```
> dfGr
  Group.1  math  geo
1       1 2.333333 1.333333
2       2 4.000000 3.333333
3       3 3.000000 2.333333
```

Um diese Daten auf der Gruppen- bzw. Makroebene (z.B. für eine Mehrebenenanalyse) wieder mit den Daten auf der Mikroebene (in der Tabelle **dfInd**) zusammen zu führen, ist die in Abschnitt 7.6 vorgestellte Funktion **merge()** zu verwenden.

7.8 Sekundärstichproben ziehen

Mit der Funktion **sample()** kann man aus einem Vektor mit einer Primärstichprobe durch Ziehen mit Zurücklegen eine Sekundärstichprobe gleicher Größe gewinnen, z.B.:

```
> prim <- 1:9
> sec <- sample(prim, size=length(prim), replace=TRUE)
[1] 5 9 6 5 2 5 3 2 9
```

Um für einen numerischen Vektor eine zufällige **Permutation** zu gewinnen, zieht man eine Stichprobe vom Originalumfang mit dem voreingestellten Wert **FALSE** für das Argument **replace**.

8 Grafik-Optionen in R

Das traditionelle Grafiksystem in **R** basiert auf dem stets installierten und auch automatisch geladenen Paket **graphics**. Daneben bietet **R** als alternatives Grafiksystem die so genannte **Gittergrafik** (engl. *grid graphics*), auf der die folgenden Optionen zur Grafikproduktion aufbauen:

- **Trellis-Grafiken**

Das zur Grundinstallation von **R** gehörige Paket **lattice** enthält die so genannten *Trellis-Grafiken* (siehe z.B. Murrell 2011).

- **Grammar of Graphics**

Das zusätzlich zu installierende, von Hadley Wickham erstellte und in einer Monographie beschriebene Paket **ggplot2** (siehe Wickham 2009, 2016) basiert auf der von Wilkinson (2005) entwickelten *Grammar of Graphics* (daher der Namensanfang *gg*). Dieselbe Grundlogik ist übrigens auch in der von SPSS unterstützten *Graphics Production Language* (GPL) realisiert.

Im Wettbewerb der **R**-Grafik-Optionen zeichnet sich eine Tendenz zur Bevorzugung des Pakets **ggplot2** ab (siehe z.B. Field 2012, S. 121). Robert Muenchen (2011, S. 444f) urteilt:

The point is that **ggplot2** gives you the broadest range of graphical options that R offers.

Auf das ohne jedes Installieren und/oder Laden von Paketen verfügbare traditionelle Grafiksystem sollte man nicht verzichten, denn:

- Viele vorhandene Lösungen nutzen die traditionelle Grafikproduktion.
- Derzeit unterstützt ggplot2 keine dreidimensionalen Grafiken (z.B. Funktionsplots)

Daher werden im Manuskript das traditionelle Grafiksystem (siehe Abschnitt 8.2) sowie die Grafikproduktion mit dem Paket **ggplot2** behandelt (siehe Abschnitt 8.3). Der Abschnitt 8.1 über Ausgabeformate und -geräte ist für beide Grafiksysteme relevant.

8.1 Ausgabeformate und -geräte

Die Ausgaben der Grafikfunktionen werden zu einem Ausgabegerät kanalisiert, das eine Bildschirm- oder Dateiausgabe produziert.

8.1.1 Verfügbare Ausgabegeräte

Zum Öffnen eines Ausgabegeräts von bestimmtem Typ ist jeweils die zuständige Funktion aufzurufen. Über die Namen dieser Funktionen und damit über die unterstützten Ausgabeformate informiert die **R**-Hilfe nach der Anweisung

```
> ?device
```

Bei **R** 3.2.5 erscheint unter Windows die folgende Geräteliste:

- **windows** The graphics device for Windows (on screen, to printer and to Windows metafile)
- **pdf** Write PDF graphics commands to a file
- **postscript** Writes PostScript graphics commands to a file
- **xfig** Device for XFIG graphics file format
- **bitmap** bitmap pseudo-device via Ghostscript (if available)
- **pictex** Writes TeX/PicTeX graphics commands to a file (of historical interest only)
- **cairo_pdf, cairo_ps** PDF and PostScript devices based on cairo graphics
- **svg** SVG device based on cairo graphics
- **png** PNG bitmap device
- **jpeg** JPEG bitmap device
- **bmp** BMP bitmap device
- **tiff** TIFF bitmap device

Während das erste Gerät in dieser Liste zur Ausgabe in ein Grafikkfenster (siehe Abschnitt 8.1.2) taugt, sind die anderen Geräte jeweils mit einer Datei verbunden.

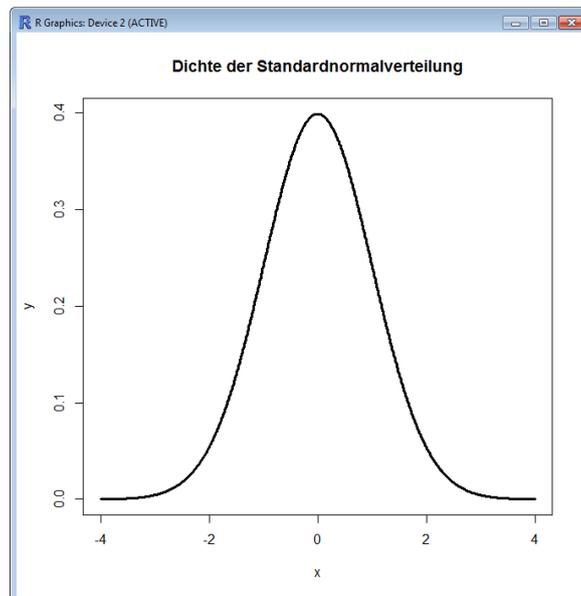
Das in der obigen Liste fehlende Gerät zur Direktausgabe in eine Datei vom Typ *Enhanced Windows Metafile (EMF)* wird (zumindest in der Windows-Version von **R**) durch die Funktion **win.metafile()** bereitgestellt.

8.1.2 Grafikkfenster

Bei Verwendung der RGui-Bedienoberfläche ist das voreingestellte Ausgabegerät ein Grafikkfenster, das beim ersten Aufruf einer Grafikfunktion automatisch initialisiert wird. Als Beispiel erstellen wir mit dem traditionellen Grafiksystem ein Liniendiagramm mit der Standardnormalverteilungsdichte:

```
> x <- seq(-4,4,0.01); y <- dnorm(x)
> plot(x, y, type="l", main="Dichte der Standardnormalverteilung", lwd=3)
```

Per **seq()** - Funktion (siehe Abschnitt 5.3.4.2.1) entsteht ein Vektor mit x-Koordinaten. Von der **dnorm()** - Funktion erhalten wir einen gleichlangen y-Vektor mit den Funktionswerten der Standardnormalverteilungsdichte zu den x-Werten. Per **plot()** - Funktion entsteht ein Liniendiagramm mit der Dichte, das in einem Grafikkfenster erscheint:



Dieses Ausgabegerät hat den Typ **windows** und kann mit der gleichnamigen Funktion auch explizit erstellt werden, z.B.

```
> windows(5, 5)
```

Dieses Vorgehen bietet die Möglichkeit, eine alternative Breite und Höhe anstelle der voreingestellten Werte von jeweils 7 Zoll (engl.: *Inch*) festzulegen (= $7 \cdot 2,54 \text{ cm} = 17,78 \text{ cm}$).

Durch jeden Aufruf der Funktion **windows()** wird ein neues Grafikkfenster erstellt und zum aktiven Ausgabegerät ernannt (vgl. Abschnitt 8.1.5 zur Verwaltung der Ausgabegeräte).

Jedes Ausgabegerät hat neben einem Typ auch eine Nummer, die bei einem Grafikkfenster in der Titelzeile erscheint (siehe Beispiel).

Ist im RGui ein Grafikkfenster aktiv, kann sein Inhalt über den Menübefehl

Datei > Speichern als

in diversen Dateiformaten gespeichert werden (**Metafile**, **Postscript**, **PDF**, **Png**, **Bmp**, **TIFF** und **Jpeg**). Die ersten drei Formate sind vektoriell organisiert, enthalten also mathematisch definierte Objekte und erlauben somit eine verlustfreie Größenänderung (vgl. Abschnitt 8.1.3). Die letzten drei Formate besitzen eine Bitmap- bzw. Rasterorganisation mit einer Auflösung von 96×96 Pixeln.

Über den Menübefehl

Datei > Kopieren in Zwischenablage als Metafile oder **als Bitmap**

kann der Inhalt des aktiven Grafikfensters im vektoriell organisierten Metafile- oder im Bitmap-Format in die Windows-Zwischenablage kopiert werden, und über den Menübefehl

Datei > Drucken

ist eine Druckausgabe möglich.

Über das Kontextmenü zum Grafikfenster sind dieselben Aktionen möglich: Dateiausgabe (mit einer Einschränkung auf die Formate Metafile und Postscript), Zwischenablagentransfer und Drucken.

Bei aktivem Grafikfenster lässt sich über den Menübefehl

History > Aufzeichnen

die Aufzeichnung der erzeugten Grafiken ein- bzw. ausschalten. Zwischen aufgezeichneten Diagrammen kann man über die Tasten  bzw.  wechseln.

Um ein Grafikfenster mit Aufzeichnung per Syntax zu öffnen, wählt man den Funktionsaufruf:

```
> windows(record=TRUE)
```

8.1.3 Ausgabe in eine Datei

Unter Windows stehen zwei Möglichkeiten zur Verfügung, um ein Diagramm in eine Datei zu schreiben:

- Man erstellt ein Grafikfenster durch die implizite oder explizite Verwendung der Funktion **windows()** und speichert seinen Inhalt in eine Datei, was über den Menübefehl

Datei > Speichern als

in den Formaten **Metafile**, **Postscript**, **PDF**, **Png**, **Bmp**, **TIFF** und **Jpeg** oder via Kontextmenü in den Formaten **Metafile** und **Postscript** gelingt. Die ersten drei Formate sind vektoriell organisiert, enthalten also mathematisch definierte Objekte und erlauben somit eine verlustfreie Größenänderung. Die letzten drei Formate besitzen eine Bitmap- bzw. Rasterorganisation mit einer Auflösung von 96 dpi (*dots per inch*), die für Publikationszwecke *nicht* ausreicht.

- Man öffnet ein Ausgabegerät, das direkt mit einer Datei verbunden ist, durch Aufruf der Format-spezifischen Funktion (siehe Abschnitt 8.1.1). Die Ausgabegeräte **postscript**, **pdf**, **xfig**, **cairo_pdf**, **cairo_ps**, **svg** und **win.metafile** führen zu einem vektoriell organisierten (also auflösungsunabhängigen) Ergebnis. Die anderen Ausgabegeräte liefern ein Bitmap-Resultat, so dass die Größe und die Auflösung passend gewählt werden müssen, um eine gute Qualität zu erzielen.

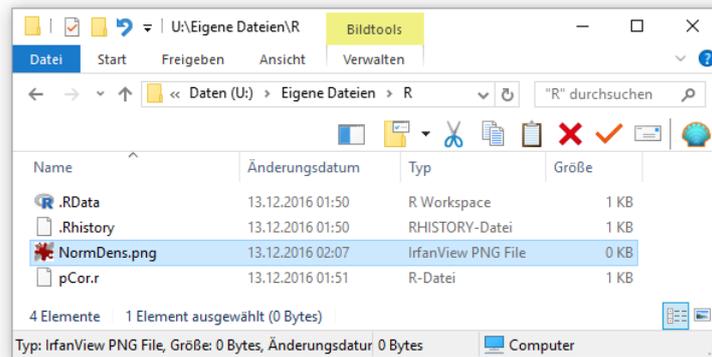
Im folgenden Beispiel entsteht eine PNG-Datei mit einer Breite und Höhe von 10 cm bei einer Auflösung von 600 dpi:

```
> png("NormDens.png", 10, 10, units="cm", res=600)
```

Nach dem Öffnen des Ausgabegeräts führt man die Grafikfunktionsaufrufe durch, z.B.:

```
> x <- seq(-4,4,0.01); y <- dnorm(x)
> plot(x, y, type="l", main="Dichte der Standardnormalverteilung", lwd=3)
```

Daraufhin wird die Ausgabedatei (mit der Größe von 0 Bytes) angelegt, die durch das RGui blockiert ist:



Nach Fertigstellung des Diagramms schließt man das Ausgabegerät, z.B.:

```
> dev.off()
windows
2
```

Daraufhin erfolgt die Ausgabe, und die Datei wird zur Verwendung durch andere Programme freigegeben. Außerdem erfährt man Typ und Nummer des nunmehr aktiven Ausgabegeräts.

8.1.4 R-Diagramme an Microsoft Office übergeben

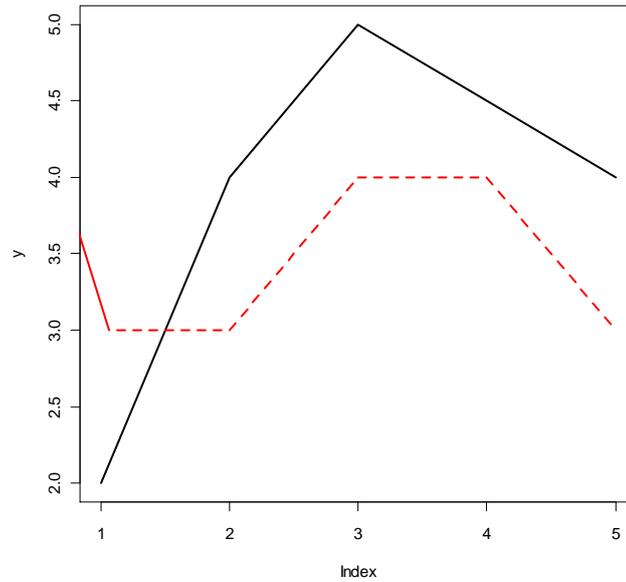
Bei der Übernahme eines **R**-Diagramms in ein MS-Office-Dokument ist unter Windows grundsätzlich das Metafile-Format gut geeignet, weil es auflösungsunabhängig eine optimale Qualität bietet und außerdem eine Bearbeitung einzelner Diagrammbestandteile in der Office-Anwendung erlaubt. Allerdings gilt diese Empfehlung nur eingeschränkt für Diagramme mit Kurven. Weil Kurven in Geradensegmente zerlegt werden, und die Anzahl der Segmente im Metafile-Format offenbar begrenzt ist, resultiert eine zu grobe Auflösung mit einem deutlich erkennbaren Treppmuster. Im Zusammenhang mit dem Grafikpaket **ggplot2** ist von Bedeutung, dass die hier möglichen Transparenzeffekte vom Metafile-Format nicht unterstützt werden.

Andere bei der **R**-Grafikproduktion mögliche Vektorformate (z.B. PDF, Postscript, SVG) lassen sich leider in Office nicht nutzen.

Aus einem **R**-Grafikfenster den Transfer zu Office im Bitmap-Format vorzunehmen, ist wegen der geringen Auflösung von 96 dpi (*dots per inch*) meist nicht sinnvoll.

Für Windows-Anwender resultieren die folgenden Empfehlungen für den Diagrammtransfer von **R** zu MS-Office:

- Bei einem Diagramm *ohne* Kurven und Transparenzeffekte führt der Transfer im Metafile-Format meist zu einem akzeptablen Ergebnis (via Zwischenablage aus dem Grafikfenster oder nach einer Dateiausgabe), z.B.:



- Bei einem Diagramm *mit* Kurven oder Transpaarenzeffekten sollte über ein **png**-Ausgabegerät eine Datei mit hoher Auflösung (mindestens 300 dpi) erstellt und anschließend in das Office-Dokument eingefügt werden, z.B.:

Diagramm erstellt über das Gerät: <code>win.metafile("NormDens.emf")</code>	Diagramm erstellt über das Gerät: <code>png("NormDens.png",7,7,units="in",res=600)</code>
<p style="text-align: center;">Dichte der Standardnormalverteilung</p>	<p style="text-align: center;">Dichte der Standardnormalverteilung</p>

- Metafile-Diagramme sorgen gelegentlich bei der PDF-Produktion aus MS-Office für Darstellungsfehler, so dass man auf ein Bitmap-Format ausweichen muss.

8.1.5 Ausgabegeräte verwalten

Über die Funktion `dev.list()` erhält man eine Liste der geöffneten Ausgabegeräte, z.B.:

```
> dev.list()
      windows png:NormDens.png
         2         3
```

Sind mehrere Ausgabegeräte offen, ist eines als aktiv ausgezeichnet und das Ziel für die Ausgabe von Grafikfunktionen. Man erfährt seinen Typ und seine Nummer über die Funktion `dev.cur()`, z.B.:

```
> dev.cur()
png:NormDens.png
      3
```

Ein aktives *Grafikfenster* zeigt in der Titelzeile die Zeichenfolge **(ACTIVE)**.

Um ein vorhandenes anderes Ausgabegerät zu aktivieren, ruft man die Funktion `dev.set()` mit der Geräte-
nummer auf, z.B.:

```
> dev.set(2)
windows
      2
```

Eine Besonderheit der Geräte **postscript** und **pdf** besteht darin, dass bei der sukzessiven Erstellung von *mehreren* Diagrammen ein *mehrseitiges* Dokument entsteht, während bei sonstigen Ausgabegeräten ein neues Diagramm das bisherige ersetzt.

Auf das aktive Grafikfenster hat der Funktionsaufruf `dev.off()` denselben Effekt wie ein Klick auf das Schließkreuz in der Titelzeile.

Um *alle* Ausgabegeräte zu schließen, verwendet man die Anweisung:

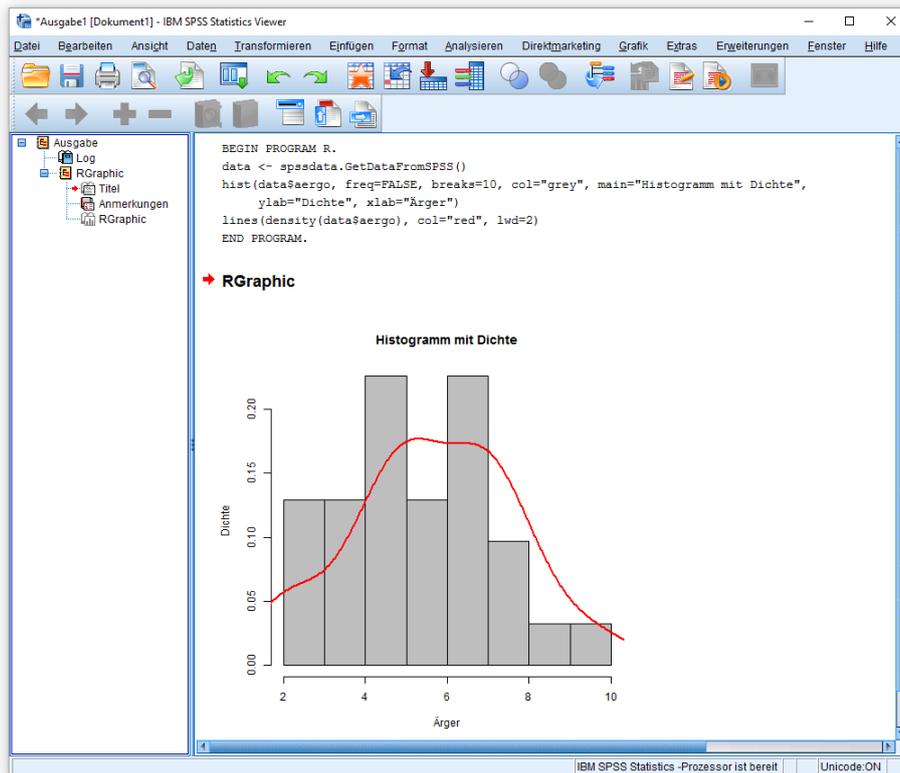
```
> graphics.off()
```

8.1.6 R-Diagramme im SPSS-Ausgabefenster

Es ist selbstverständlich möglich, **R**-Diagramme per SPSS-Syntaxfenster zu erstellen, und wir werden im Manuskript noch mehrfach von dieser Möglichkeit mehrfach Gebrauch machen. Im folgenden Beispiel wird mit Funktionen aus dem traditionellen **R**-Grafiksystem ein Histogramm mit Dichteschätzung zu einer Variablen der SPSS-Arbeitsdatei erstellt:

```
BEGIN PROGRAM R.
data <- spssdata.GetDataFromSPSS()
hist(data$aergo, freq=FALSE, breaks=10, col="grey", main="Histogramm mit Dichte",
      ylab="Dichte", xlab="Ärger")
lines(density(data$aergo), col="red", lwd=2)
END PROGRAM.
```

Das Ergebnis landet im SPSS-Ausgabefenster



und kann von dort via Windows-Zwischenablage in andere Anwendungen (z.B. Word) weiterbefördert werden.

Die Qualität der **R**-Diagramme im SPSS-Ausgabefenster reicht für den Forschungsalltag, aber in der Regel nicht für Publikationszwecke. Bei höheren Qualitätsansprüchen sollten die Hinweise im Abschnitt 8.1.4 beachtet werden. Für die Ausgabe eines Diagramms in eine Datei kann man das RGui oder das SPSS-Syntaxfenster benutzen, z.B.:

```

BEGIN PROGRAM R.
data <- spssdata.GetDataFromSPSS()
png("u:/eigene dateien/r/NormDens.png", 10, 10, units="cm", res=600)
hist(data$aergo, freq=FALSE, breaks=10, col="grey", main="Histogramm mit Dichte",
      ylab="Dichte", xlab="Ärger")
lines(density(data$aergo), col="red", lwd=2)
END PROGRAM.
    
```

8.2 Das traditionelle R-Grafiksystem

In diesem Abschnitt wird eine kleine Auswahl der in **R** verfügbaren traditionellen Optionen zur Darstellung von Daten und mathematischen Funktionen vorgestellt.

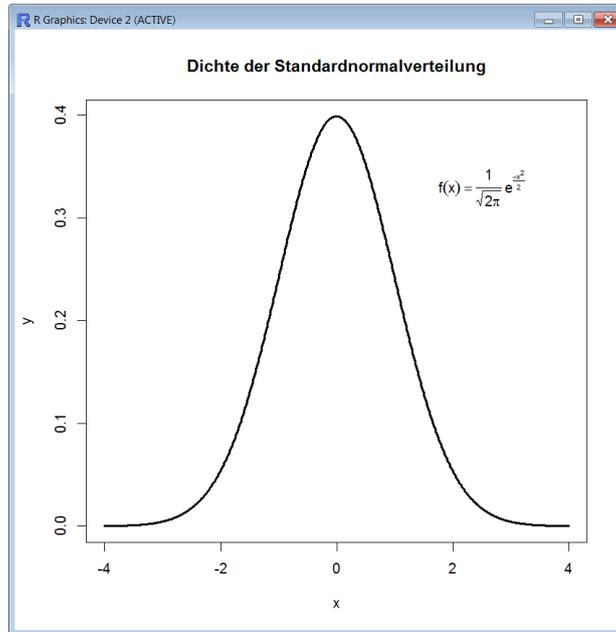
8.2.1 High- und Low-Level - Grafikfunktionen

Die traditionelle **R**-Grafik kennt High-Level - Funktionen, die ein vollständiges Diagramm erstellen (z.B. **plot()**) und Low-Level - Funktionen, die ein Diagramm um ein Element (z.B. eine Beschriftung) erweitern (z.B. **text()**). Das schafft Flexibilität und sorgt dafür, dass sich die traditionelle Grafik gegen starke Konkurrenz im **R**-Universum immer noch behaupten kann.

Mit den folgenden Anweisungen

```
> x <- seq(-4,4,0.01)
> y <- dnorm(x, mean = 0, sd = 1)
> plot(x,y, type="l", main="Dichte der Standardnormalverteilung", lwd=3)
> nd <- expression(f(x) == frac(1,sqrt(2*pi))~e^{frac(-x^2,2)})
> text(2.5, 0.33, nd)
```

wird per **plot()** - Funktion ein Liniendiagramm mit der Standardnormalverteilungsdichte gezeichnet und anschließend per **text()** - Funktion die Definition der Dichte ergänzt:



Mit der **expression()** - Funktion lässt sich mathematische Notation in Diagramme einfügen. Das Beispiel stammt aus Ligges (2007, S. 170).

8.2.2 Grafikparameter

Die traditionelle **R**-Grafik kennt ca. 70 Grafikparameter, die sich für ein Ausgabegerät oder für einen Grafikfunktionsaufruf setzen lassen:

- Sollen Parameter mit Gültigkeit bei allen Grafikfunktionsaufrufen für das aktive Ausgabegerät gesetzt werden, ist die Funktion **par()** mit der Syntax

par(name=wert, ...)

zu verwenden. Mit dem Schließen des Ausgabegeräts sind die gesetzten Grafikparameter hinfällig. Ein **par()** - Aufruf ohne Argumente protokolliert die aktuellen Werte der Grafikparameter für das aktive Ausgabegerät.

- Viele Parameter lassen sich auch als Transferargumente in einem konkreten Grafikfunktionsaufruf (High- oder Low-Level) mit lokaler Gültigkeit verwenden.

Anschließend werden einige Parameter vorgestellt, die entweder für eine generelle Einstellung per **par()** in Frage kommen oder für mehrere Grafikfunktionen relevant sind. Später folgen noch Parameter, die vor allem für spezielle Grafikfunktionen (z.B. **plot()**) von Interesse sind. Die folgende Liste eignet sich weniger zum Studieren als zum Nachschlagen:

- **col**

Die Zeichenfarbe kann über einen Farbnamen oder einen RGB-Wert aus drei Hexadezimalzahlen von 00 bis FF (optional ergänzt durch einen Transparenzwert im selben Wertebereich) festgelegt werden, z.B.:

```
> boxplot(x, col="red")
> boxplot(x, col="#ff0000") # Rot
> boxplot(x, col="#ff000080") # Rot, halb-transparent
```

Über die Funktion **colors()** erhält man einen Vektor mit den 657 verfügbaren Farbnamen.¹

- **col.axis, col.lab, col.main, col.sub**

Mit diesen Parametern wird die Farbe für bestimmte Bestandteile eines Diagramms gewählt: Teilstrichbeschriftungen (**col.axis**), Achsenbeschriftungen (**col.lab**), Überschriften erster und zweiter Ordnung (**col.main, col.sub**).

- **bg**

Mit dem Parameter **bg** wird die Hintergrundfarbe für ein Ausgabegerät gesetzt. Voreinstellung für **bg** ist Weiß. Einige Grafikfunktionen (z.B. **points()**) haben ein Argument mit demselben Namen, aber unterschiedlicher Bedeutung.

- **fg**

Mit dem Parameter **fg** wird die Vordergrundfarbe gesetzt. Geschieht dies in einer Grafikfunktion, sind vor allem Achsen und Rahmen betroffen, wobei sich nicht alle Grafikfunktionen gleich verhalten. Geschieht es in einem **par()** - Aufruf, wird außerdem der Parameter **col** auf denselben Wert gesetzt. Voreinstellung für **fg** ist Schwarz.

- **family**

Mit diesem Parameter wählt man eine Schriftfamilie. Generell verfügbar sind:

Familienname	Unter Windows abgebildet auf
mono	TT Courier New
serif	TT Times New Roman
sans	TT Arial

Über die Funktion **windowsFonts()** lassen sich andere im Windows-Betriebssystem des Rechners installierte Schriftarten einbinden (siehe **R**-Hilfe).

- **font, font.axis, font.lab, font.main, font.sub**

Mit diesem Parameter wählt man eine Schriftauszeichnung:

Wert	Auszeichnung
1	normal
2	fett
3	<i>kursiv</i>
4	<i>fett-kursiv</i>
5	Verwendet die Schriftart Symbol

Diese Einstellung kann für die gesamte Grafik (**font**) oder für die Teilstrichbeschriftungen (**font.axis**), die Achsenbeschriftungen (**font.lab**), sowie die Überschriften erster und zweiter Ordnung (**font.main, font.sub**) vorgenommen werden.

¹ In der folgenden PDF-Datei (abgerufen am 13.12.2016) sind die in **R** verfügbaren Farbnamen mit Musterflächen dokumentiert:

<http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf>

- **ps**
Dieser Parameter legt die generelle Schriftgröße für ein Ausgabegerät in Point Size - Einheiten fest (eine Einheit = 1/72 Zoll, voreingestellter **ps**-Wert: 12). Wie sich daraus über Skalierungsfaktoren die Größen für spezielle Beschriftungen ableiten lassen, wird anschließend beschrieben.
- **cex, cex.axis, cex.lab, cex.main, cex.sub**
Durch einen positiven Skalierungsfaktor kann die Größe von Symbolen (**cex**), Teilstrichbeschriftungen (**cex.axis**), Achsenbeschriftungen (**cex.lab**), sowie Überschriften erster und zweiter Ordnung (**cex.main, cex.sub**) im Vergleich zum Standard (Wert = 1) reduziert (Wert < 1) oder erhöht werden (Wert > 1).
- **lwd**
Durch diesen positiven Skalierungsfaktor kann die Linienstärke im Vergleich zum Standard (Wert = 1) reduziert (Wert < 1) oder erhöht werden (Wert > 1), z.B.:

```
> boxplot(x, lwd=2)
```


Durch Verwendung des Parameters in einer High-Level - Grafikfunktion oder in einem **par()** - Aufruf lässt sich die Linienstärke des Diagramminhalts beeinflussen. Durch Verwendung in einer Low-Level - Grafikfunktion (z.B. **axis()**, siehe Abschnitt 8.2.4.2) lassen sich die Linienstärken anderer Diagrammbestandteile beeinflussen.

8.2.3 Beschriftungen

Mit den folgenden Argumenten lassen sich in High-Level - Grafikfunktionen Beschriftungen einfügen:

- **main, sub**
Titel und Untertitel
- **xlab, ylab**
Beschriftungen für die X- bzw. Y-Achse
Um die voreingestellten Beschriftungen abzuschalten, weißt man den Argumenten eine leere Zeichenfolge oder den Wert **NA** zu, z.B.:

```
> plot(x, y, type="S", xlab=NA, ylab=NA)
```

8.2.4 Die generische Funktion plot()

Mit der generischen High-Level - Grafikfunktion **plot()** erstellt man für k Koordinatenpaare $((x, y)$ -Punkte) ein Liniendiagramm (siehe Abschnitt 8.2.5.1) oder ein Streudiagramm (siehe Abschnitt 8.2.5.2).

8.2.4.1 Argumente

Die Funktion **plot()** kennt u.a. die folgenden Argumente:

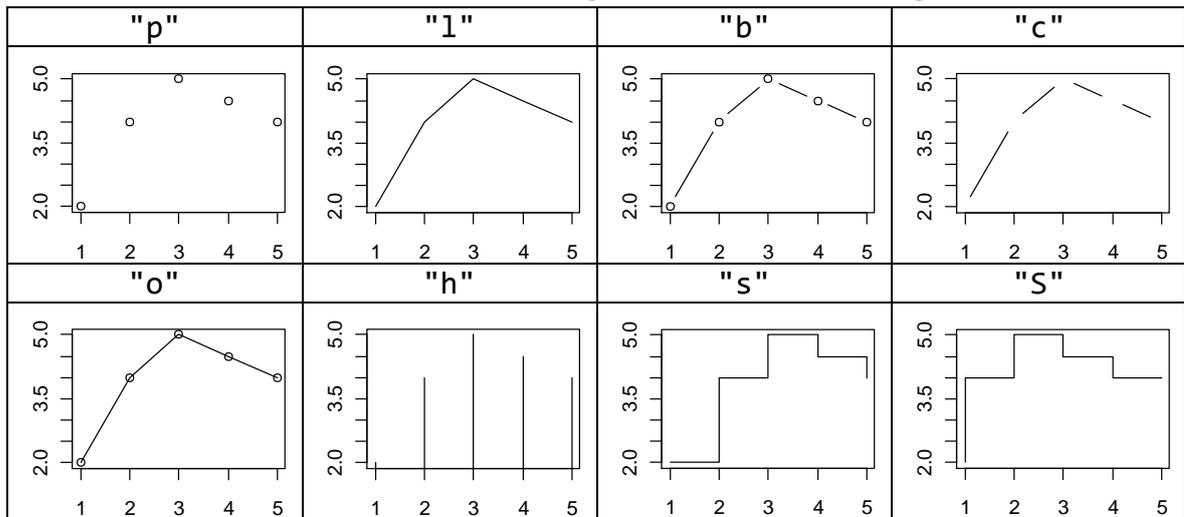
- **x**
 k -elementiger Vektor mit den x -Koordinaten der Punkte
Lässt man den **x**-Vektor weg, verwendet **R** die Indexnummern der **y**-Elemente (1, 2, 3, ...).
- **y**
 k -elementiger Vektor mit den y -Koordinaten der Punkte

• **type**

Zur Gestaltung von Punktmarkierungen und/oder Linien stehen folgende **Plot-Typen** zur Verfügung:

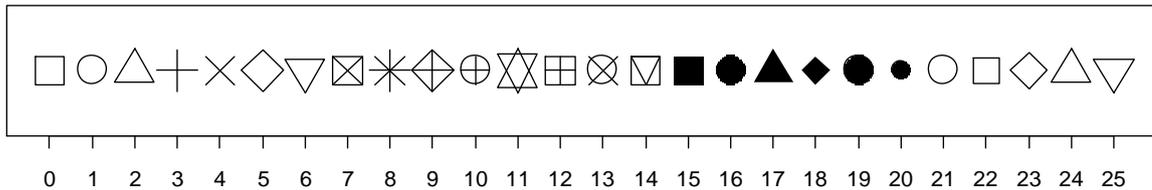
- "p" Nur Punkte (= Voreinstellung)
- "l" Nur Linien
- "b" Beides (Punkte und Linien)
- "c" Linien mit Lücken an den x-Positionen (nur die Linien aus Typ "b")
- "o" Beides (Punkte und Linien, „overplotted“)
- "h" Senkrechte Linien von den y-Werten bis zur X-Achse
- "s", "S" Treppenstufen in zwei Varianten
- "n" leeres Diagramm zur Vorbereitung für spätere Low-Level - Ausgaben

Um bei den folgenden Plot-Typ - Demonstrationen Platz zu sparen, wurde als Ausgabegerät ein Grafikenfenster mit einer Breite und Höhe von lediglich 3 Zoll verwendet (vgl. Abschnitt 8.1.2):



• **pch**

Das Symbol für die Datenpunkte kann aus der folgenden Palette¹



über seine Nummer gewählt werden, z.B.:

```
> x <- 1:5; y <- c(2, 4, 5, 4.5, 4)
> plot(x, y, pch=4)
```

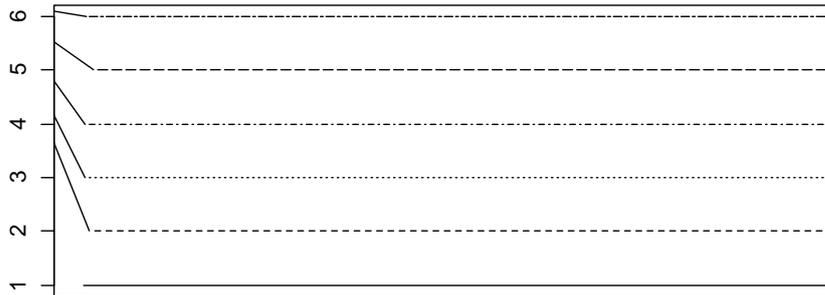
Statt für alle Datenpunkte dasselbe Symbol zu verwenden, kann man dem Argument **pch** einen Vektor übergeben, der für jeden Datenpunkt ein eigenes Symbol enthält (siehe Abschnitt 8.2.5.2.3).

¹ Die Abbildung wurde folgendermaßen mit **R** erstellt:

```
> windows(10, 3)
> plot(0:25, rep(1,26), pch=0:25, cex=3, xlab="", ylab="", yaxt="n")
> axis(side = 1, at = 0:25)
```

- **lty**

Der Linientyp kann über seine Nummer aus der folgenden Palette¹



gewählt werden, z.B.:

```
> plot(x, y, type="l", lty=4)
```

- **xlim, ylim**

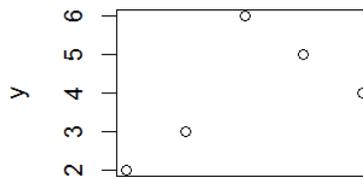
Begrenzung des Darstellungsbereichs für die X- bzw. Y-Achse, z.B.:

```
> plot(x, y, type="s", xlim=c(1,3))
```

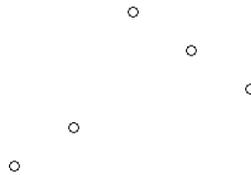
- **xaxt, yaxt, axes**

Ist die Erstellung individueller Achsen über die Low-Level - Funktion **axis()** geplant (siehe Abschnitt 8.2.4.2), schaltet man die zu ersetzenden Standardvarianten ab:

- Der Wert "n" für das Argument **xaxt** bzw. **yaxt** verhindert die Teilstriche und Teilstrichbeschriftungen für die X- bzw. Y-Achse. Im folgenden Beispiel ist die X-Achse deaktiviert:

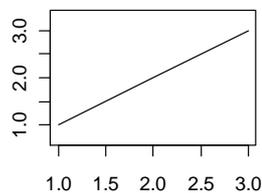


- Der Wert "n" für das Argument **axes** verhindert *beide* Achsen komplett (Teilstriche, Teilstrichbeschriftungen, Linien, Rahmen), z.B.:



- **asp**

Dieses Argument bestimmt den Quotienten aus der Länge einer Y-Achseneneinheit und der Länge einer X-Achseneneinheit. Im folgenden Beispiel ist der Abstand zwischen den X-Werten 1 und 2 doppelt so groß wie der Abstand zwischen den Y-Werten 1 und 2:



¹ Die Abbildung wurde folgendermaßen mit **R** erstellt:

```
> windows(7,4)
> plot(1:6, 1:6, type="n", xlab="", ylab="", xaxt="n")
> for (i in 1:6) lines(1:6, rep(i,6), lty=i)
```

Neben ihren eigenen Argumenten akzeptiert die Funktion **plot()** auch viele Grafikparameter (siehe Abschnitt 8.2.2) als Transferargumente, z.B. zur Änderung der Zeichenfarbe:

```
> plot(x, y, col="red")
```

Statt für alle Datenpunkte dieselbe Farbe zu verwenden, kann man dem Argument **col** einen Vektor übergeben, der für jeden Datenpunkt eine eigene Farbe enthält (siehe Abschnitt 8.2.5.2.3).

8.2.4.2 Ergänzende Low-Level - Grafikfunktionen

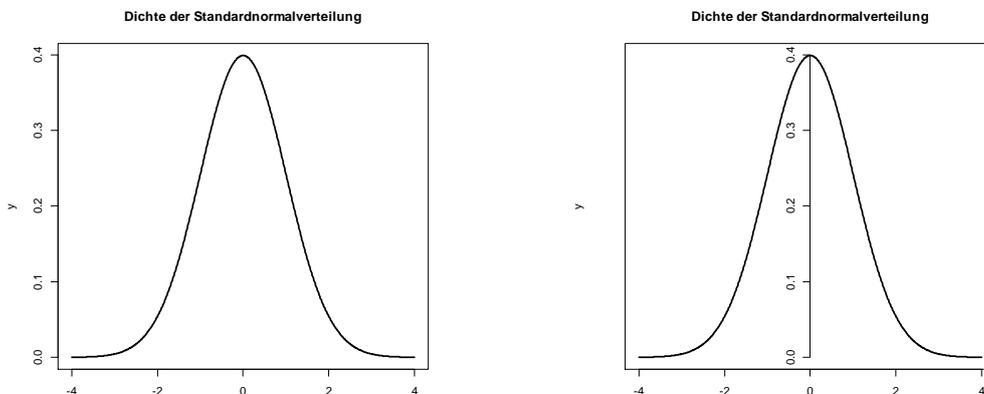
Zur Ergänzung zur **plot()** - Ausgabe kommen u.a. die folgenden Low-Level - Grafikfunktionen in Frage:

8.2.4.2.1 Achsengestaltung mit der Funktion **axis()**

Um eine spezielle Achsengestaltung zu erreichen, schaltet man die Standardvariante über das **plot()** - Argument **xaxt** bzw. **yaxt** ab (siehe Abschnitt 8.2.4.1) und erstellt mit der **axis()** - Funktion ein individuelles Exemplar, wobei (neben generellen Grafikparametern, siehe Abschnitt 8.2.2) die folgenden Argumente verfügbar sind:

- **side**
Über eine Zahl wird der Erscheinungsort der Achse festgelegt (1 = unten, 2 = links, 3 = oben, 4 = rechts).
- **at**
Ein numerischer Vektor mit den Positionen für die Teilstriche, z.B.:

```
axis(side = 1, at = 0:25)
```
- **tck**
Über einen Anteil der Zeichenflächengröße legt man die Länge der Teilstriche fest, wobei positive Werte auf der Innenseite und negative Werte auf der Außenseite der Achse erscheinen (Voreinstellung = -0,02). Mit dem Wert 1 erhält man Gitterlinien.
- **labels**
Mit diesem Zeichenketten-Vektor kann man die Teilstrichbeschriftungen festlegen. Unterlässt man es, kommen die Werte im Vektor **at** zum Einsatz.
- **las**
Sollen bei einer senkrechten Achse (**side** gleich 2 oder 4) die Teilstrichbeschriftungen parallel zur Achse (Wert = 0) oder orthogonal zur Achse (Wert = 1) geschrieben werden?
- **pos**
Per Voreinstellung landet eine Achse am Rand des Ausgaberechtecks (siehe linkes Beispiel).



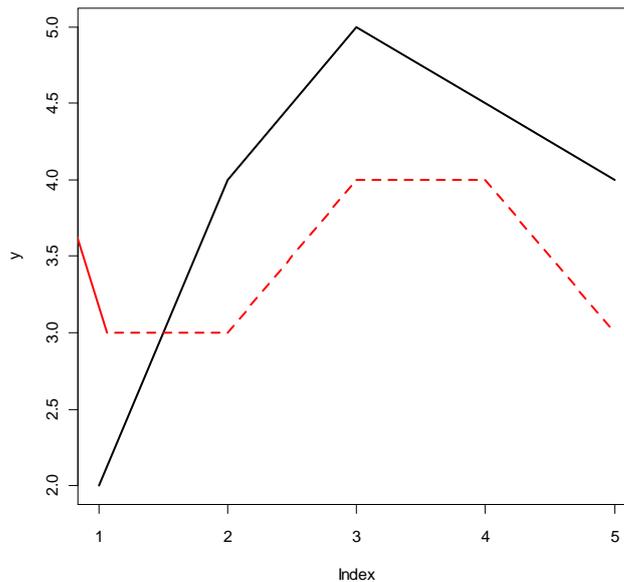
Mit dem Argument **pos** lässt sich der Schnittpunkt mit der orthogonalen Achse festlegen, was im rechten Beispiel für die Y-Achse geschehen ist (**pos=0**).

8.2.4.2.2 Linienzüge ergänzen mit der Funktion `lines()`

Mit der Funktion `lines()` ergänzt man zusätzliche Linienzüge, wobei die Daten (Argumente `x` und `y`) und der Linientyp (Argument `lty`) genauso konfiguriert werden wie bei der `plot()` - Funktion. Als Plot-Typ (Argument `type`) ist bei der `lines()` - Funktion der meist passende Wert `"l"` voreingestellt. Die folgenden Anweisungen

```
> x <- 1:5; y <- c(2, 4, 5, 4.5, 4); y2 <- c(3,3,4,4,3)
> plot(x, y, type="l", lwd=2)
> lines(x, y2, col="red", lty=2, lwd=2)
```

produzieren dieses Diagramm:



Bei einem Streudiagramm (siehe Abschnitt 8.2.5.2) verwendet man die `lines()` - Funktion z.B. dazu, eine lokal optimierte Modellprognose einzuzichnen:

```
> lines(lowess(daten$gewicht ~ daten$groesse), lty=2, col="blue")
```

8.2.4.2.3 Legenden bilden mit der Funktion `legend()`

Wenn ein Diagramm *mehrere* Linien oder mehrere Datengruppen (Symbole) enthält, ist eine Legende zur Erläuterung der Linien bzw. Gruppen empfehlenswert. Die dafür zuständige Low-Level - Grafikfunktion `legend()` kennt u.a. die folgenden generellen Argumente (unabhängig vom Grafiktyp):

- **x, y**
Die Koordinaten der linken, oberen Ecke des Legendenrahmens
Statt über ein (x, y) - Zahlenpaar kann man den Ort auch über ein Schlüsselwort aus der folgenden Liste festlegen: **"bottomright"**, **"bottom"**, **"bottomleft"**, **"left"**, **"topleft"**, **"top"**, **"topright"**, **"right"**, **"center"**.
- **title**
Eine Zeichenfolge mit dem Titel der Legende
- **legend**
Zeichenfolgen-Vektor mit den Etiketten zu den Linien bzw. Symbolen

Sollen Linien erläutert werden, die sich nach Typ, Stärke und/oder Farbe unterscheiden, sind folgende Argumente zu verwenden:

- **lty**
Ein Vektor mit den Linientypen
- **lwd**
Ein Vektor mit den Linienstärken
- **col**
Ein Vektor mit den Linienfarben

Sollen Datengruppen erläutert werden, die sich nach Symbol und/oder Farbe unterscheiden, sind folgende Argumente zu verwenden:

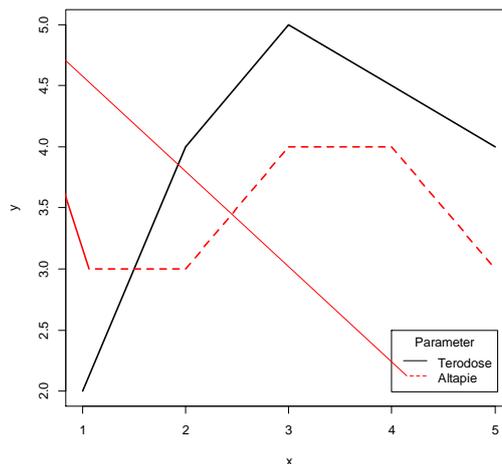
- **pch**
Ein Vektor mit den Symbolen
- **col**
Ein Vektor mit den Linienfarben

Um Datengruppen *und* Linien gemeinsam in einer Legende zu erläutern, können z.B. die Argumente **lty**, **pch** und **col** gemeinsam verwendet werden.

Das in Abschnitt 8.2.4.2.2 erstellte Liniendiagramm erhält durch den folgenden Funktionsaufruf

```
> legend(4, 2.5, title="Parameter", legend=c("Terodose", "Altapie"),
+ lty=c(1,2), col=c("black", "red"))
```

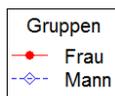
eine Legende:



In Abschnitt 8.2.5.2.3 wird für ein gruppiertes Streudiagramm eine Legende zur Erläuterung von Symbolen erstellt:



Nach Erweiterung dieses Streudiagramms um gruppenspezifische Regressionsgeraden werden in der Legende zu jeder Gruppe ein Symbol *und* eine Linie angezeigt:



Nach dem Kommando

```
> ?legend
```

liefert die **R**-Hilfe zahlreiche weitere Details zur Legendenbildung.

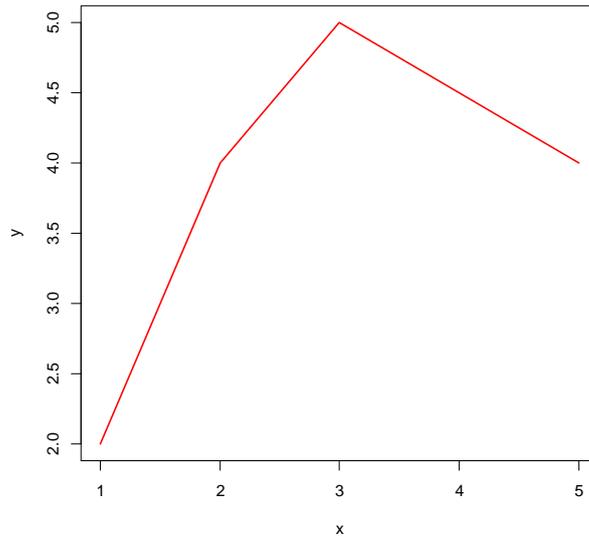
8.2.5 Wichtige Diagrammtypen

8.2.5.1 Liniendiagramm

Wenn sich bei einem `plot()` - Aufruf alle Punkte auf einer Linie befinden, und das `type`-Argument einen passenden Wert erhält (z.B.: "l"),

```
> x <- 1:5; y <- c(2, 4, 5, 4.5, 4)
> plot(x, y, type="l", col="red", lwd=2)
```

dann resultiert ein Liniendiagramm, z.B.:



8.2.5.2 Streudiagramm

In diesem Abschnitt nutzen wir die Funktion `plot()` dazu, Streudiagramme zu erstellen. Zunächst entstehen Varianten, die bei weitgehend äquivalentem Ergebnis auch mit SPSS möglich sind (bis Abschnitt 8.2.5.2.4). Schließlich werden die erlernten Techniken aber auch dazu verwendet, spezielle Streudiagramme zu produzieren, die mit SPSS nicht (in derselben Qualität) zu erstellen sind.

Wir verwenden anschließend die in Abschnitt 7.2.1 beschriebenen Beispieldaten, die z.B. aus der **R**-Binärdatei `ggg.RData` (zu finden an der im Vorwort vereinbarten Stelle, Unterordner **Datendateien mit R lesen**) mit der folgenden Anweisung in die Datentabelle `ggg` eingelesen werden können:

```
> load("ggg.RData")
> ggg
  geschlecht  gröÙe  gewicht
1      Mann    186    82.0
2      Mann    178    72.0
3      Mann    182    75.5
4      Frau    160    65.0
5      Frau    168    66.0
6      Frau     NA    76.0
7      Frau    165    55.0
8      Mann    179    76.5
9      Frau    158    50.5
10     Mann    175    80.0
11     Frau    176    62.0
13     Mann    176     NA
```

8.2.5.2.1 Einfaches Streudiagramm

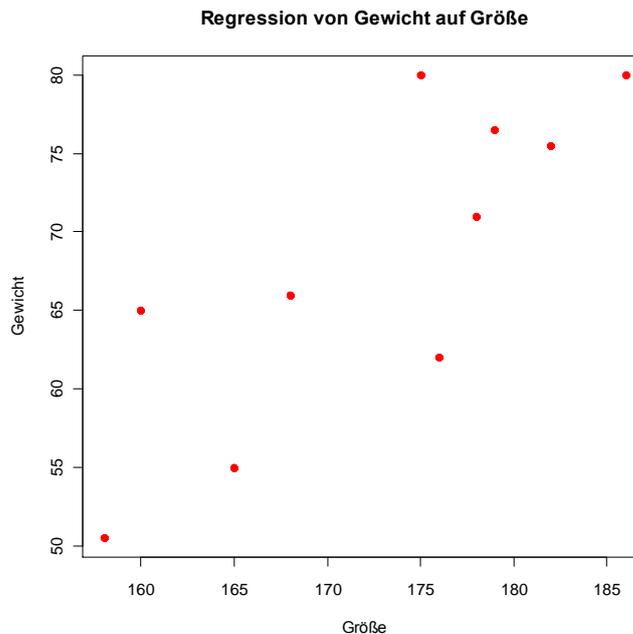
Die in Abschnitt 8.2.4 beschriebene **plot()** - Funktion erstellt beim voreingestellten Wert **p** für das Argument **type** ein zweidimensionales Streudiagramm für Punkte mit Koordinaten in den Vektoren **x** und **y**, die als erstes bzw. zweites Argument zu übergeben sind. In folgenden **plot()** - Aufruf

```
> plot(ggg$größe, ggg$gewicht, main="Regression von Gewicht auf Größe",
+ xlab="Größe", ylab="Gewicht", pch=19, col="red")
```

werden zur Gestaltung eines Streudiagramms einige Argumente verwendet, die in den Abschnitten 8.2.2 (über Grafikparameter), 8.2.3 (über Beschriftungen) und 8.2.4 (über die **plot()** - Funktion) vorgestellt worden sind:

- **main**
Hauptüberschrift
- **xlab, ylab**
Achsenbeschriftungen
- **pch**
Markierungsart für die Datenpunkte
- **col**
Zeichenfarbe

Das Ergebnis:



Die **plot()** - Funktion erlaubt zur Spezifikation der beiden Variablen auch die Modellsyntax (vgl. Abschnitt 9.2), wobei auf das Kriterium (die Y-Achsen - Variable) nach einer Tilde (~) der Regressor folgt, z.B.:

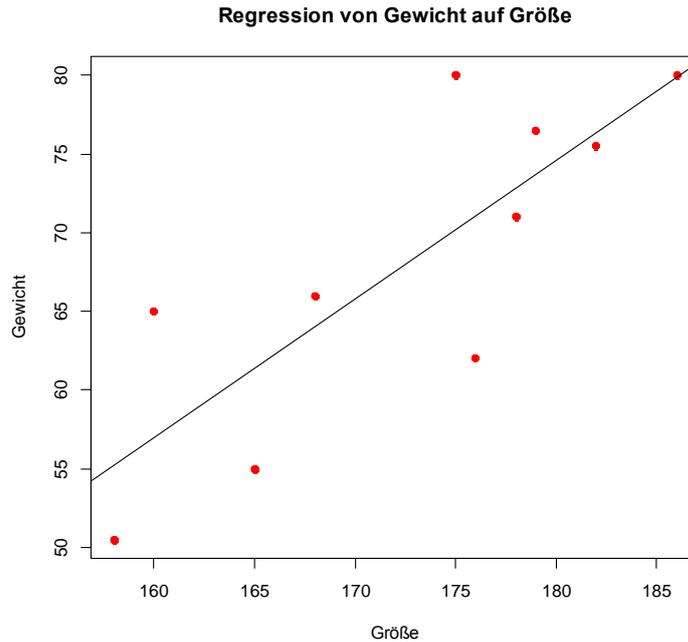
```
> plot(ggg$gewicht ~ ggg$größe, main="Regression von Gewicht auf Größe",
+ xlab="Größe", ylab="Gewicht", pch=19, col="red")
```

8.2.5.2.2 Regressionsgerade und sonstige Linien ergänzen

Um eine Regressionsgerade in das Streudiagramm einzeichnen zu können, lassen wir die lineare Regression von Gewicht auf Größe durch die Funktion **lm()** schätzen (zur Modellformulierung siehe Abschnitt 9.2). Die Low-Level - Grafikfunktion **abline()**

```
> abline(lm(gewicht ~ größe, data=ggg))
```

zeichnet unter Verwendung der `lm()` - Ergebnisse die Regressionsgerade:



Zur Anzeige einer nichtlinearen (z.B. quadratischen) Anpassungslinie eignen sich die Funktionen `lines()` (vgl. Abschnitt 8.2.4) und `predict()`, z.B.:

```
> lines(ggg$größe, predict(lm(gewicht~größe+I(größe^2), data=ggg)))
```

Als Y-Koordinaten werden die von `predict()` gelieferten Prognosen eines linearen Modells verwendet, das einen (im konkreten Beispiel eigentlich überflüssigen) quadratischen Term enthält. Im konkreten Fall klappt der `lines()` - Aufruf allerdings *nicht*, weil die `predict()` - Rückgabe und der Vektor `ggg$größe` wegen fehlender Werte eine verschiedene Länge haben.¹ Um zu gültigen Argumenten für die Funktion `lines()` zu kommen, wird ein neuer Vektor mit Werten im Bereich des Regressors erzeugt:

```
> dfx <- data.frame(größe = seq(155, 190, 1))
```

Auf den Vektor `dfx$größe` (in der Datentabelle `dfx`) wird das Modell aus den realen Daten angewendet, und die `lines()` - Funktion erhält schließlich 2 gleichlange Vektoren:

```
> qmodGG <- lm(gewicht~größe+I(größe^2), data=ggg)
> lines(dfx$größe, predict(qmodGG, newdata=dfx))
```

Mit Hilfe der Funktionen `lines()` und `lowess()` lässt sich eine lokal optimierte Modellprognose einzeichnen. Im Beispiel steigt der technische Aufwand etwas, weil die Vektoren `ggg$größe` und `ggg$gewicht` jeweils einen `NA`-Wert enthalten. Alle unvollständigen Fälle müssen entfernt werden, was mit einer neuen Datentabelle und der `na.omit()` - Funktion (siehe Abschnitt 5.3.5) gelingt:

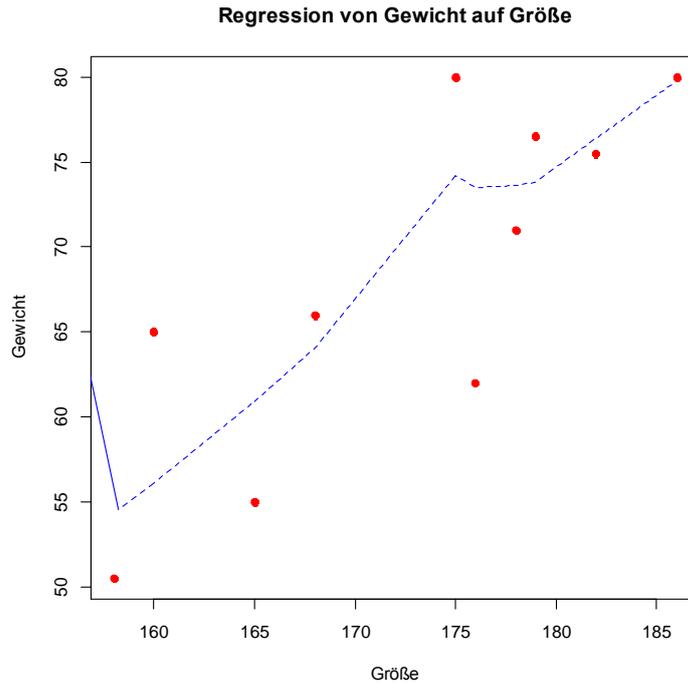
```
> datKom <- na.omit(ggg)
```

Die `lowess()` - Rückgabe ist eine Liste mit `x`- und `y`-Vektor, so dass sie als alleiniges Positionierungsargument für die `lines()` - Funktion taugt:

```
> lines(lowess(datKom$gewicht ~ datKom$größe), lty=2, col="blue")
```

Das Ergebnis:

¹ Der Vektor `ggg$größe` enthält 12 Elemente (das Element mit dem Wert `NA` zählt mit), die `predict()` - Rückgabe enthält hingegen nur 10 Elemente, weil bei jeweils einem Fall die Größen- bzw. Gewichtsausprägung fehlte.



Mit der Low-Level - Grafikfunktion **abline()** kann man auch eine durch die Argumente **a** und **b** definierte Linie zeichnen, z.B. eine „normative“ Regressionsgerade zur Idealgewichtsempfehlung „Größe - 100 - 10%“, die in mathematischer Formulierung zu folgenden Koeffizienten führt:¹

```
> abline(a=-90, b=0.9, col="blue")
```

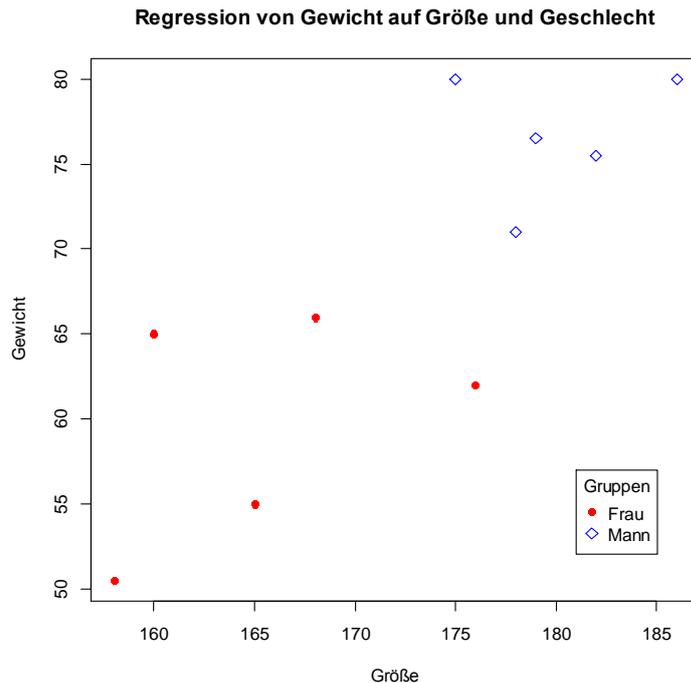
Für das Ergänzen von horizontalen oder vertikalen Linien ist die Funktion **abline()** ebenfalls zuständig. Den Argumenten **h** bzw. **v** ist als Wert die Achsenposition zu übergeben, z.B.:

```
> abline(h=68)
```

8.2.5.2.3 Gruppiertes Streudiagramm

In der geladenen Datentabelle **ggg** ist auch der Faktor **geschlecht** vorhanden, und nun soll das folgende gruppierte Streudiagramm erstellt werden, das die Geschlechtszugehörigkeit anhand der Symbole erkennen lässt:

¹ „Größe - 100 - 10%“ = $(\text{Größe} - 100) * 0,9 = -90 + 0,9 * \text{Größe}$



Über die in Abschnitt 7.3.1.2 beschriebene Rekodierung erstellen wir den Vektor `syms`, der für jeden Fall eine Symbolnummer enthält. Wir verwenden für Frauen das Symbol Nummer 16 (●) und für Männer das Symbol Nummer 5 (◇):

```
> syms <- numeric(length(ggg$geschlecht))
> syms[ggg$geschlecht == "Frau"] <- 16
> syms[ggg$geschlecht == "Mann"] <- 5
```

Analog gelangen wir zum **character**-Vektor `cols` mit der Farbe Rot für die Frauen und der Farbe Blau für die Männer:

```
> cols <- character(length(ggg$geschlecht))
> cols[ggg$geschlecht == "Frau"] <- "red"
> cols[ggg$geschlecht == "Mann"] <- "blue"
```

Mit den Vektoren `syms` bzw. `cols` als Werten für die Argumente `pch` bzw. `col` erstellen wir das Diagramm neu:

```
> plot(ggg$größe, ggg$gewicht,
+ main="Regression von Gewicht auf Größe und Geschlecht",
+ xlab="Größe", ylab="Gewicht", pch=syms, col=cols)
```

Schließlich erstellen wir noch eine Legende mit Hilfe der Low-Level - Grafikfunktion **legend()**:

```
> legend(181, 57, title="Gruppen", legend=c("Frau", "Mann"),
+ pch=c(16, 5), col=c("red", "blue"))
```

Wenn man damit einverstanden ist, die Symbole mit den Nummern 1 bzw. 2 für Frauen bzw. Männer zu verwenden (= interne Werte für die `geschlecht`-Faktorstufen) und auf Farben keinen Wert legt, dann lässt sich der Aufwand für die geschlechtsspezifischen Markierungen reduzieren:

```
> plot(ggg$größe, ggg$gewicht,
+ main="Regression von Gewicht auf Größe und Geschlecht", xlab="Größe",
+ ylab="Gewicht", pch=as.numeric(ggg$geschlecht))
```

Man erstellt aus dem `geschlecht`-Faktor mit der Funktion **as.numeric()** einen numerischen Vektor und verwendet diesen als Wert für das **plot()** - Argument **pch**.

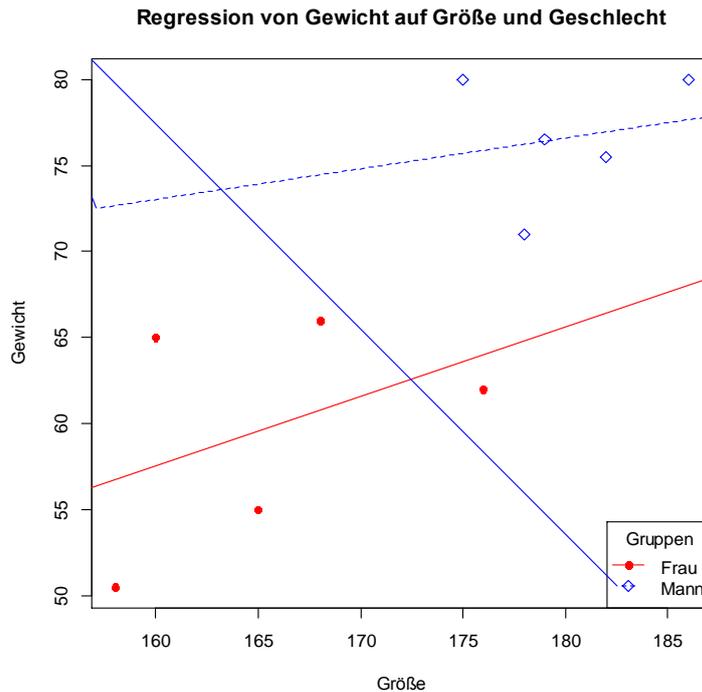
Um geschlechtsbedingte Regressionsgeraden einzuzeichnen, erweitert man den in Abschnitt 8.2.5.2.2 vorgestellten `abline()` - Funktionsaufruf um eine Fallauswahl über Indexvektoren:

```
> abline(lm(gewicht[geschlecht=="Frau"]~größe[geschlecht=="Frau"],data=ggg),lty=1,col="red")
> abline(lm(gewicht[geschlecht=="Mann"]~größe[geschlecht=="Mann"],data=ggg),lty=2,col="blue")
```

Mit einer um die Linientypen erweiterten Legende in der rechten unteren Ecke

```
> legend("bottomright", title="Gruppen", legend=c("Frau", "Mann"),
+       pch=c(16, 5), lty=c(1, 2), col=c("red", "blue"))
```

sieht das gruppierte Streudiagramm jetzt so aus:



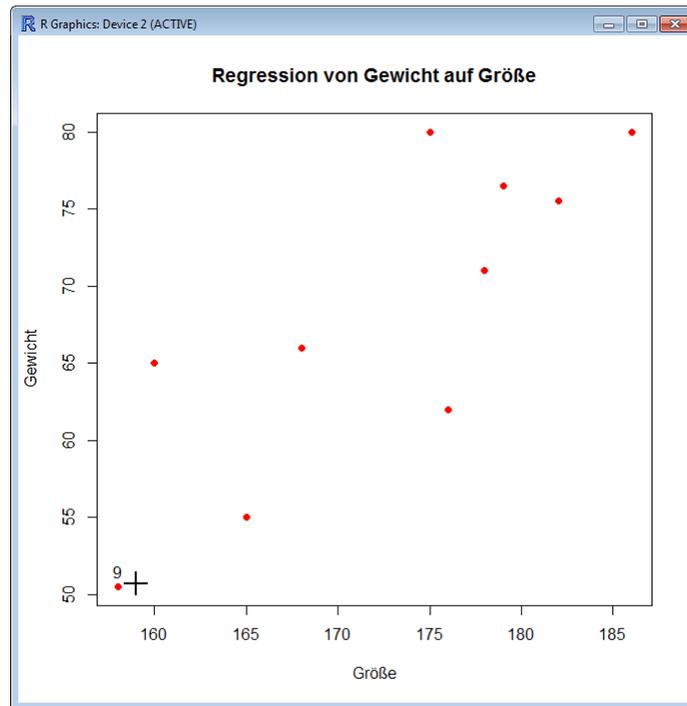
Das klassische Grafiksystems bietet noch weitere Optionen für Streudiagramme (z.B. Konfidenzintervalle, Streudiagramm-Matrizen), über die z.B. Muenchen (2011, S. 480ff) informiert.

8.2.5.2.4 Fallidentifikation

Befindet sich im aktiven Grafikfenster ein Streudiagramm mit den Variablen `ggg$größe` und `ggg$gewicht`, dann kann man nach dem folgenden `identify()` - Funktionsaufruf

```
> identify(ggg$größe, ggg$gewicht)
```

die Datenpunkte im Grafikfenster per Mausklick beschriften, z.B.:

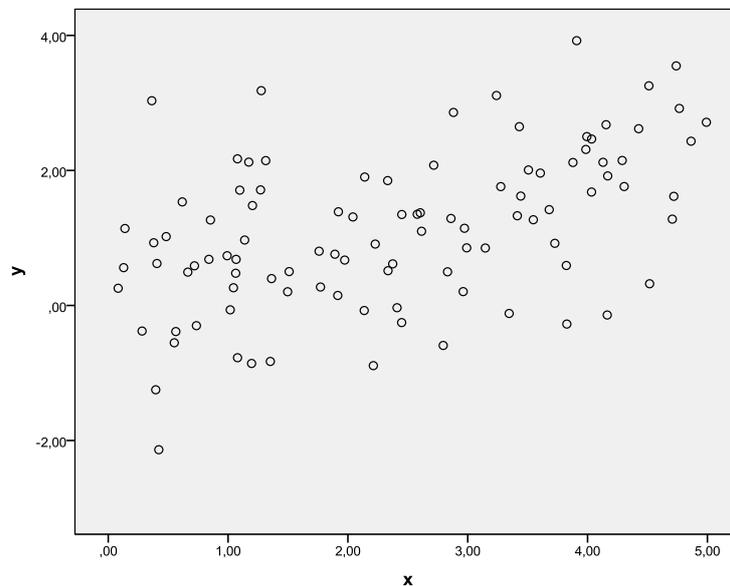


Um den Fallidentifikationsmodus zu beenden, ...

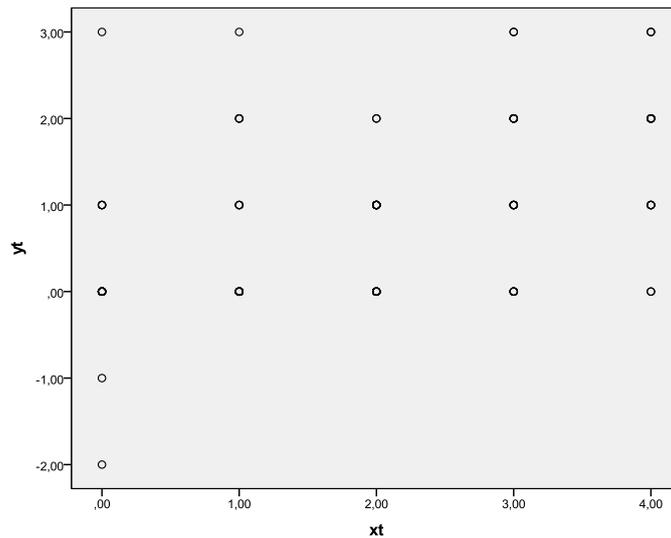
- drückt man die **Esc**-Taste
- oder wählt das Item **Stopp** aus dem Kontextmenü zum Grafikfenster.

8.2.5.2.5 Jitter-Darstellung

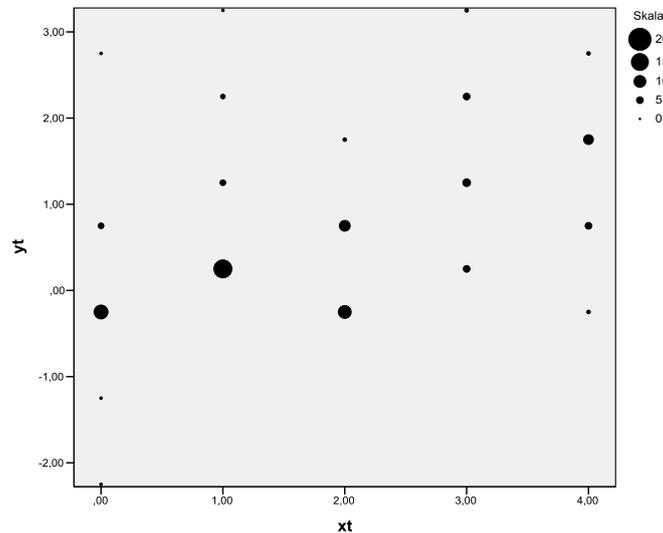
Das folgende (mit SPSS erstellte) Streudiagramm zeigt zwei deutlich miteinander korrelierte metrische Variablen:



Sind statt der metrischen Variablen nur vergrößerte Messungen vorhanden (mit wenigen Ausprägungen, vgl. Abschnitt 3.3), dann ergeben sich im Streudiagramm mehrfach besetzte Punkte, und die Interpretation ist erschwert, z.B.



Der Zusammenhang wird deutlicher, wenn die Größe der Punkte von der Häufigkeit abhängt:

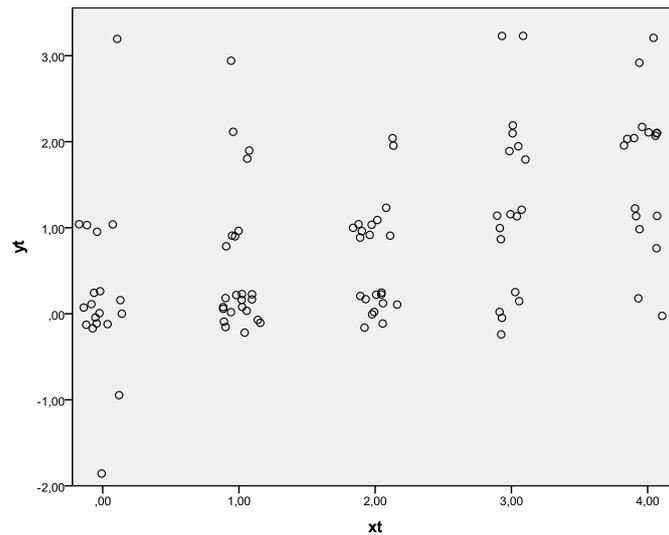


Bei der Jitter-Technik zur Lösung des Problems addiert man kleine Zufallsschwankungen zu den Beobachtungswerten, um die Punkte auseinander zu ziehen. In SPSS ist diese Technik mit Hilfe der GPL-Syntax realisierbar, doch lässt sich der Jitter-Grad nicht einstellen, so dass im Beispiel kein überzeugendes Ergebnis resultiert:¹

¹ Verwendete Syntax:

```

GGRAPH
  /GRAPHDATASET NAME="graphdataset" VARIABLES=xt yt MISSING=LISTWISE REPORTMISSING=NO
  /GRAPHSPEC SOURCE=INLINE.
BEGIN GPL
  SOURCE: s=userSource(id("graphdataset"))
  DATA: xt=col(source(s), name("xt"))
  DATA: yt=col(source(s), name("yt"))
  GUIDE: axis(dim(1), label("xt"))
  GUIDE: axis(dim(2), label("yt"))
  ELEMENT: point.jitter(position(xt*yt))
END GPL.
    
```

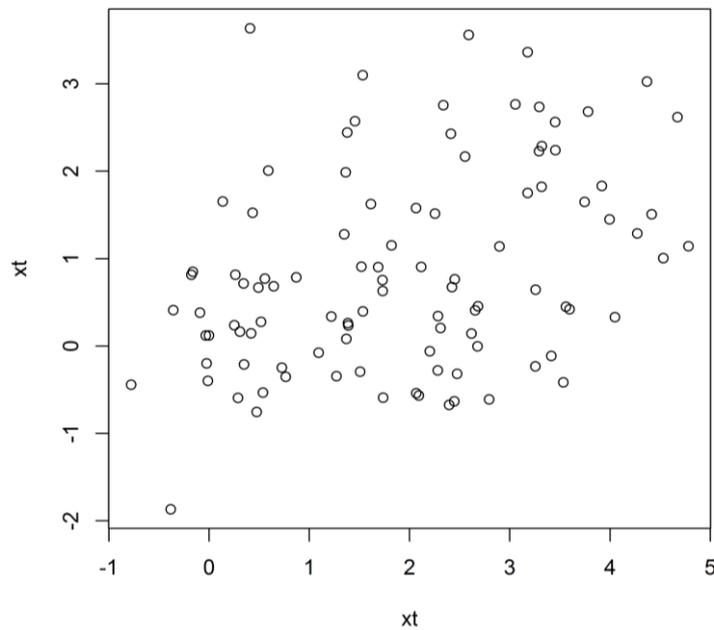


In der **R**-Standardgrafik erlaubt die **jitter()** - Funktion einen einstellbaren „Verwacklungsgrad“,

```
BEGIN PROGRAM R.
png("u:/eigene dateien/r/jitter.png", 15, 15, units="cm", res=600)
data <- spssdata.GetDataFromSPSS()
plot(jitter(data$xt,4),jitter(data$yt,4),main="Streudiagramm mit Jitter",xlab="xt",ylab="yt")
END PROGRAM.
```

und die **plot()** - Funktion liefert mit den so behandelten Daten ein besser interpretierbares Bild:¹

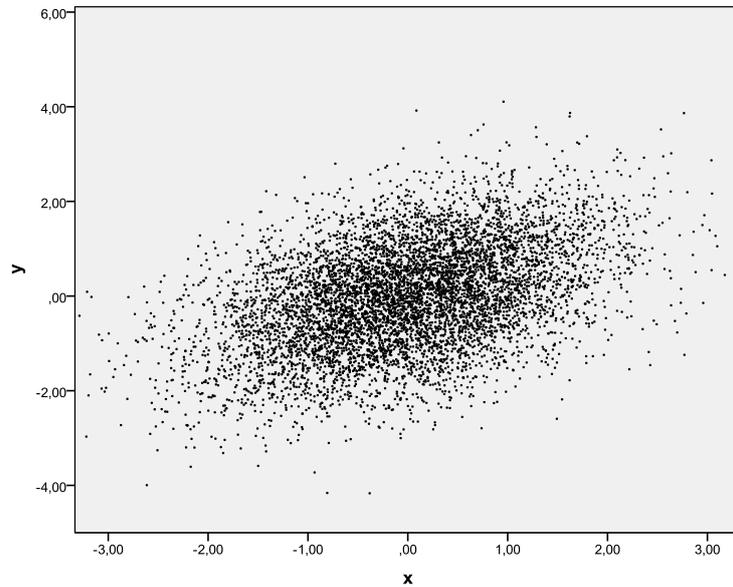
Streudiagramm mit Jitter



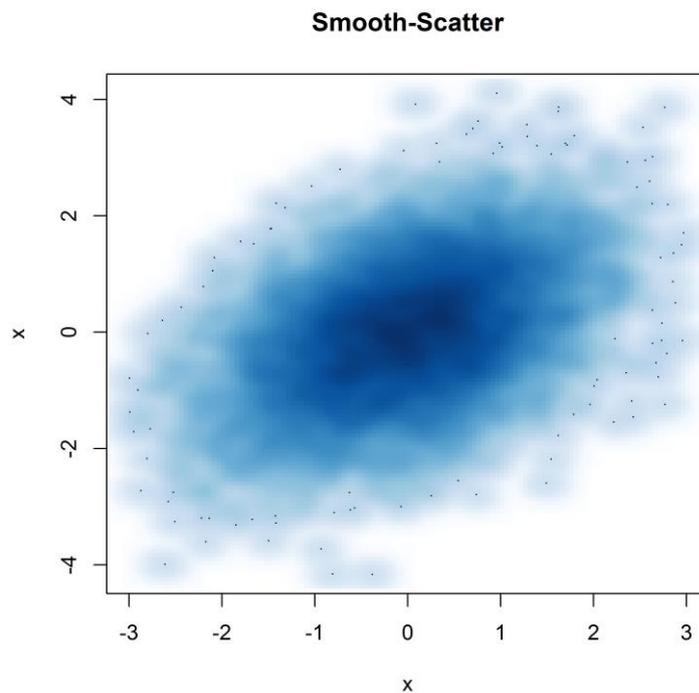
¹ Die zur Demonstration verwendeten Daten stammen aus der SPSS-Datendatei **Streudiagramm mit Jitter.sav**, die sich an der im Vorwort vereinbarten Stelle im Unterordner **R-Grafik\Traditionelle Grafik\Streudiagramme** befindet.

8.2.5.2.6 Farbliche Dichtedarstellung bei großen Stichproben

Sind sehr viele Datenpunkte vorhanden, ist ein Streudiagramm auch bei Wahl von winzigen Markierungen nicht ideal, z.B.:



Die **R**-Standardgrafik bietet für diesen Fall mit der Funktion **smoothScatter()** eine zumindest ästhetisch interessante Alternative, wobei die Wahrscheinlichkeitsdichte durch die Farbintensität ausgedrückt wird, z.B.:



Das Beispiel wurde über ein SPSS-Syntaxfenster angefordert:¹

¹ Die zur Demonstration verwendeten Daten stammen aus der SPSS-Datendatei **SmoothScatter.sav**, die sich an der im Vorwort vereinbarten Stelle im Unterordner **R-Grafik\Traditionelle Grafik\Streudiagramme** befindet.

```
BEGIN PROGRAM R.
png("u:/eigene dateien/r/smooth.png", 15, 15, units="cm", res=600)
data <- spssdata.GetDataFromSPSS()
smoothScatter(data$x,data$y,xlim=c(-3,3),main="Smooth-Scatter",xlab="x",ylab="x")
END PROGRAM.
```

8.2.5.3 Boxplot

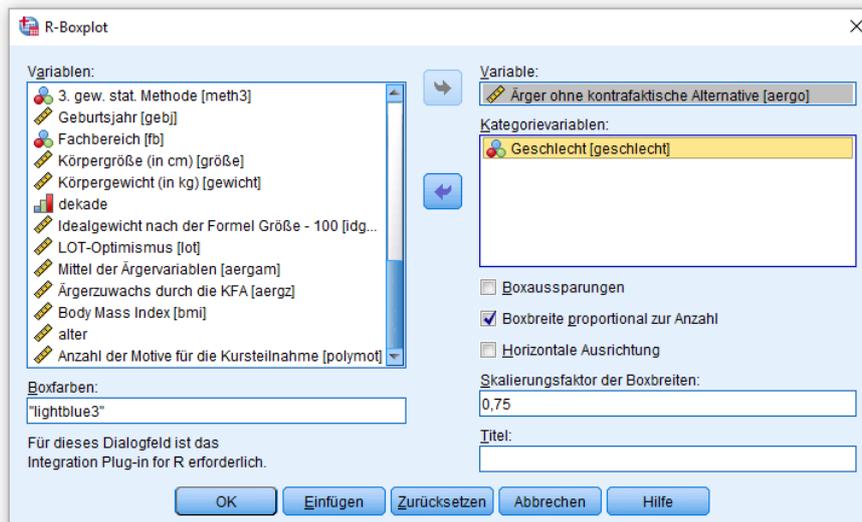
Zum Erstellen von Boxplots bieten **R** und SPSS weitgehend äquivalente Möglichkeiten:

- Darstellung einer einzelnen Variablen
- Mehrere Variablen nebeneinander
- Mehrere Gruppen nebeneinander (ein- oder zweifaktorielles Design)

Installiert man in SPSS 24 gemäß Abschnitt 2.2.3 das Erweiterungsbundle **Rboxplot**, dann stehen anschließend die wichtigsten Boxplot-Optionen der traditionellen **R**-Grafik in SPSS nach dem Menübefehl

Grafik > R-Boxplot

über die folgende Dialogbox zur Verfügung:¹

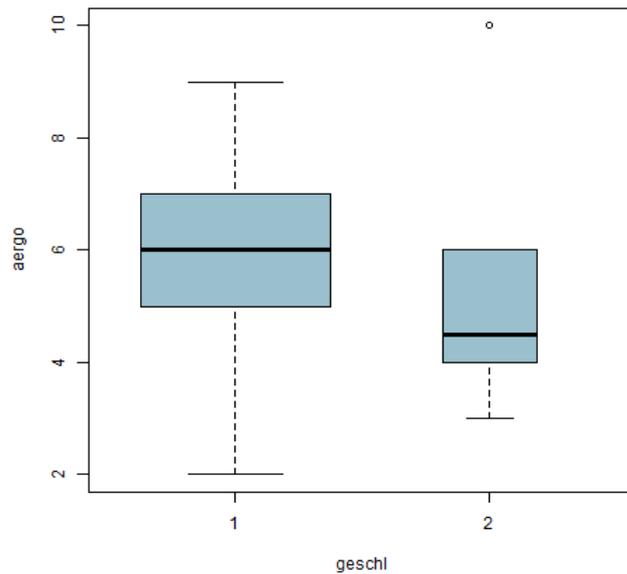


Hier bestehen u.a. die folgenden Einstellmöglichkeiten:

- **Boxbreite proportional zur Anzahl**
Die Breite lässt sich zur Anzeige der Teilstichprobengröße nutzen.
- **Skalierungsfaktor der Boxbreiten**
Damit lässt sich die generelle Breite der Boxen beeinflussen.
- **Boxfarben**
Hier ist ein **R**-Farbname einzutragen.

Im folgenden Beispiel wird deutlich, dass die beiden Teilstichproben unterschiedlich groß sind, wobei aber exakte Angaben fehlen:

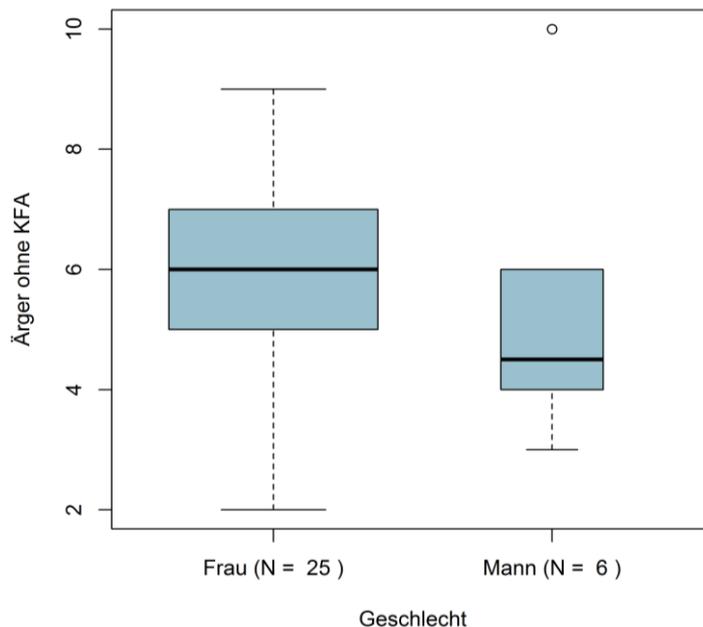
¹ Die Variablen stammen aus der Beispieldatei **kfa.sav**, die an der im Vorwort vereinbarten Stelle zu finden ist.



Bei direkter Nutzung der **R**-Funktion **boxplot()** über **R**-Anweisungen im SPSS-Syntaxfenster

```
BEGIN PROGRAM R.
png("u:/eigene dateien/r/box.png", 15, 15, units="cm", res=600)
data <- spssdata.GetDataFromSPSS(factorMode="labels")
levels(data$geschlecht)[1] <- paste(levels(data$geschl)[1], "(N = ",
length(data$geschlecht[data$geschlecht=="Frau"]), ")")
levels(data$geschlecht)[2] <- paste(levels(data$geschlecht)[2], "(N = ",
length(data$geschlecht[data$geschlecht=="Mann"]), ")")
boxplot(data$aergo ~ data$geschlecht,col="lightblue3", varwidth=TRUE, boxwex=0.75,
xlab="Geschlecht", ylab="Ärger ohne KFA")
END PROGRAM.
```

lässt sich mit Hilfe von Grafikparametern und **boxplot()** - Argumenten ein besseres Ergebnis erzielen:



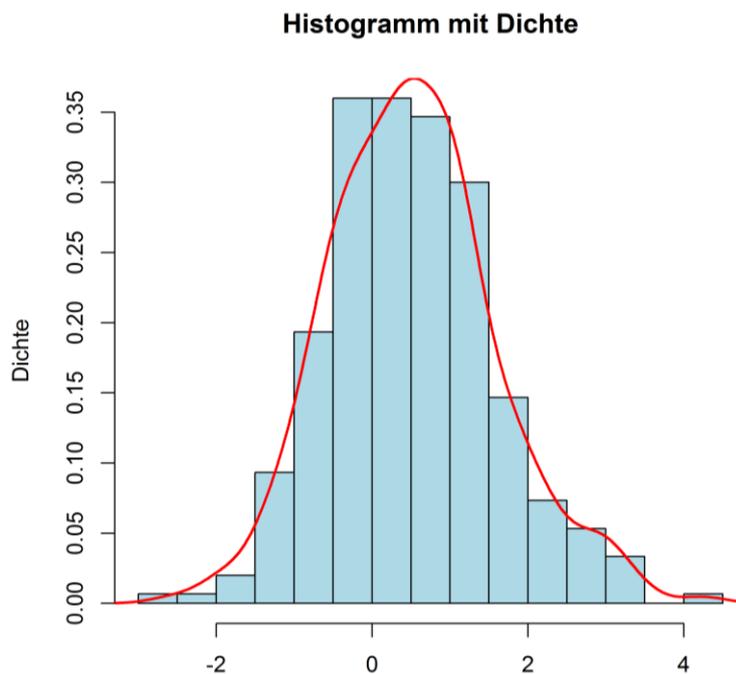
Zur Anzeige der Fallzahlen wurden per `levels()` - Funktion (siehe Abschnitt 5.3.4.3.1) die Etiketten zu den Ausprägungen des Faktors `data$geschl` modifiziert, wobei die neuen Zeichenfolgen mit Hilfe der `R`-Funktion `paste()` aus den alten und den berechneten Teilstichprobenumfängen entstanden sind.

Hier zeigt sich exemplarisch, dass die Grafikproduktion in `R` viele Optionen bietet, wenn man den Aufwand beim Erstellen der Syntax nicht scheut.

8.2.5.4 Histogramm mit Dichteschätzung

Mit den Grafikfunktionen `hist()`, `density()` und `lines()`

lässt sich (wie schon in Abschnitt 8.1.6 demonstriert) zu einer Variablen ein Histogramm mit geschätzter Dichtefunktion erstellen:



Diese Darstellung wurde für eine SPSS-Variablen angefordert (vgl. Abschnitt 4.1) und mit publikations-tauglicher Auflösung in eine PNG-Datei geschrieben: ¹

```
BEGIN PROGRAM R.
data <- spssdata.GetDataFromSPSS()
png("u:/eigene dateien/r/NormDens.png", 15, 15, units="cm", res=600)
hist(data$y, freq=FALSE, breaks=10, col="lightblue", main="Histogramm mit Dichte",
      ylab="Dichte", xlab=NA)
lines(density(data$y), col="red", lwd=2)
END PROGRAM.
```

Die `hist()` - Funktion gehört zu den High Level - Grafikfunktionen in `R`, die ein komplettes Diagramm produzieren. Im obigen `hist()` - Aufruf werden neben den generellen Grafikparametern `col`, `main`, `xlab` und `ylab` für die Zeichenfarbe bzw. für Beschriftungen (vgl. Abschnitt 8.2.2) noch folgende Argumente genutzt:

¹ Die zur Demonstration verwendeten Daten stammen aus der SPSS-Datendatei **Histogramm mit Dichte.sav**, die sich an der im Vorwort vereinbarten Stelle im Unterordner **R-Grafik\Traditionelle Grafik\Histogramm mit Dichte** befindet.

- **freq=FALSE**
Statt absoluter Häufigkeiten sollen *relative* zur Beschriftung der Y-Achsen-Teilstriche verwendet werden.
- **breaks**
Mit diesem Parameter beeinflusst man die Anzahl der Intervalle auf der X-Achse, wobei zu kleine oder zu große Werte zu einem wenig aussagekräftigen Histogramm führen.

Der Low Level - Grafikfunktion **lines()** - Funktion wird im Beispiel als erstes Argument die Funktion **density()** zur Definition des Linienverlaufs zu übergeben. Mit dem Grafikparameter **lwd** wird die Linienstärke beeinflusst.

8.2.5.5 Mehrere Diagramme kombinieren

Sollen zu Vergleichszwecken mehrere Diagramme in *einer* Abbildung kombiniert werden, definiert man mit dem Grafikparameter **mfrow** eine Platzierungsmatrix, deren Zellen durch die anschließenden Grafikfunktionsaufrufe gefüllt werden. Sollen z.B. zwei Diagramme übereinander erscheinen, verwendet man den folgenden Aufruf der **par()** - Funktion (vgl. Abschnitt 8.2.2):

```
> par(mfrow=c(2,1))
```

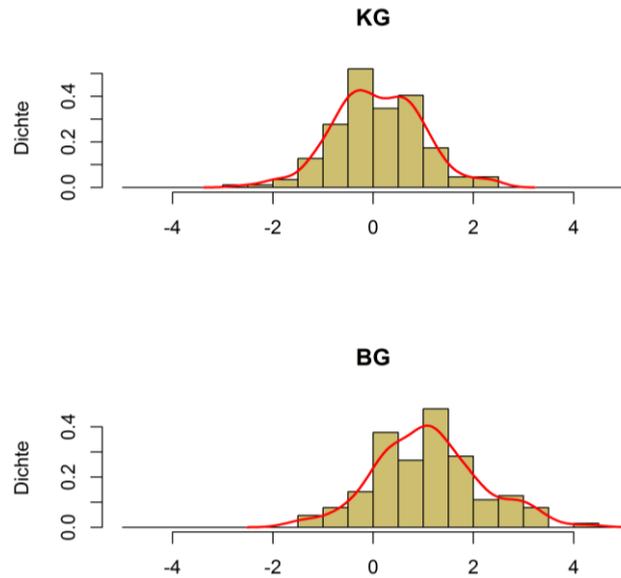
Zur Demonstration verwenden wir die künstlichen Daten aus Abschnitt 8.2.5.4 und fügen in SPSS noch die Variable **treat** zur Gruppeneinteilung hinzu, die einen Effekt auf die abhängige Variable **y** ausübt:

```
do if uniform(1) < 0.5.
  compute treat = 0.
else.
  compute treat = 1.
  compute y = y + 1.
end if.
```

Mit den folgenden, aus einem SPSS-Syntaxfenster abgeschickten **R**-Anweisungen werden zwei übereinander stehende Histogramme mit empirischer Dichte erstellt:

```
BEGIN PROGRAM R.
data <- spssdata.GetDataFromSPSS()
g0 <- data$y[data$treat==0]
g1 <- data$y[data$treat==1]
png("u:/eigene dateien/r/NormDens2.png", 15, 15, units="cm", res=600)
par(mfrow=c(2,1))
hist(g0, freq=FALSE, breaks=seq(-5,5,0.5), col="lightgoldenrod3", main="KG", xlab="", ylab="Dichte")
lines(density(g0), col="red", lwd=2)
hist(g1, freq=FALSE, breaks=seq(-5,5,0.5), col="lightgoldenrod3", main="BG", xlab="", ylab="Dichte")
lines(density(g1), col="red", lwd=2)
par(mfrow=c(1,1))
END PROGRAM.
```

Wegen der identischen Klassendefinitionen (Argument **breaks**) wird der „Behandlungseffekt“ gut sichtbar:

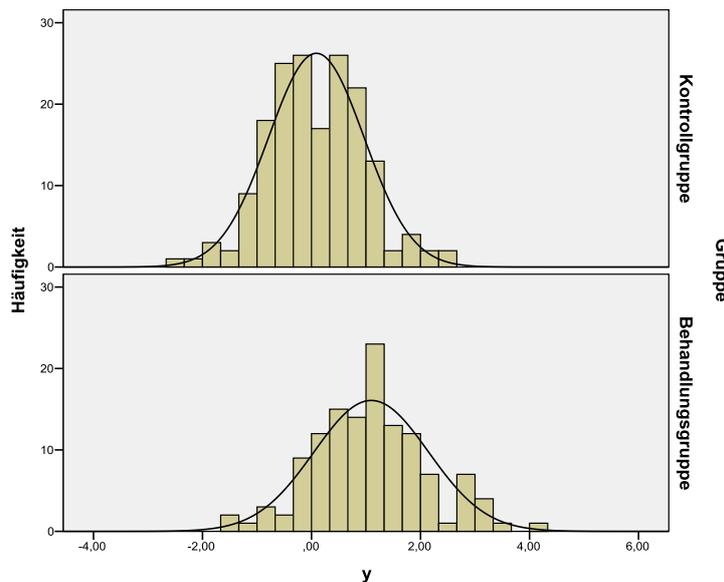


Am Ende wird der Grafikparameter **mfrow** wieder auf seinen Standardwert **c(1,1)** zurückgesetzt.

In SPSS lässt sich mit dem Kommando

```
GRAPH
  /HISTOGRAM(NORMAL)=y
  /PANEL ROWVAR=treat ROWOP=CROSS.
```

ein analoges Diagramm erstellen, wobei die Dichteschätzung eine Normalverteilung voraussetzt:

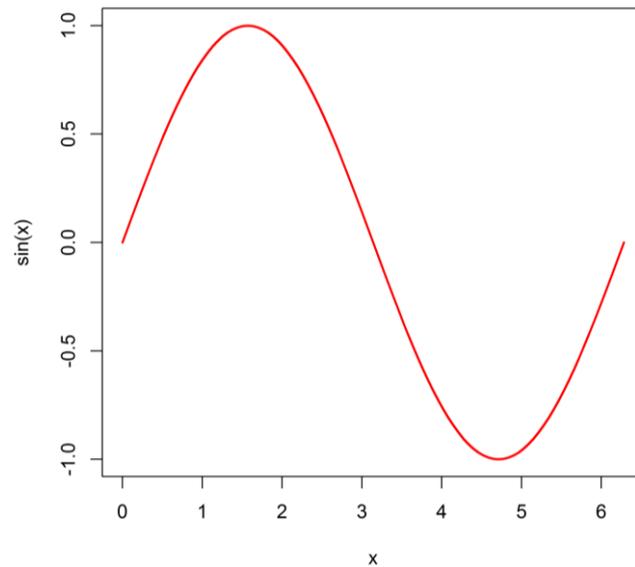


8.2.5.6 Funktionsplots

Mit der High Level - Grafikfunktion **curve()** kann man den Graphen einer Funktion von einer Veränderlichen plotten, z.B. den Sinus im Intervall $[0, 2\pi]$:

```
> curve(sin(x), 0, 2*pi, col="red", lwd=2)
```

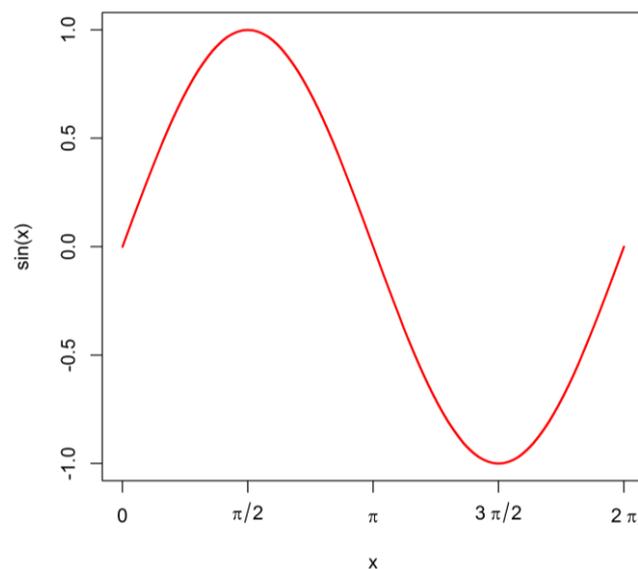
Beim Ergebnis des ersten Versuchs stört, dass die X-Achsenbeschriftung keinen Bezug zu den kritischen Punkten im Sinus-Verlauf hat:



Daher schalten wir im `curve()` - Aufruf die X-Achse aus und erzeugen über die Low-Level - Grafikfunktion `axis()` (siehe Abschnitt 8.2.4.2) eine Achse mit Markierungen an den passenden Stellen ($0, \pi/2, \pi, 3\pi/2, 2\pi$):

```
> curve(sin(x), 0, 2*pi, col="red", lwd=2, xaxt="n")
> axis(side=1, at=c(0, pi/2, pi, 3*pi/2, 2*pi),
+ labels=c("0", expression(pi/2), expression(pi), expression(3~pi/2),
+ expression(2~pi)))
```

Das `labels`-Argument erhält einen `character`-Vektor mit den gewünschten Teilstrichbeschriftungen, wobei durch `expression()` - Aufrufe für die mathematische Typographie gesorgt wird:



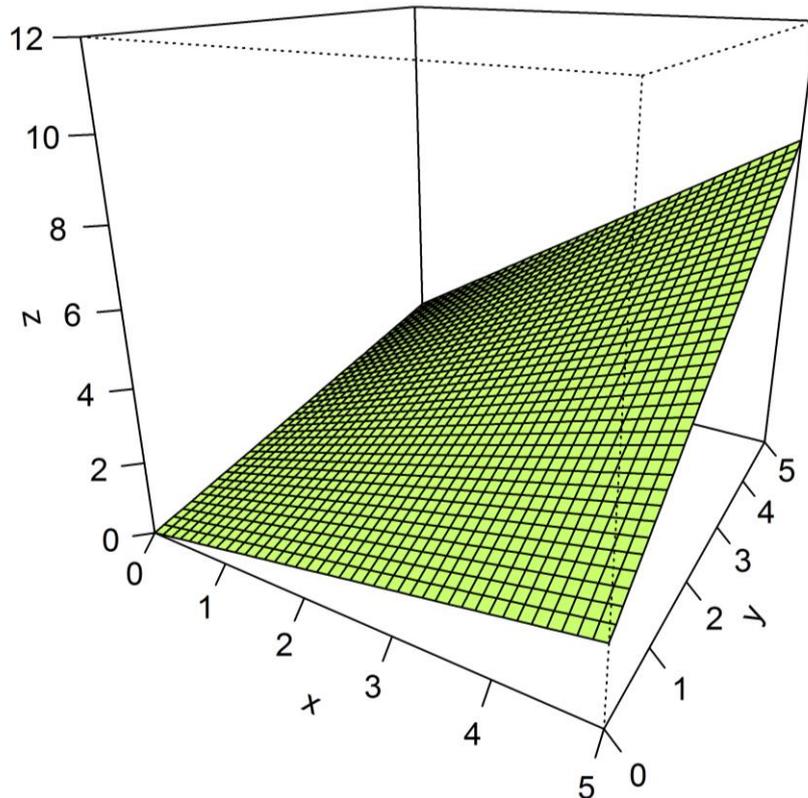
Als Beispiel für einen 3D-Plot soll die Reaktionsoberfläche einer multiplen Regression mit zwei interagierenden metrischen Regressoren (siehe Baltés-Götz 2017) mit der Funktion `persp()` gezeichnet werden:

```

> png("u:/eigene dateien/r/ModReg.png", 15, 15, units="cm", res=600)
> x <- seq(0, 5, length=40)
> y <- x
> f <- function(x, y) {
>   0.4*x + 0.4*y + 0.2*x*y
> }
> z <- outer(x, y, f)
> persp(x, y, z, zlim=range(0, 12), theta=30, phi=20, col="darkolivegreen1",
+       ticktype="detailed")
> dev.off()

```

Weil für das PNG-Ausgabegerät eine hohe Auflösung gewählt wurde, ist das Ergebnis recht ansehnlich:



8.3 Das Grafikpaket ggplot2

Das von Hadley Wickham erstellt und in einer Monographie beschriebene Paket **ggplot2** (siehe Wickham 2009, 2016) basiert auf der von Wilkinson (2005) entwickelten *Grammar of Graphics* (daher die beiden Anfangsbuchstaben *gg* im Paketnamen). Seine anspruchsvolle, aber gut durchdachte Konstruktion macht die Anweisungen zur Definition von Diagrammen relativ übersichtlich. Im Ergebnis ist im Vergleich zur traditionellen **R**-Grafik bei einfachen Diagrammen der Aufwand etwas höher, bei komplexen Diagrammen aber deutlich geringer.

Weil die GG-Grafik nicht mit der traditionellen Technik arbeitet, sind die in Abschnitt 8.2.2 beschriebenen Grafikparameter hier bedeutungslos.

Um die Funktionen in **ggplot2** nutzen zu können, müssen Sie das Paket einmalig installieren:

```
> install.packages("ggplot2")
```

Neben dem Paket **ggplot2** werden in der Regel etliche weitere, von **ggplot2** vorausgesetzte Pakete installiert. Vor jeder Verwendung muss **ggplot2** geladen werden:

```
> library(ggplot2)
```

Das Paket **ggplot2** bietet zum Erstellen eines Diagramms zwei Optionen:

- Den einfachen Weg über die Funktion **qplot()**
Bei den Argumenten und Voreinstellungen dieser auf einfache Bedienbarkeit getrimmten Funktion (*Quick-Plot*) hat man sich an der traditionellen Graphik-Funktion **plot()** orientiert (vgl. Abschnitt 8.2.4).
- Den flexiblen Weg über die Funktion **ggplot()**
Dabei wird ein Diagramm inkrementell aus Schichten aufgebaut, was für große Flexibilität bei der Erstellung individueller Lösungen sorgt.

Es ist durchaus möglich, ein Diagramm-Objekt mit **qplot()** zu initialisieren und mit der Flexibilität von **ggplot()** weitere Schichten zu ergänzen. Im Manuskript wird bevorzugt die Funktion **ggplot()** verwendet.

Derzeit beschränkt sich das Paket **ggplot2** auf 2D-Darstellungen (Wickham 2016, Abschnitt 3.6), so dass sich z.B. keine Reaktionsoberflächen von Funktionen darstellen lassen (vgl. Abschnitt 8.2.5.6). Auf 3D-Effekte in genuin zweidimensionalen Diagrammen (z.B. in einem Balkendiagramm zur Verteilung einer nominalskalierten Variablen) wird bewusst verzichtet.

Zur vertieften und systematischen Einarbeitung in die Verwendung des **ggplot2** - Pakets eignen sich die folgenden Bücher:

- In Wickham (2016) liegt der Schwerpunkt auf der Logik der Grafikproduktion, wobei aber auch die Anwendungsmöglichkeiten demonstriert werden.
- Chang (2013) liefert mit seinem „Kochbuch“ zahlreiche Lösungen für typische Aufgaben.

Eine komplette technische Dokumentation ist hier zu finden:

<http://docs.ggplot2.org/current/>

Es ist auf eine aktuelle Dokumentation zu achten, weil aufgrund der dynamischen Weiterentwicklung des **ggplot2**-Pakets so manches Beispiel aus einem älteren Lehrbuch mit einer modernen Installation nicht mehr klappt.

In der **R**-Version 3.2.5 ist das **ggplot2**-Paket in der Version 2.2.0 enthalten.

8.3.1 Grammatik eines ggplot2-Plots

Wir verwenden anschließend meist die in Abschnitt 7.2.1 beschriebenen Beispieldaten, die z.B. aus der **R**-Binärdatei **ggg.RData** (zu finden an der im Vorwort vereinbarten Stelle, Unterordner **Datendateien mit R lesen**) mit der folgenden Anweisung in die Datentabelle **ggg** eingelesen werden können:

```
> load("ggg.RData")
```

8.3.1.1 Plot-Objekte, Schichten und Geome

Ein **ggplot2**-Plot besteht aus übereinander liegenden **Schichten**. Eine Schicht enthält eine Datenvisualisierung, die als **Geom** (*geometric object*) bezeichnet wird. Ein Geom kann z.B. (x,y) - Datenpunkte, Balken oder Linien enthalten.

Bevor der Schichtenaufbau losgehen kann, muss ein **Plot-Objekt** angelegt werden. Dazu rufen wir die Funktion **ggplot()** auf und benennen im ersten Argument eine Datentabelle mit den darzustellenden Variablen. Sind Variablen auf *mehreren* Diagrammschichten mit bestimmten Darstellungsaspekten (z.B. mit der X- oder Y-Position) verbunden, gibt man diese Verbindung schon bei der Plot-Initialisierung bekannt, wobei ein Aufruf der Funktion **aes()** zu verwenden ist (vgl. Abschnitt 8.3.1.2). Im folgenden Beispiel

werden die Variablen `größe` bzw. `gewicht` aus der Datentabelle `ggg` der X- bzw. Y-Position des Diagramms zugeordnet, weil sie im `aes()` - Aufruf als erstes bzw. zweites Argument auftreten:

```
> sd <- ggplot(ggg, aes(größe, gewicht))
```

Wird bei der Plot-Initialisierung darauf verzichtet, eine voreingestellte Datentabelle bzw. Rollenzuweisungen für Variablen zu benennen, müssen diesen Angaben später im Zusammenhang mit den Schichtdefinitionen nachgeholt werden.

Ein Plot-Objekt enthält keine Referenz auf die Datentabelle, sondern eine *Kopie*, so dass spätere Änderungen der Datentabelle ohne Effekt auf das Plot-Objekt bleiben.

Wie ein (impliziter) `print()` - Aufruf für das eben erzeugte Plot-Objekt zeigt,

```
> sd
```

ist mangels vorhandener Schichten noch keine Grafikproduktion möglich:

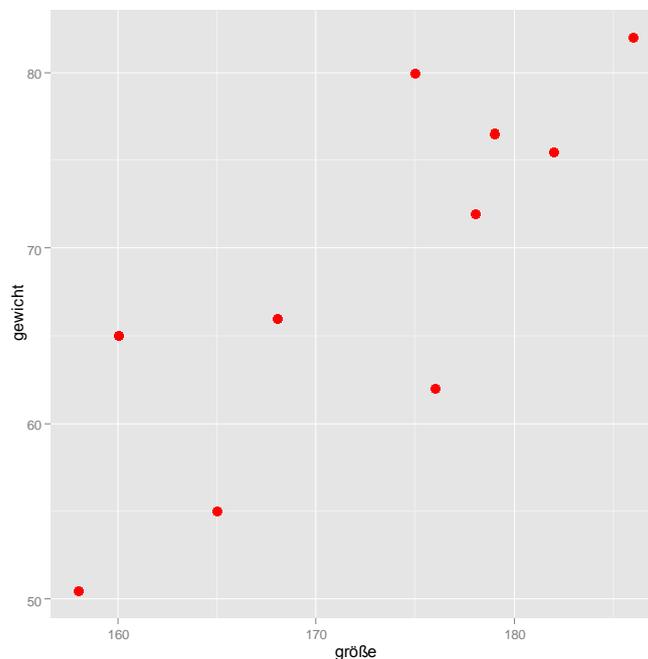
```
Fehler: No layers in plot
```

Um für ein Plot-Objekt eine Schicht mit einem bestimmten Geom zu erstellen, ruft man eine Funktion mit der gewünschten Darstellungsart im Namen auf, z.B. `geom_point()` für eine Schicht mit Datenpunkten, die (x,y) - Wertepaare darstellen. Die Rückgabe der Schicht-generierenden Funktion wird mit dem Operator `+` zum Plot-Objekt hinzugefügt, z.B.:

```
> sd + geom_point(colour="red", size=3)
```

Optional kann man mit ihren Argumenten z.B. schichtspezifische Abbildungen der ästhetischen Attribute des Geoms (siehe Abschnitt 8.3.1.2) auf Variablen oder feste Werte sowie eine schichtspezifische Datentabelle festlegen, wenn keine passenden Voreinstellungen auf Plot-Ebene bestehen. Im Beispiel werden die ästhetischen Attribute `colour` und `size` jeweils auf einen festen Wert abgebildet.

Die „Summe“ aus dem Plot-Objekt und der darstellbaren Schicht erscheint in einem Grafikkfenster:



Wie sich mit der Funktion `str()` nachweisen lässt, hat ein `ggplot2`-Objekt den Datentyp *Liste*, z.B.:

```
> str(sd)
List of 9
 $ data      :'data.frame':   12 obs. of  3 variables:
  ..$ geschlecht: Factor w/ 2 levels "Frau","Mann": 2 2 2 1 1 1 1 2 1 2 ...
  ..$ gröÙe     : int [1:12] 186 178 182 160 168 NA 165 179 158 175 ...
  ..$ gewicht   : num [1:12] 82 72 75.5 65 66 76 55 76.5 50.5 80 ...
  $ layers     : list()

. . .

$ plot_env   :<environment: R_GlobalEnv>
$ labels     :List of 2
  ..$ x: chr "gröÙe"
  ..$ y: chr "gewicht"
- attr(*, "class")= chr [1:2] "gg" "ggplot"
```

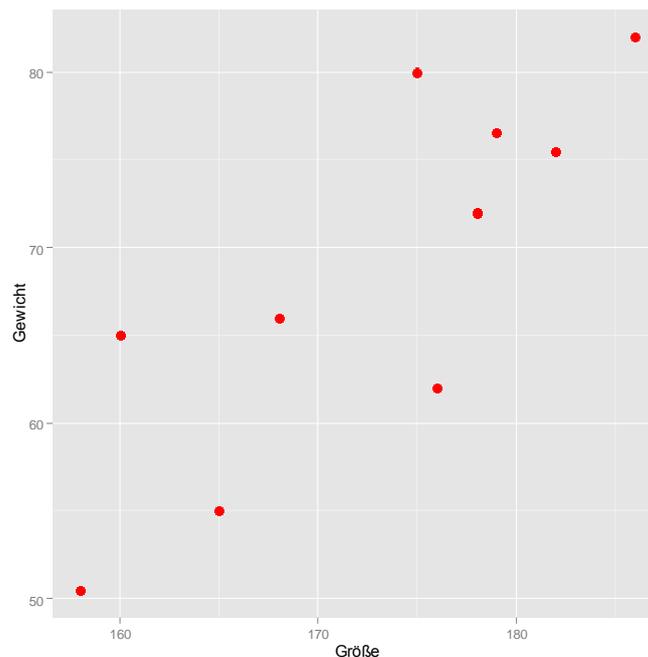
Es ist z.B. möglich, einzelne Bestandteile dieser Liste (z.B. die Achsenbeschriftungen) zu modifizieren:

```
> sd$labels$x<-"GröÙe";sd$labels$y<-"Gewicht"
```

und das Geom (die darstellbare Schicht) anschließend zu ergänzen,

```
> sd+geom_point(colour="red",size=3)
```

um eine modifizierte Grafikausgabe anzufordern:



Wir lernen allerdings in Abschnitt 8.3.1.4.6 eine einfachere Methode zur Änderung der Achsenbeschriftungen kennen.

Weitere Beispiele für Geome bzw. generierende Funktionen sind:

- **geom_histogram()** erstellt ein Histogramm
- **geom_line()** erstellt einen Linienverlauf
- **geom_boxplot()** erstellt Boxplots für die Kategorien eines Faktors
- **geom_text()** erstellt Beschriftungen (z.B. zu den Datenpunkten eines Streudiagramms)

Auf der Webseite <http://docs.ggplot2.org/current/> werden alle in der **ggplot2**-Version 2.2.1 verfügbaren Geome beschrieben.

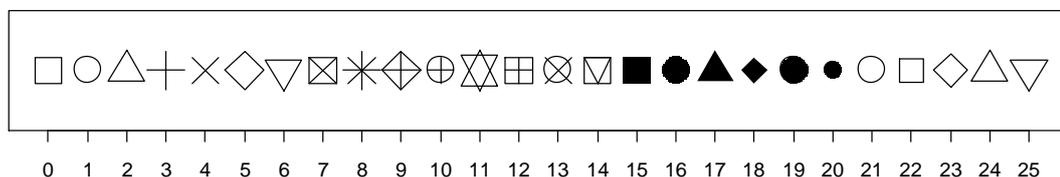
Ein Diagramm muss mindestens eine Schicht (ein Geom) enthalten und kann beliebig komplex aufgebaut sein. Man verwendet den Operator `+`, um Schichten zu ergänzen.

8.3.1.2 Aesthetics

Ein Geom besitzt visuelle Attribute (sogenannte **Aesthetics**), die entweder die Ausprägungen einer Variablen darstellen oder auf einen konstanten Wert gesetzt werden. Das Geom mit den (x,y) - Datenpunkten eines Streudiagramms (von der Funktion `geom_point()` erstellt) unterstützt z.B. die folgenden Aesthetics (mit den **ggplot2**-Namen in Klammern), wobei die beiden ersten Attribute obligatorisch zu versorgen sind:

- X-Koordinaten der Datenpunkte (**x**)
- Y-Koordinaten der Datenpunkte (**y**)
- Form der Datenpunkte (**shape**)

Die folgende Palette wird sowohl von der traditionellen Grafik (vgl. Abschnitt 8.2.4.1) als auch von **ggplot2** genutzt:



- Gesamt-, Rand- bzw. Füllfarbe (**colour**, **fill**)
Bei den Symbolen 15 bis 20 haben der Rand und das Innere dieselbe Farbe, welche durch das Attribut **colour** festgelegt wird. Bei den anderen Symbolen werden die Randfarbe durch das Attribut **colour** und die Innenfarbe durch das Attribut **fill** festgelegt.
- Transparenzgrad der Datenpunkte (**alpha**)
- Größe der Datenpunkte (**size**)

Auch jedes andere Geom besitzt eine Menge von *obligatorischen* sowie eine Menge von *unterstützten* Attributen.

Eine Liste mit Geomen zusammen mit ihren jeweiligen obligatorischen bzw. verfügbaren ästhetischen Attributen findet sich z.B. in Wickham (2009, Tabelle 4.3 in Abschnitt 4.6). Selbstverständlich kann man sich auch über die **R**-Hilfe informieren und z.B. mit

```
> ?geom_point
```

u.a. die ästhetischen Attribute zum Streudiagramm abrufen.

Die Datentabelle mit den zu visualisierenden Variablen sowie die Attribut-Abbildungen können auf Plot-Ebene festgelegt werden und sind dann für alle enthaltenen Schichten bzw. Geome gültig. Im folgenden Beispiel wird die `ggplot()` - Funktion benutzt, um ein Plot-Objekt anzulegen (vgl. Abschnitt 8.3.3):

```
> sd <- ggplot(ggg, aes(größe, gewicht))
```

So vermeidet man es, die Spezifikation für mehrere Schichten wiederholen zu müssen. Wird ein Attribut auf Plot- und Schichtebene abgebildet, dann dominiert die lokale Zuordnung.

Ein ästhetisches Attribut kann auf einen festen Wert abgebildet werden, z.B. die Form von (x,y) - Datenpunkten:

```
> sd + geom_point(shape=3)
```

Weit wichtiger für den intendierten Zweck eines Diagramms ist es jedoch, ästhetische Attribute auf zu visualisierende Variablen abzubilden. Dazu ist die Funktion `aes()` zu verwenden, die als Argumente eine Liste von (Aesthetic = Variable) - Zuordnungen erhält, wobei auch Funktionen von Variablen erlaubt

sind. Im folgenden Beispiel erzeugt die Funktion `geom_smooth()` zu einem gruppierten Streudiagramm eine Schicht mit Anpassungsfunktion und Konfidenzstreifen, wobei die Füllfarbe von der Variablen `geschlecht` abhängig gemacht wird, während der Transparenzgrad datenunabhängig festgelegt wird:

```
> sd + geom_point() + geom_smooth(method="lm", aes(fill=geschlecht), alpha=0.1)
```

Während die konstante Wertzuweisung nur für Schichten bzw. Geome möglich ist, kann die Variablengebundene, via `aes()` realisierte Zuweisung auf eine Schicht oder das gesamte Plot-Objekt angewendet werden.

Häufig arbeiten die Geome zu den Schichten eines Diagramms mit der voreingestellten (auf Plot-Ebene) vereinbarten Datentabelle und den dazu definierten Abbildungen von Attributen (Aesthetics) auf Variablen. Es ist jedoch erlaubt, dass ein Geom eine eigene Datentabelle mit eigenen Aesthetics-Abbildungen verwendet.

8.3.1.3 Statistische Transformationen

Für eine Schicht bzw. das dort präsentierte Geom kann eine **statistische Transformation** (Kurzbezeichnung: **stat**) der darzustellenden Variablen erforderlich sein. Im eben erwähnten Beispiel mit `geom_smooth()` als generierender Funktion entsteht eine neue Datentabelle, indem für gleichabständige Stützstellen die Werte einer Anpassungsfunktion berechnet werden. Ein anderes Beispiel ist die Bildung von Intervallen (`stat = "bin"`) für ein Histogramm (`geom_histogram()`). Wird die identische Transformation als **stat**-Spezialfall einbezogen, kann man sagen, dass zu jedem Geom eine statistische Transformation gehört (siehe Tabelle 4.3 in Abschnitt 4.6 von Wickham 2009). Außerdem hat jedes Geom eine voreingestellte statistische Transformation, und zu jeder Transformation gehört ein voreingestelltes Geom.

Mit der Funktion `args()` lässt sich die Voreinstellung für ein Geom bzw. für eine statistische Transformation in Erfahrung bringen, z.B.:

```
> args(geom_bar)
function (... , stat = "bin", ...)
> args(stat_bin)
function (... , geom = "bar", ...)
```

Durch eine Änderung der Voreinstellung sind eigenständige Diagrammkreationen möglich.

Bei der Erstellung einer Schicht kann man sich zwischen zwei letztlich äquivalenten Vorgehensweisen entscheiden:

- Man verwendet eine **geom**-Funktion und spezifiziert nötigenfalls die statistische Transformation per **stat**-Argument, z.B. (vgl. Abschnitt 8.3.4.3.2)


```
> ggplot(ggg, aes(x=geschlecht, y=gewicht)) + geom_bar(stat="summary", fun.y=mean)
```
- Man verwendet eine **stat**-Funktion und spezifiziert nötigenfalls das gewünschte Geom per **geom**-Argument. Das Resultat aus dem letzten Beispiel lässt sich auch so anfordern:


```
> ggplot(ggg, aes(x=geschlecht, y=gewicht)) + stat_summary(geom="bar", fun.y=mean)
```

Bei den meisten Schichten ist es bequemer, mit einer **geom**-Funktion zu starten.

Ist eine statistische Transformation zu konfigurieren, können die fälligen Argumente für die zuständige Funktion (z.B. `stat_bin()` bei einer Klassenbildung) an die generierende Geom-Funktion (z.B. `geom_histogram()`) übergeben und durchgeschleust werden, z.B.:

```
> ggplot(ggg, aes(gewicht)) + geom_histogram(binwidth = 5)
```

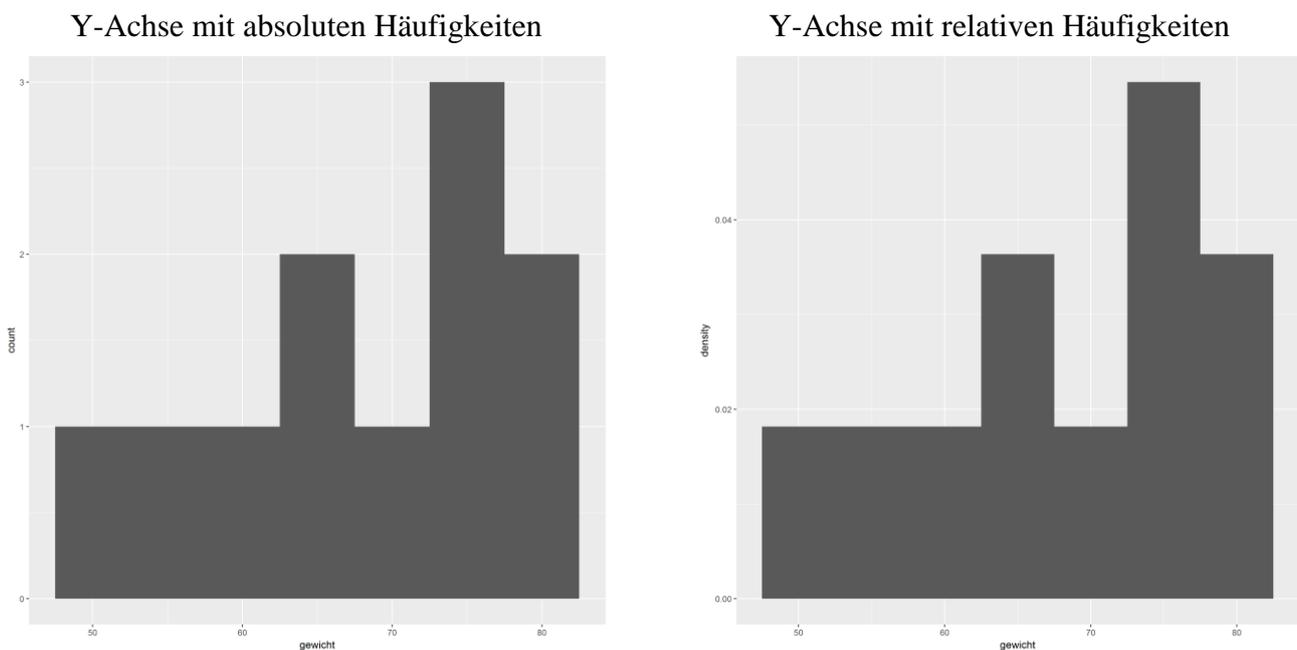
Sollte unklar sein, welche Argumente die zu einem Geom gehörige statistische Transformation unterstützt, kann man ...

- per `args()` - Funktion die voreingestellte statistische Transformation des Geoms ermitteln (siehe oben)
- und per **R**-Hilfe die Argumente der Transformationsfunktion anzeigen lassen (z.B. `?stat_bin`).

Gelegentlich muss man von einer statistischen Transformation erstellte Variablen explizit ansprechen und z.B. einem ästhetischen Attribut zuweisen. Das geschieht im folgenden Beispiel mit der Variablen **density**, welche von der bei `geom_histogram()` voreingestellten Transformation `stat_bin()` produziert wird (vgl. Abschnitt 8.3.4.1). Um Verwechslungen mit Dataframe-Variablen zu vermeiden, sind vor und hinter die Namen von Transformationsergebnissen jeweils zwei Punkte zu setzen. Im folgenden Beispiel wird die Y-Achse auf das Transformationsergebnis **density** abgebildet:

```
> ggplot(ggg, aes(gewicht)) + geom_histogram(binwidth=5, aes(y=..density..))
```

Im Ergebnis werden auf der Y-Achse anstelle absoluter Häufigkeiten relative angezeigt:



8.3.1.4 Skalen, Achsen und Legenden

Ist eine Variable mit einem ästhetischen Attribut verbunden, dann bildet eine so genannte *Skala* die Variablenausprägungen auf die Attributausprägungen ab. Ist z.B. eine Variable mit der X-Achse verbunden, dann müssen ihre Ausprägungen auf Positionswerte abgebildet werden, die das Grafiksystem versteht. Weil die **ggplot2**-Grafik auf dem **grid**-Paket basiert, sind für die X- bzw. Y-Positionen Werte im Intervall [0, 1] erlaubt. Skalen sind aber nicht nur für die Positionsattribute relevant, sondern auch für alle anderen ästhetischen Attribute (z.B. Farbe, Form, Größe, Linientyp).

Zu den wesentlichen Aufgaben der Skalen gehört auch die Gestaltung von Achsen und Legenden, die gemeinsam als **Guides** bezeichnet werden und dem Betrachter die Interpretation von Diagrammen ermöglichen:

- Für Positionsattribute übernehmen Achsen die Interpretationshilfe.
- Die restlichen ästhetischen Attribute werden durch die Legende erläutert.

Zwar ist bei jedem datenabhängigen ästhetischen Attribut eine Skala im Spiel, doch können oft die von **ggplot2** erstellten Voreinstellungen beibehalten werden. Zur Modifikation von Skalen dienen Funktionen mit einem Namen nach dem folgenden Schema:

scale_<aesthetic>_<type>

Als **type**-Werte sind z.B. **discrete**, **continuous** und **manual** erlaubt. So resultieren Funktionsnamen wie

- `scale_x_continuous()`
- `scale_fill_manual()`

Eine **scale**-Funktion wird per `+`-Operator zu einem Diagramm hinzugefügt, z.B.:

```
> sd <- ggplot(ggg, aes(größe, gewicht))
> sd + geom_point(colour="red", size=3) +
+ scale_x_continuous("Größe", limits=c(150,190), breaks=seq(150,190,by=5), minor_breaks=NULL)
```

8.3.1.4.1 Generelle Argumente der **scale**-Funktionen

Den **scale**-Funktionen sind die folgenden Argumente gemeinsam:

- **name**
Mit dem ersten Argument kann man den Titel einer Achse (bei **x** oder **y** als Aesthetic) oder einer Legende (bei **colour**, **fill**, **size**, **shape** oder **linetype** als Aesthetic) festlegen. Neben einfachen Texten können über **expression()** - Aufrufe Formeln mit mathematischer Typographie realisiert werden.
- **limits**
Dieses Argument legt die zu berücksichtigenden Werte der dargestellten Variablen fest. Bei einer metrischen Variablen ist ein numerischer Vektor mit dem minimalen und dem maximalen zu berücksichtigenden Wert anzugeben, z.B.:
`limits=c(150,190)`
Bei einer diskreten Variablen ist ein **character**-Vektor mit den zu berücksichtigenden Werten anzugeben, wobei auch die Reihenfolge (von Balken etc.) beeinflusst wird, z.B.:
`limits=c("Frau", "Mann")`
Fälle mit einem ausgeschlossenen Wert bei *irgendeiner* dargestellten Variablen werden vom Diagramm ausgeschlossen.
- **breaks**
Für eine metrische Variable werden in einem numerischen Vektor die Hauptunterteilungspunkte für die Achse (z.B. beim ästhetischen Attribut **x**) bzw. die Legendenwerte (z.B. beim ästhetischen Attribut **size**) festgelegt. Eine Achse zeigt per Voreinstellung zwischen den etikettierten Hauptunterteilungspunkten noch mittig positionierte Nebenunterteilungen ohne Etikett (abzuschalten mit: **minor_breaks=NULL**).
- **labels**
Ein Vektor mit Etiketten zu den Hauptunterteilungspunkten bzw. zu den Legendeneinträgen. Bei metrischen Achsenvariablen ist es meist überflüssig, die automatisch angezeigten Zahlen durch Etiketten zu ersetzen. Stattdessen kann es sinnvoll sein, durch eine Formatierungsfunktion als Wert für das **labels**-Argument die Darstellung der Zahlen zu beeinflussen. Die folgende Formatierungsfunktion
`comma <- function(x) {format(x, decimal.mark="," , trim=TRUE)}`
sorgt für eine Anzeige mit Dezimal*komma*, das exakt unter dem Hauptunterteilungspunkt erscheint. anschließend ist die Verwendung im Rahmen einer Skalierungsfunktion zu sehen:
`scale_x_continuous("Spannung", breaks=seq(10,14,by=0.5), labels=comma)`

8.3.1.4.2 Positionierungsskalen

Im folgenden Streudiagramm

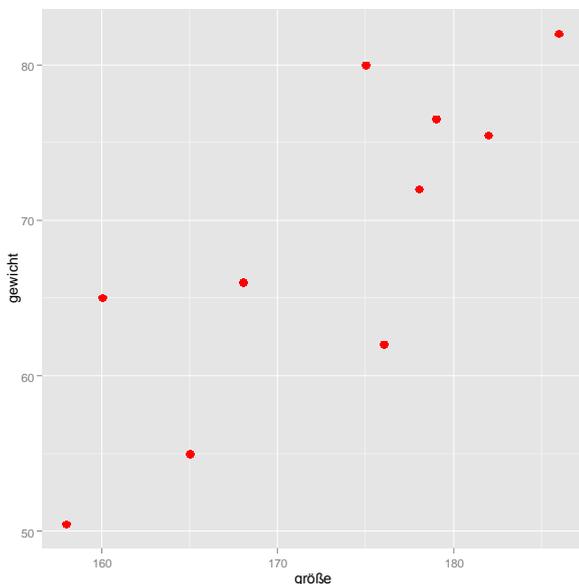
```
> sd <- ggplot(ggg, aes(größe, gewicht))
> sd + geom_point(colour="red", size=3) +
+ scale_x_continuous("Größe", limits=c(150,190), breaks=seq(150,190,by=5), minor_breaks=NULL)
```

wird mit der Funktion **scale_x_continuous()** die Skala zur X-Achse modifiziert:

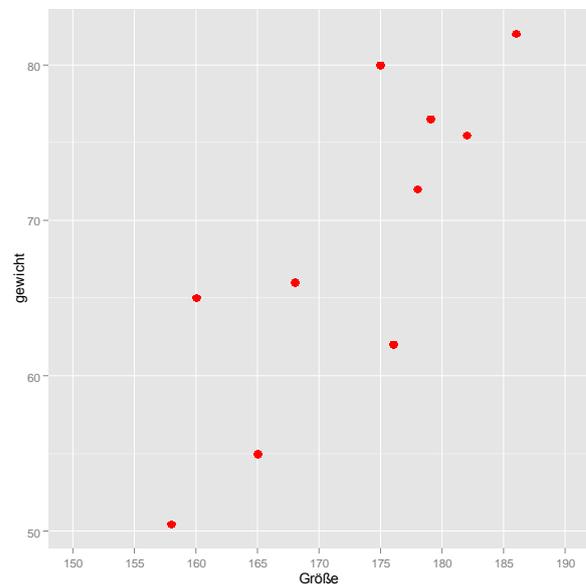
- Der Titel wird geändert.
- Der darzustellende Wertebereich wird festgelegt.
- Per **seq()** - Funktion entsteht ein Vektor mit Hauptunterteilungspunkten im 5er-Abstand als Wert für das **breaks**-Argument.
- Die Nebenunterteilungspunkte (**minor_breaks**) werden abgeschaltet.

Die Skala wird dem Plot-Objekt über den Operator + zugewiesen und ersetzt dabei die voreingestellte Skala für das betroffenen ästhetische Attribut.

Ergebnis mit der voreingestellten x-Skala



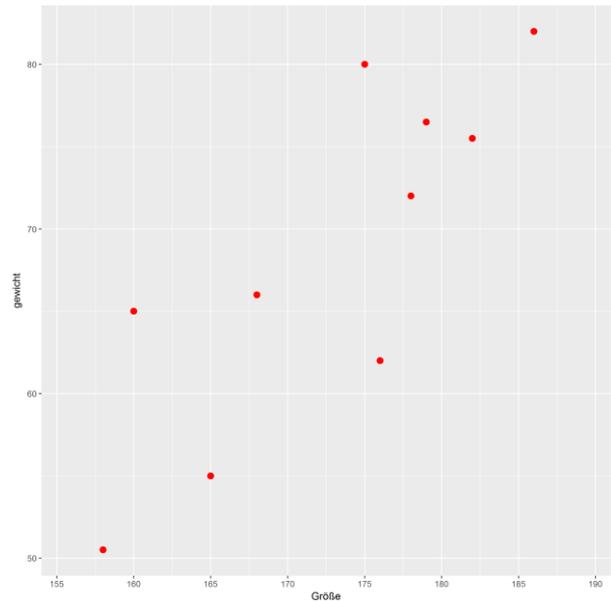
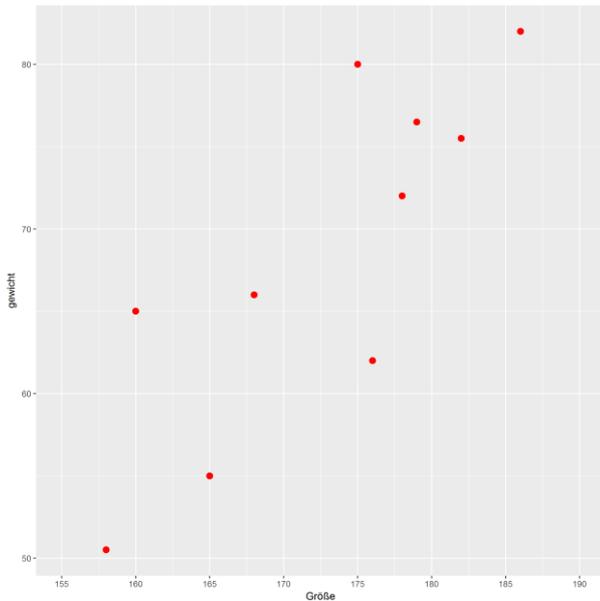
Ergebnis mit der modifizierten x-Skala



Per Voreinstellung halten die Positionierungsskalen einen Randabstand ein. Über das **expand**-Argument lässt sich dieses Verhalten steuern, wobei ein numerischer Vektor mit zwei Elementen anzugeben ist. Mit dem ersten Element wird eine multiplikative und mit dem zweiten Element eine additive Erweiterung der Randzone bewirkt, z.B.:

```
> sd <- ggplot(ggg, aes(größe, gewicht)) + geom_point(colour="red", size=3)
> sd + scale_x_continuous("Größe", limits=c(155,190), breaks=seq(155,190,by=5), expand=c(0,1))
```

Links ist der voreingestellte und rechts ist der modifizierte (verkleinerte) Randabstand für die x-Skala zu sehen:



Bei den kontinuierlichen Positionierungsskalen lässt sich für die abgebildete Variable über das **trans**-Argument eine Transformation vornehmen, z.B.:

```
> sd + scale_x_continuous(trans = "log10")
```

Neben etlichen eingebauten Transformationen (siehe Wickham 2009, S. 97) können auch selbstdefinierte verwendet werden. Zur Vereinfachung der Diagrammspezifikation existierten Skalierungsfunktionen mit einer eingebauten Transformation, z.B.:

```
> sd + scale_x_log10()
```

Über Positionierungsskalen für Datums- und Zeitvariablen informiert Wickham (2009, S. 100f).

8.3.1.4.3 Farbskalen

Bei der Abbildung von *metrischen* Variablen auf Füll- und Linienfarben (also auf die ästhetischen Attribute **fill** und **colour**) unterstützt **ggplot2** drei verschiedene **Gradienten**, die sich bei der Definition von Zwischenfarben unterscheiden:¹

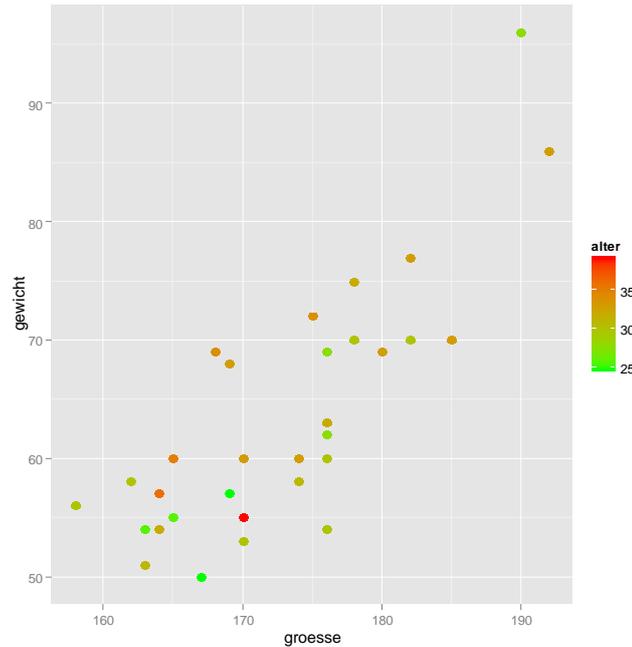
- **scale_fill_gradient()**, **scale_colour_gradient()**

Die Argumente **low** und **high** definieren den Farbverlauf, z.B.:

```
> csd <- ggplot(kfa, aes(x=größe, y=gewicht, colour=alter)) + geom_point(size=3)
> csd + scale_colour_gradient(low="green", high="red")
```

Das im Beispiel resultierende Streudiagramm zeigt die drei metrischen Variablen Größe, Gewicht und Alter:

¹ Bei einem Streudiagramm (**geom_point()**, siehe Abschnitt 8.3.1.2) hängt die Bedeutung der Attribute **colour** und **fill** von der Form der Punkte ab.



Zu den angegebenen Funktionsnamen existieren äquivalente Schreibweisen:

- Statt **colour** wird auch **color** akzeptiert
- Statt **gradient** kann man **continuous** schreiben.
- **scale_fill_gradient2()**, **scale_colour_gradient2()**
Zusätzlich zu den Argumenten **low** und **high** für die Endfarben des Verlaufs kann mit dem Argument **mid** eine Zwischenfarbe festgelegt werden.
- **scale_fill_gradientn()**, **scale_colour_gradientn()**
Es kann ein Vektor mit Farben angegeben werden.

Bei der Abbildung von *kategorialen* Variablen auf Füll- und Linienfarben (also auf die ästhetischen Attribute **fill** und **colour**) unterstützt **ggplot2** drei Verfahren:

- **scale_fill_hue()**, **scale_colour_hue()**
Diese Funktionen sind per Voreinstellung im Einsatz und liefern gleichmäßig im Sinne des HCL-Farbraums (Hue, Chroma, Luminance) verteilte Farbtöne. Nach Wickham (2009, S. 106) sind die Ergebnisse akzeptabel, sofern nicht mehr als 8 Farben beteiligt sind.
- **scale_fill_brewer()**, **scale_colour_brewer()**
Über diese Skalierungsfunktionen können per **palette**-Argument die Farbpaletten des **Color-Brewer** - Online-Werkzeugs (<http://colorbrewer2.org/>) gewählt werden. Für kategoriale Variablen empfiehlt Wickham (2009, S. 106) die Paletten **Set 1** und **Dark2**, z.B.:

```
> csd <- ggplot(kfa, aes(x=größe, y=gewicht, colour=geschlecht)) + geom_point(size=3)
> csd + scale_colour_brewer(palette="Dark2")
```

Mit der folgenden Anweisung

```
> RColorBrewer::display.brewer.all()
```

erhält man eine Beschreibung aller Brewer-Paletten.

- **scale_fill_manual()**, **scale_colour_manual()**
Statt einen Farbwahlalgorithmus oder eine vorgefertigte Palette zu verwenden, kann man alle Farben individuell bestimmen (siehe Abschnitt 8.3.1.4.4).

Weitere Details zur Verwendung von Farben in **ggplot2** finden sich bei Chang (2013, S. 251ff) und Wickham (2009, S. 102ff).

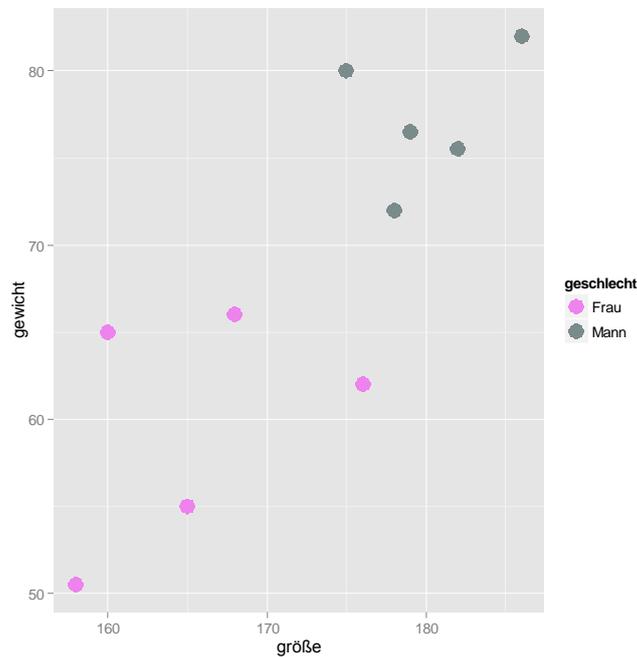
8.3.1.4.4 Manuelle Skalen

Ist ein diskretes Merkmal (z.B. repräsentiert durch einen Faktor) mit einem ästhetischen Attribut verbunden, dann werden den einzelnen Ausprägungen sukzessive die Darstellungsvarianten (z.B. Farben oder Symbole) aus einer vorgegebenen Palette zugeordnet. Über das **values**-Argument einer **scale**-Funktion kann man die Voreinstellungen durch individuelle Werte ersetzen.

Im ersten Beispiel werden in einem Streudiagramm für das zur Unterscheidung von zwei Gruppen verwendete Attribut **colour** über das **values**-Argument der Funktion **scale_colour_manual()** individuelle Farben vereinbart:¹

```
> gsd <- ggplot(ggg,aes(größe,gewicht,colour=geschlecht)) + geom_point(size=5)
> gsd + scale_colour_manual(values=c("violet","lightcyan4"))
```

Dabei kommt das von **geom_point()** per Voreinstellung verwendete Symbol zum Einsatz (ein gefüllter Kreis):



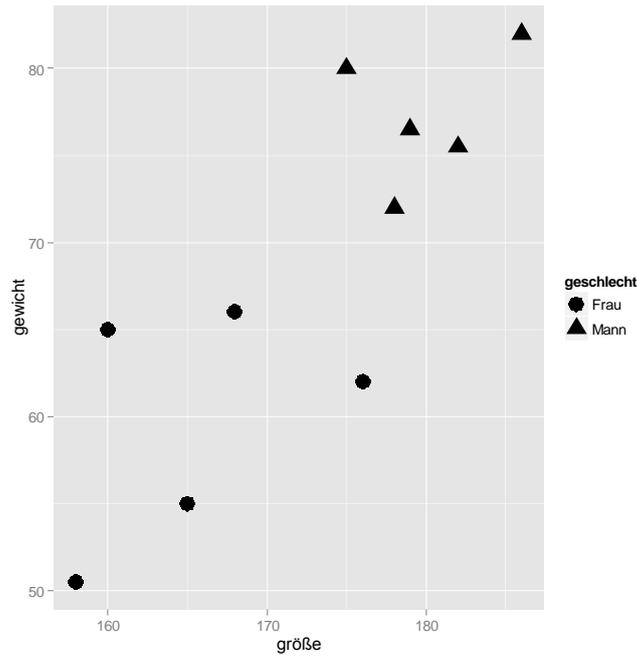
Dient das **shape**-Attribut zur Unterscheidung der Geschlechtsgruppen,

```
> ggplot(ggg,aes(größe,gewicht,shape=geschlecht)) + geom_point(size=5)
```

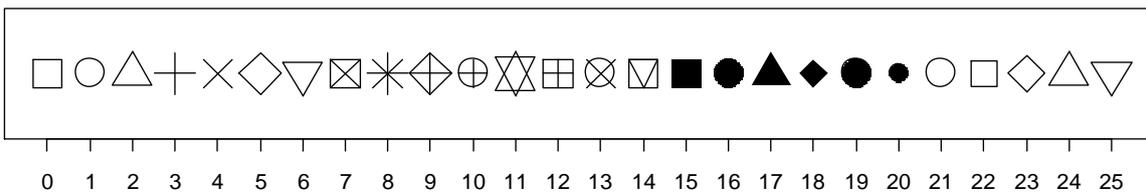
verwendet **geom_point()** per Voreinstellung die Symbole ● und ▲:

¹ In der folgenden PDF-Datei (abgerufen am 06.08.2017) sind die in **R** verfügbaren Farbnamen mit Musterflächen dokumentiert:

<http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf>



Dies sind die Symbole 16 und 17 aus der folgenden Palette, die sowohl von der traditionellen Grafik (vgl. Abschnitt 8.2.4.1) als auch von **ggplot2** genutzt wird:



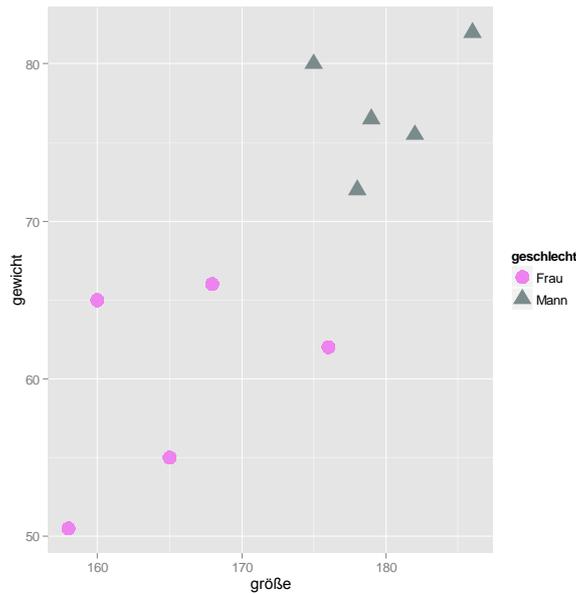
Über das **values**-Argument der Funktion **scale_shape_manual()** lassen sich alternative Symbole wählen:

```
> gsd <- ggplot(ggg,aes(größe,gewicht,shape=geschlecht)) + geom_point(size=5)
> gsd + scale_shape_manual(values=c(15,22))
```

Selbstverständlich ist es möglich, Form *und* Farbe der Symbole simultan zur Gruppenunterscheidung zu verwenden:

```
> gsd <- ggplot(ggg,aes(größe,gewicht,colour=geschlecht,shape=geschlecht))
> gsd + geom_point(size=5) + scale_colour_manual(values=c("violet","lightcyan4"))
```

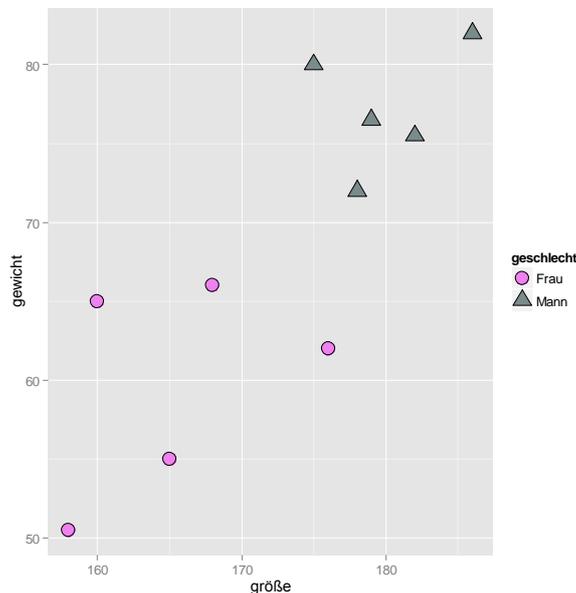
Im Beispiel kommen die voreingestellten Symbole mit individuellen Farben zum Einsatz:



Bei den Symbolen mit Nummern außerhalb des Bereichs von 15 bis 20 (siehe obige Palette) kann die Rand- und die Füllfarbe separat gesteuert werden. In der folgenden Variante des letzten Diagramms

```
> gsd <- ggplot(ggg, aes(größe, gewicht, fill=geschlecht, shape=geschlecht))
> gsd <- gsd+geom_point(size=5)
> gsd + scale_fill_manual(values=c("violet", "lightcyan4"))+scale_shape_manual(values=c(21,24))
```

kommen die Symbole 21 und 24 zum Einsatz, wobei die datengebundenen Füllfarben durch das voreingestellte Schwarz umrandet werden:



8.3.1.4.5 Legende

Eine Legende wird von **ggplot2** bei Bedarf automatisch eingefügt, wenn für ein ästhetisches Attribut die Rückübersetzung in Datenwerte erläutert werden muss. Das ist z.B. erforderlich, wenn in einem gruppierten Streudiagramm Markierungen mit unterschiedlichen Farben auftreten:

```
> ggplot(ggg, aes(größe, gewicht, colour=geschlecht)) + geom_point(size=3)
```

Sind *mehrere* Schichten bzw. ästhetische Attribute zu erläutern, strebt **ggplot2** die minimale Anzahl von Legenden an, so dass z.B. in einem gruppierten Streudiagramm mit gruppenspezifischen Anpassungsfunktionen (vgl. Abschnitt 8.3.3.4) *eine* Kombilegende mit Symbolen und Linien entsteht:

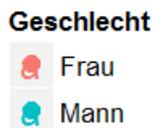


Ist der zu einem Geom erscheinende Legendenbeitrag unerwünscht, lässt er sich über das Argument **show.legend** der Geom-Funktion abschalten, z.B. beim folgenden Geom zur Anzeige von Beschriftungen für Datenpunkte in einem gruppierten Streudiagramm (vgl. Abschnitt 8.3.3):

```
> gsd <- ggplot(ggg, aes(größe, gewicht, colour=geschlecht)) + geom_point(size=3)
> gsd + geom_text(aes(y=gewicht-1, label=rownames(ggg)), show.legend=FALSE)
```

Im Beispiel wird so vermieden, dass aufgrund der Text-Geoms an die Legendensymbole ein unvorteilhaftes „a“ angeheftet wird:

Mit unerwünschtem Legendenbeitrag



Korrekte Lösung



Um den Titel einer Legende zu ändern, kann man die zum ästhetischen Attribut gehörige **scale**-Funktion aufrufen und das **name**-Attribut auf den gewünschten Wert setzen, z.B.:

```
> gsd + scale_colour_discrete(name="Geschlecht")
```

In Abschnitt 8.3.1.4.6 wird eine bequemere Alternative zur **scale**-Funktion vorgestellt für den Fall, dass lediglich der Titel zu ändern ist.

Sind mehrere ästhetische Attribute per Legende zu erläutern, dann müssen die Legendentitel identisch gewählt werden, damit eine Kombilegende erscheint. Im folgenden Beispiel

```
> gsd <- ggplot(ggg, aes(größe, gewicht, colour=geschlecht, shape=geschlecht))
> gsd <- gsd + geom_point(size=5) + scale_colour_hue(name="Geschlecht")
> gsd + scale_shape_discrete(name="Gender")
```

wird diese Regel missachtet, so dass zwei getrennte Legenden entstehen:



Zur Positionierung der Legende siehe Abschnitt 8.3.1.7.2.

8.3.1.4.6 Hilfsfunktionen für Wertebereiche und Beschriftungen

Geht es lediglich um eine Änderung der zu berücksichtigenden Wertebereiche geht, dann sind die folgenden Hilfsfunktionen einfacher zu handhaben als die **scale**-Funktionen:

- **xlim()**
Mit dieser Funktion legt man den die zu berücksichtigenden Werte der X-Achsenvariablen fest, was bei einer metrischen

```
> sd + xlim(150,190)
```

und bei einer diskreten Variablen möglich ist:

```
> balken + xlim("L", "M", "B")
```
- **ylim()**
Mit dieser Funktion legt man die zu berücksichtigenden Werte der Y-Achsenvariablen fest, z.B.:

```
> sd + ylim(50,90)
```

Fälle außerhalb dieser Grenzen werden ausgeschlossen, was sich auch auf andere Schichten auswirkt.

Geht es lediglich um eine Änderung der Beschriftung, dann sind die folgenden Funktionen einfacher zu handhaben als die **scale**-Funktionen:

- **xlab()**
Mit dieser Funktion lässt sich die Beschriftung der X-Achse ändern, z.B.:

```
xlab("Größe")  
> sd + xlab("Größe")
```
- **ylab()**
Mit dieser Funktion lässt sich die Beschriftung der Y-Achse ändern, z.B.:

```
> sd + ylab("Gewicht")
```
- **ggtitle()**
Mit dieser Funktion kann man einen Diagrammtitel ergänzen, wobei ein Zeilenwechsel mit der Escape-Sequenz `\n` (New Line) veranlasst wird, z.B.:

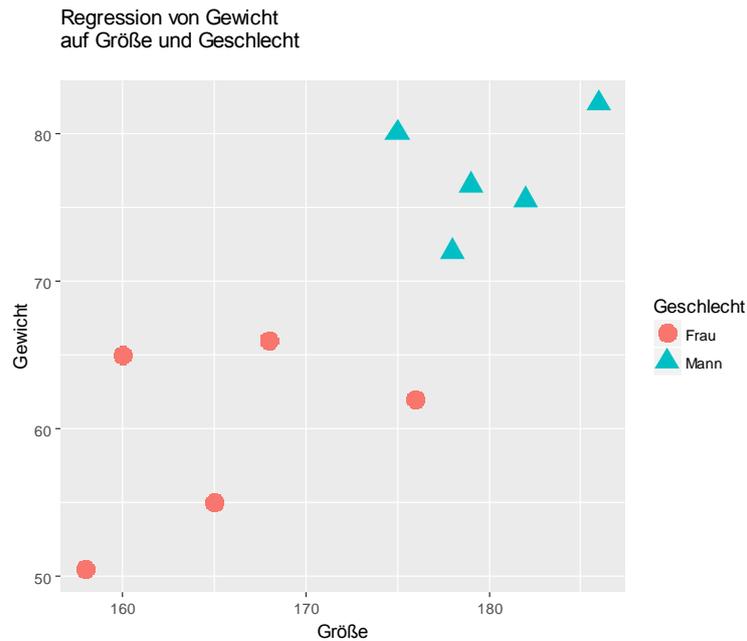
```
> sd + ggtitle("Regression von Gewicht\nauf Größe und Geschlecht\n")
```
- **labs()**
Mit dieser Funktion kann man die Beschriftungen von Achsen und Legenden ändern, z.B.:

```
> sd + labs(x="Größe", y="Gewicht", colour="Geschlecht")
```

Eine Legendenbeschriftung muss dem zuständigen ästhetischen Attribut zugewiesen werden. Auch ein Diagrammtitel lässt sich vereinbaren, z.B.:

```
> sd + labs(title="Regression von Gewicht\nauf Größe und Geschlecht\n")
```

In der Version 2.2.0 von **ggplot2** wird der Diagrammtitel per Voreinstellung *links* ausgerichtet, z.B.:



Das lässt sich über die Funktion **theme()** ändern (vgl. Abschnitt 8.3.1.7.2):

```
> gsd <- ggplot(ggg, aes(größe, gewicht, colour=geschlecht, shape=geschlecht))
> gsd <- gsd + geom_point(size=5) + scale_colour_hue(name="Geschlecht")
> gsd <- gsd + scale_shape_discrete(name="Geschlecht") + labs(x="Größe", y="Gewicht")
> gsd <- gsd + ggtitle("Regression von Gewicht\nauf Größe und Geschlecht\n")
> gsd + theme(plot.title = element_text(hjust = 0.5))
```

Über **expression()** - Aufrufe können Formeln mit mathematischer Typographie in die Beschriftungen aufgenommen werden, z.B.

```
> dia + ylab(expression(f2))
```

8.3.1.5 Positionsanpassungen

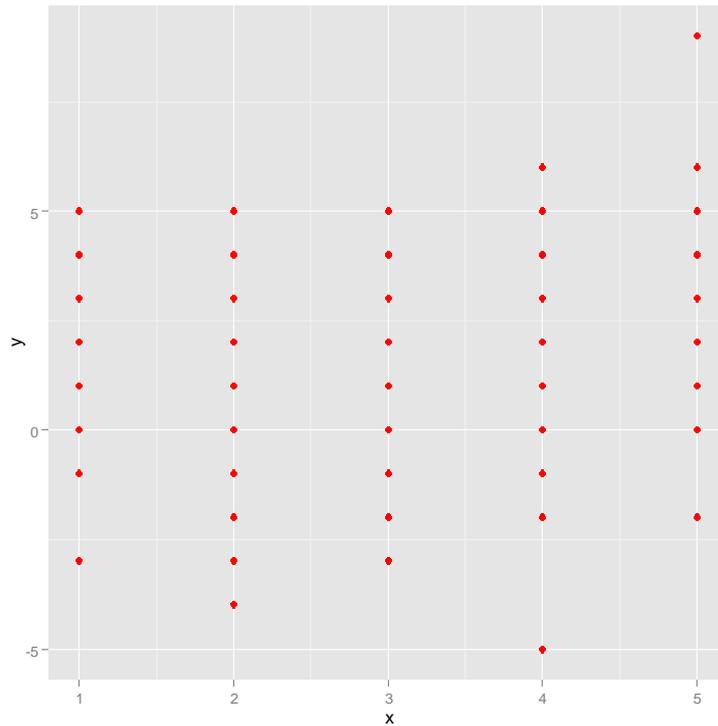
Um die Überlappung ihrer Elemente zu verhindern, kann eine Schicht **Positionsanpassungen** vornehmen (siehe Wickham 2009, Abschnitt 4.8). Bei einem Streudiagramm hilft manchmal die einfache Jitter-Technik, wobei kleine Zufallswerte zu den Daten addiert werden (vgl. Abschnitt 8.2.5.2.5). Im folgenden Beispiel mit simulierten Daten

```
> n <- 200
> set.seed(18)
> x <- round(runif(n, 1, 5), digits=0)
> e <- round(rnorm(n, 0, 2), digits=0)
> y <- 0.6*x + e
> df <- data.frame(x,y)
```

kann über das mit folgender Anweisung

```
> ggplot(df, aes(x,y)) + geom_point(colour="red")
```

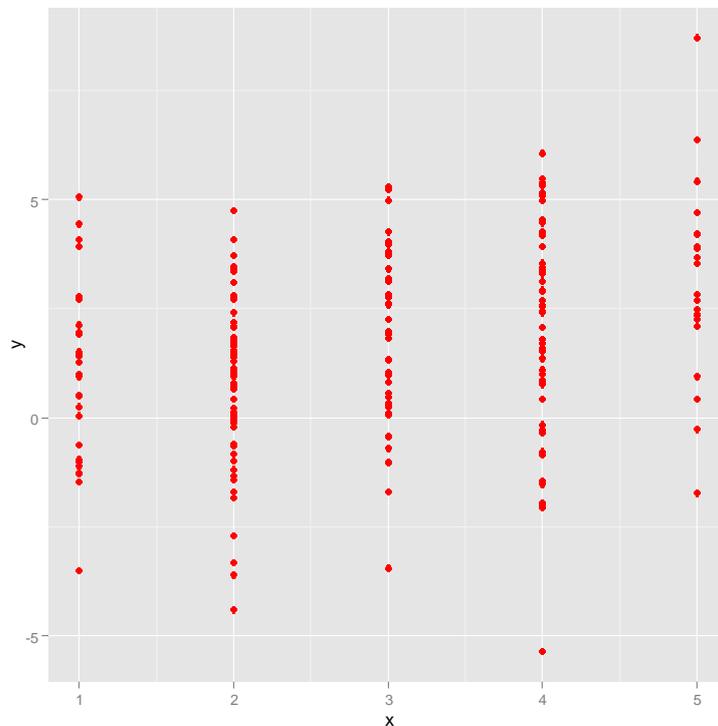
erstellte Streudiagramm die Stärke des Effekts von **x** auf **y** relativ schlecht beurteilt werden:



Daher wird über das **position**-Argument der Funktion **geom_point()** für eine Positionsanpassung gesorgt:

```
> ggplot(df,aes(x,y)) + geom_point(colour="red",position=position_jitter(width=0,height=0.5))
```

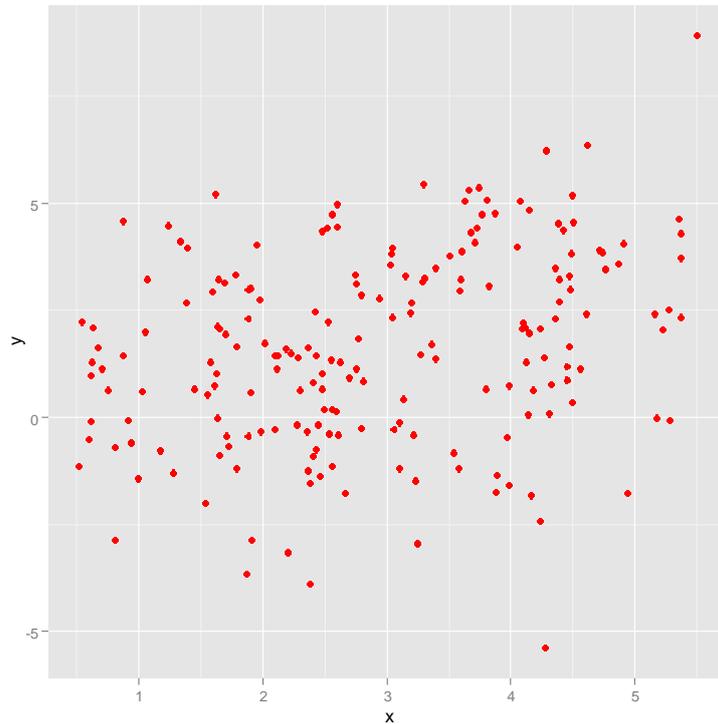
Es kommt die Funktion **position_jitter()** zum Einsatz, wobei die „Verwacklung“ in X- und Y-Richtung über entsprechende Argumente gesteuert werden kann. Im Beispiel bleiben die Regressorwerte unverändert, während die Kriteriumswerte eine Darstellungskorrektur erhalten. Nun sind die Dichteverhältnisse der bedingten Verteilungen und die Stärke der Beziehung besser zu beurteilen:



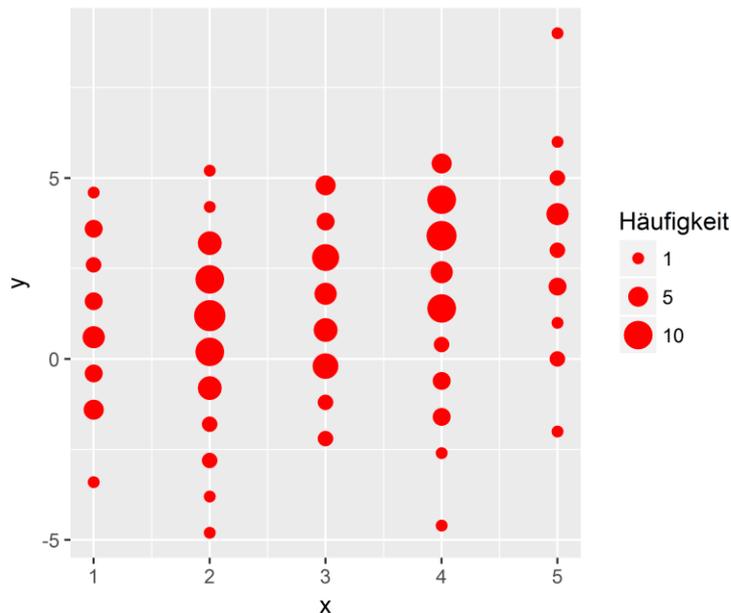
Werden auch die X-Positionen verwackelt,

```
> ggplot(df,aes(x,y))+geom_point(colour="red",position=position_jitter(width=0.5,height=0.5))
```

dann resultiert ein ausdrucksstarkes Diagramm, das allerdings von der Realität mit 5-stufig erfasstem Regressor abweicht:



Eine Alternative zur Jitter-Technik besteht im Beispiel darin, bei mehrfach besetzten Wertekombinationen die Punktgröße entsprechend zu erhöhen. Zur Realisation des folgenden Diagramms



wird in der Punkteschicht die statistische Transformation **stat_sum()** an Stelle der voreingestellten Identität verwendet und die dabei entstehende Ergebnisvariable **n** auf das ästhetische Attribut **size** abgebildet:

```
> psd <- ggplot(data=df, aes(x, y))
> psd <- psd + stat_sum(aes(size=..n..), geom="point", colour="red")
> psd + scale_size_continuous(breaks=c(1,5,10)) + labs(size="Häufigkeit")
```

Mit der Funktion **scale_size_continuous()** wird festgelegt, welche Werte in der Legende erscheinen sollen.

8.3.1.6 Facetten

Die in **ggplot2** realisierte *Grammar of Graphics* bietet die Möglichkeit der Facettierung, wobei die Gesamtstichprobe nach einer oder nach mehreren Variablen aufgeteilt und für jede Teilstichprobe ein Diagramm (mit allen Schichten) erstellt wird.

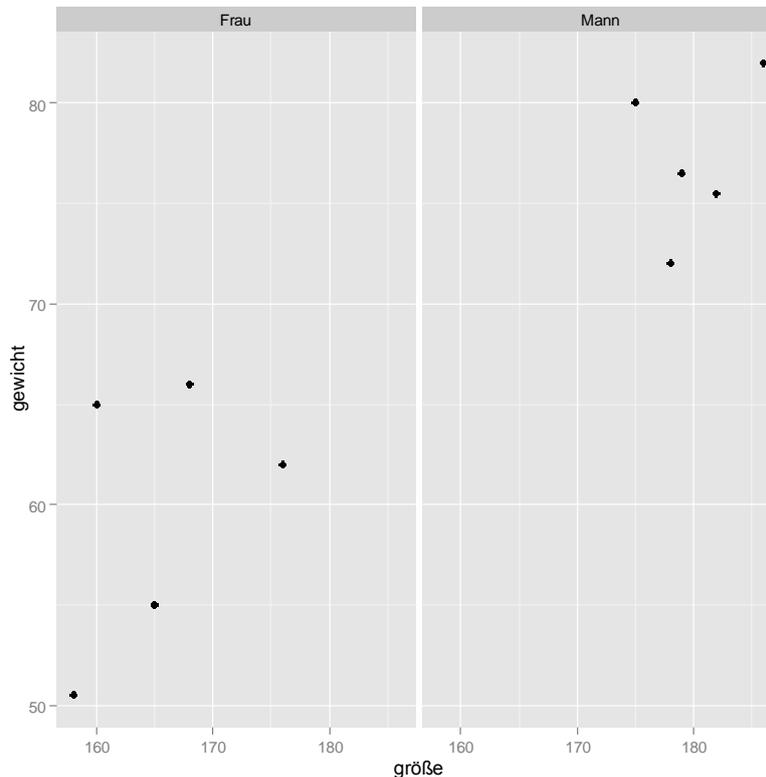
Um die Anordnung der einzelnen Diagramme zu kontrollieren, stehen zwei Facettierungs-Funktionen zur Verfügung:

- **facet_wrap()**
Diese Funktion produziert eine Sequenz von Diagrammen, die per „Zeilenumbruch“ auf der rechteckigen Ausgabefläche angeordnet werden. Meist gibt man nur *eine* steuernde Variable an.
- **facet_grid()**
Diese Funktion erlaubt es, durch *zwei* steuernde Variablen eine Anordnungsmatrix zu definieren.

Im folgenden Beispiel

```
> ggplot(ggg, aes(größe, gewicht)) + geom_point() + facet_wrap(~ geschlecht)
```

wird ein einfaches Streudiagramm (vgl. Abschnitt 8.3.3) mit den Variablen **größe** und **gewicht** definiert, wobei die Funktion **facet_wrap()** dafür sorgt, dass separate Diagramme für Frauen und Männer (die Ausprägungen des Faktors **geschlecht**) erstellt werden:



Um das Anordnungsdesign zu definieren, gibt man im Argument von **facet_wrap()** eine Tilde und anschließend die steuernde Variable an.

Ausführliche Informationen zur Facettierung sind im Kapitel 7 von Wickham (2009) zu finden.

8.3.1.7 Themes

Über das Themes-System von **ggplot2** lassen sich datenunabhängige Aspekte im Erscheinungsbild eines Diagramms kontrollieren (z.B. Schriftarten und Farben von Texten). Man kann ...

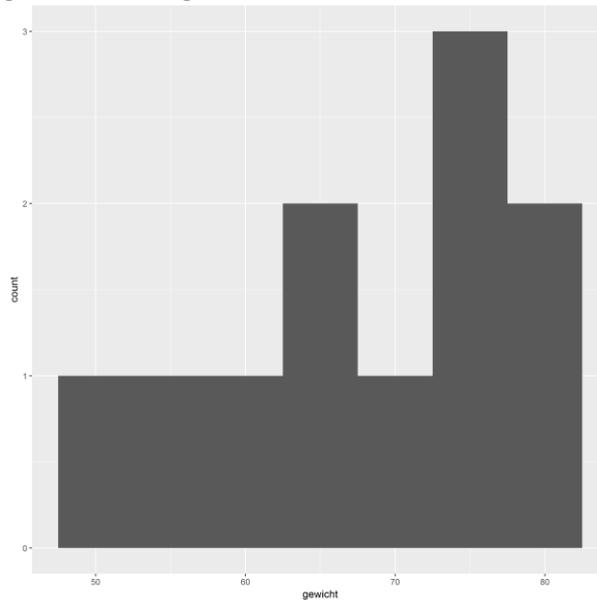
- das Standard-Theme komplett durch eine Alternative (ein anderes Einstellungspaket) ersetzen
- und/oder einzelne Elemente eines Themes modifizieren

8.3.1.7.1 Themes im Ganzen

ggplot2 besitzt eingebaute Themes, die per Funktionsaufruf zu wählen sind, z.B.:

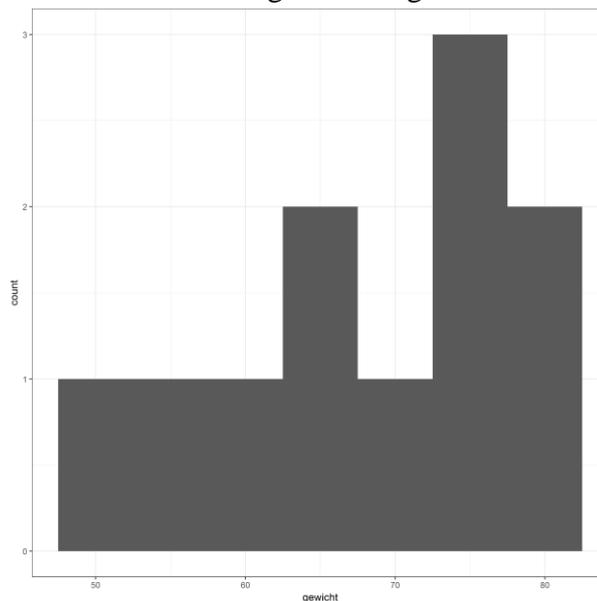
- **theme_grey()**

Dieses Theme mit einem grauen Hintergrund und weißen Gitterlinien ist voreingestellt, z.B.:



- **theme_bw()**

Dieses Theme verwendet einen weißen Hintergrund und graue Gitterlinien, z.B.:



Über weitere Themes informiert die Dokumentation zum Paket **ggplot2**.

Um ein Theme auf ein einzelnes Diagramm anzuwenden, „addiert“ man es zum Plot-Objekt, z.B.:

```
> histo + theme_bw()
```

Um ein Theme für *alle* im weiteren Verlauf einer **R**-Sitzung angefertigten **ggplot2**-Diagramme zu wählen, verwendet man die Funktion **theme_set()**, welche das bisherige Theme als Rückgabe liefert, so dass eine spätere Wiederherstellung möglich ist:

```
> actTheme <- theme_set(theme_bw())
```

Den Funktionen zur Theme-Auswahl ist gemeinsam, dass mit den Argumenten **base_size** bzw. **base_family** für die Beschriftungen eine Basisgröße und eine Schriftartenfamilie gewählt werden kann, z.B.:

```
> histo + theme_bw(base_size=14, base_family="serif")
```

In der Basisschriftgröße (Voreinstellung: 20pt) erscheinen die Achsentitel. Davon abgeleitete Größen sind:

- Größe des Titels: 120% der Basisgröße
- Größe der Teilstrichbeschriftungen: 80% der Basisgröße

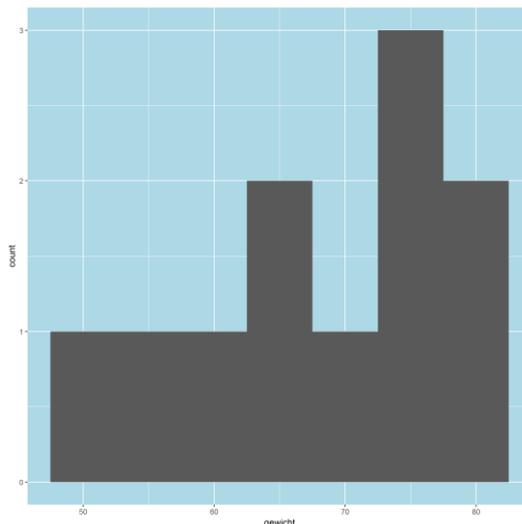
Zu den möglichen Schriftartenfamilien siehe Abschnitt 8.2.2.

8.3.1.7.2 Elemente von Themes modifizieren

Mit der Funktion **theme()**, welche die veraltete Funktion **opts()** ersetzt, kann man Elemente eines Themes modifizieren. Im folgenden Beispiel wird über das **theme()** - Argument **panel.background**

```
> histo + theme(panel.background = element_rect(fill = "lightblue"))
```

die Hintergrundfarbe eines Diagramms mit dem Standard-Theme verändert:



Auf der nach **?theme** erscheinenden Hilfeseite werden ca. 50 **theme()** - Argumente beschrieben, z.B.:

- **panel.background**
Über die Elementfunktion **element_rect()** lässt sich der Hintergrund des Zeichenbereichs gestalten, z.B. die Rand- und die Füllfarbe über die Argumente **colour** und **fill**:

```
> plot + theme(panel.background=element_rect(colour="black", fill="lightblue"))
```

- **plot.title, axis.title, axis.title.x, axis.title.y**

Für Diagramm- und Achsentitel lassen sich mit Hilfe der Elementfunktion **element_text()** u.a. die Schriftartenfamilie (Argument **family**), die Farbe (Argument **colour**), die Größe (Argument **size**) und die Ausrichtung (Argumente **hjust** und **vjust**) beeinflussen. Man kann die beiden Achsen gemeinsam (durch **axis.title**) oder separat (durch **axis.title.x** bzw. **axis.title.y**) ansprechen. Im folgenden Beispiel wird für einen horizontal zentrierten Diagrammtitel in blauer Farbe gesorgt:

```
> plot + theme(plot.title = element_text(hjust = 0.5, colour = "blue"))
```

- **axis.line, axis.line.x, axis.line.y**

Über die Elementfunktion **element_line()** lassen sich Linienattribute der Achsen gestalten. Man kann die beiden Achsen gemeinsam (durch **axis.line**) oder separat (durch **axis.line.x** bzw. **axis.line.y**) ansprechen. Im folgenden Beispiel wird die Linienstärke beider Achsen über das **element_line()** - Argument **size** auf den Wert 1 gesetzt (Voreinstellung: 0):

```
> plot + theme(axis.line = element_line(size = 1))
```

- **axis.text, axis.text.x, axis.text.y**

Über die Elementfunktion **element_text()** lassen sich die Teilstrichbeschriftungen gestalten (siehe obige Beschreibung zu den Achsenbeschriftungen). Man kann die beiden Achsen gemeinsam (durch **axis.text**) oder separat (durch **axis.text.x** bzw. **axis.text.y**) ansprechen. Im folgenden Beispiel wird für die Teilstrichbeschriftungen beider Achsen die Schriftgröße 12 eingestellt (in der Maßeinheit Punkt, 1 Punkt = 1/72 Zoll):

```
> plot + theme(axis.text = element_text(size = 12))
```

- **axis.ticks, axis.ticks.x, axis.ticks.y**

Über die Elementfunktion **element_line()** lassen sich die Achsenteilstriche gestalten. Man kann die beiden Achsen gemeinsam (durch **axis.ticks**) oder separat (durch **axis.ticks.x** bzw. **axis.ticks.y**) ansprechen. Im folgenden Beispiel wird für beide Achsen die Linienstärke der Teilstriche über das **element_line()** - Argument **size** auf den Wert 2 gesetzt:

```
> plot + theme(axis.ticks = element_line(size = 2))
```

- **axis.ticks.length**

Über ein **unit**-Objekt lässt sich die Länge der Achsenteilstriche setzen, z.B. auf 2 Millimeter:

```
> plot + theme(axis.ticks.length = unit(2, "mm"))
```

- **legend.position**

Zur Wahl einer Legendenposition können dem Argument **legend.position** der **theme()** - Funktion folgende Werte zugewiesen werden:

- Mit **"left"**, **"right"**, **"bottom"** oder **"top"** setzt man die Legende neben das Diagramm. z.B.:

```
> plot + theme(legend.position="bottom")
```

- Über den Wert **"none"** schaltet man die Legende ab.

- Über einen numerischen Vektor mit einer X- und einer Y-Position im Intervall [0, 1] erhält man eine Legende innerhalb des Diagramms, z.B.:

```
> plot + theme(legend.position=c(0.1, 0.9))
```

Über die Funktion **element_blank()** lässt sich die Ausgabe bestimmter Elemente unterdrücken, um z.B. ein Diagramm ohne Gitterlinien zu erstellen:

```
> plot + theme(panel.grid = element_blank())
```

Ein modifiziertes Theme lässt sich zur späteren Anwendung auf Diagramme in einem **R**-Objekt speichern, z.B.:

```
> redTitle <- theme_grey() + theme(plot.title = element_text(colour = "red"))
> histo + redTitle
```

8.3.2 ggplot2 - Diagramm in eine Datei sichern

Zum Sichern eines **ggplot2** - Diagramms in eine Datei steht neben den in Abschnitt 8.1 beschriebenen Optionen die **ggsave()** - Funktion zur Verfügung. Im ersten Parameter gibt man den Dateinamen an, wobei das Dateiformat (bzw. das zu verwendende Ausgabegerät) automatisch aus der Namenerweiterung abgeleitet wird, z.B.:

```
> ggsave("kfa.png")
```

Man kann den Namen des zu sichernden Diagramms weglassen, wenn das zuletzt angezeigte Diagramm gemeint ist.

Über die optionale Argumente **width** und **height** lassen sich Breite und Höhe in der Einheit Zoll (engl.: *Inch*) festlegen (1 Zoll = 2,54 cm), z.B.:

```
> ggsave("kfa.png", width=5, height=5)
```

Per **units**-Argument lässt sich auch eine alternative Maßeinheit verwenden, z.B. cm:

```
> ggsave("kfa.png", width=12, height=12, units="cm")
```

Bei einem Bitmap-Format (z.B. PNG) kann die Auflösung über das Argument **dpi** festgelegt werden, z.B.:

```
> ggsave("kfa.png", width=12, height=12, units="cm", dpi=600)
```

Ein Bitmap-Format mit einer Auflösung von 600 dpi ist z.B. dann zu empfehlen, wenn ein Diagramm unter Windows an ein Textverarbeitungsprogramm (Libre-, MS- oder Open-Office) übergeben werden soll (siehe Wickham 2009, Abschnitt 8.3).

Das unter Windows populäre Metafile-Format EMF ist für **ggplot2**-Grafiken *nicht* geeignet:

- Es unterstützt keine Transparenz, so dass z.B. die von **geom_smooth()** erstellten Konfidenzintervalle (siehe Abschnitt 8.3.3) verloren gehen.
- Im Vergleich zu anderen Vektorformaten (z.B. SVG) sind Kurven sehr grob aufgelöst.

Das traditionelle **R**-Verfahren zur Grafikausgabe in eine Datei ist z.B. dann gegenüber der Funktion **ggsave()** zu bevorzugen, wenn mehrere Diagramme auf einzelnen Seiten einer PDF-Datei abgelegt werden sollen.

8.3.3 Streudiagramme

In diesem Abschnitt wird zunächst über Zwischenschritte ein gruppiertes Streudiagramm erstellt. Dabei werden die in Abschnitt 8.3.1 zur Illustration des logischen Aufbaus eines **ggplot2**-Plots etwas verstreut präsentierten Informationen über Streudiagramme in einem Arbeitsablauf präsentiert. Wir greifen wieder auf die schon mehrfach (z.B. in Abschnitt 8.3.1) verwendeten Beispieldaten mit den Variablen **geschlecht**, **größe** und **gewicht** zurück. Später werden noch weitere Streudiagrammvarianten vorgestellt.

8.3.3.1 Plot-Objekt anlegen

Wir legen zunächst durch einen Aufruf der Funktion **ggplot()** ein Plot-Objekt an,

```
> sd <- ggplot(ggg, aes(größe, gewicht))
```

und vereinbaren dabei ...

- im ersten Argument (Name: **data**) die Datentabelle mit den zu visualisierenden Variablen
- und im zweiten Argument (Name: **mapping**) per **aes()** - Aufruf eine Plot-globale Verknüpfung von ästhetischen Attributen mit Variablen: Die Attribute **x** und **y** werden auf die Variablen **größe** und **gewicht** abgebildet.

Die Namen der Argumente können entfallen, weil die Wertvergabe in Definitionsreihenfolge erfolgt (Positionsargumente).

8.3.3.2 Einfaches Streudiagramm

Bevor nicht mindestens eine Schicht erstellt worden ist, kann das eben erstellte Plot-Objekt nicht angezeigt werden. Daher machen wir uns daran, Diagrammschichten durch **geom**-Funktionsaufrufe zu erzeugen und per **+** - Operator mit dem Plot-Objekt zu verknüpfen. Das folgende Kommando ergänzt eine Schicht mit den (x, y) - Datenpunkten und zeigt das Diagramm (durch einen impliziten **print()** - Aufruf) an:

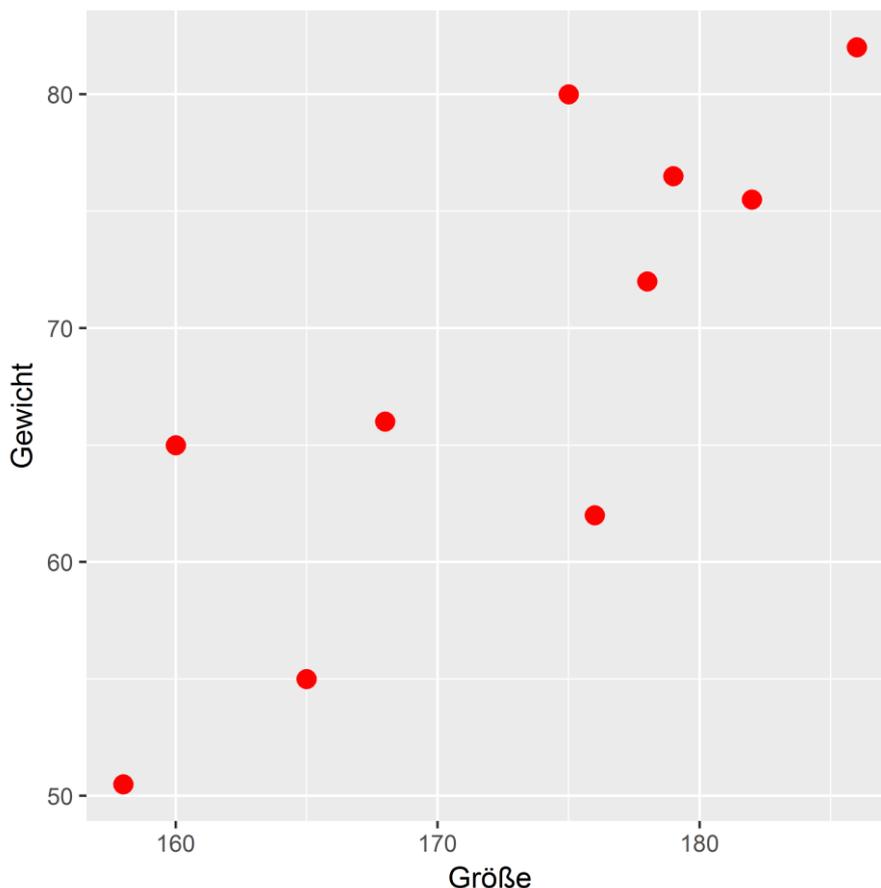
```
> sd + geom_point(colour="red", size=3)
```

Für die Datenpunkte wird die Zeichenfarbe Rot gewählt und außerdem der Durchmesser erhöht. Dabei wird jeweils ein ästhetisches Attribut auf einen festen Wert (statt auf eine Variable) abgebildet, und die Gültigkeit beschränkt sich auf die aktuelle Schicht.

Neben den **geom**-Funktionen, die jeweils eine neue Schicht anlegen, kennt das **ggplot2**-Paket Funktionen für Detailänderungen, die ebenfalls per **+** - Operator auf ein Plot-Objekt angewendet werden. Im folgenden Beispiel werden die Achsenbeschriftungen per **labs()** - Funktion modifiziert:

```
> sd + geom_point(colour="red", size=3) + labs(x = "Größe", y = "Gewicht")
```

Das Ergebnis:



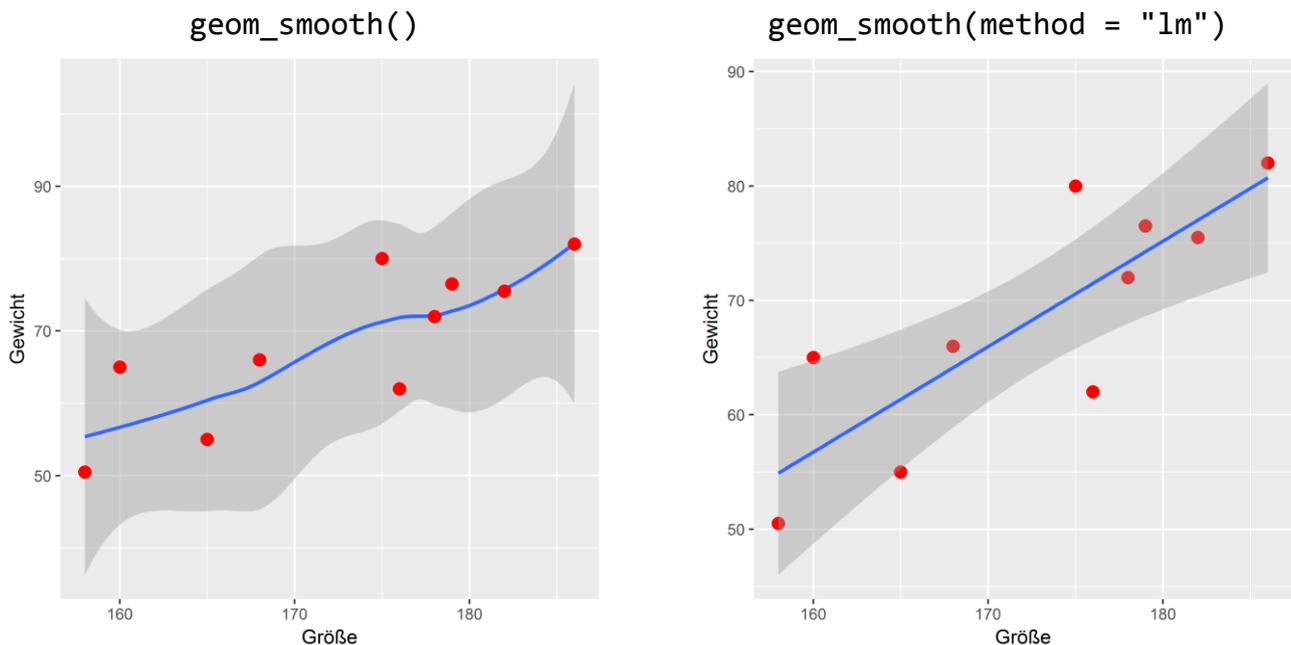
Alle Diagramme im Abschnitt 8.3.3 wurden mit dem folgenden Kommando (vgl. Abschnitt 8.3.2) in eine PNG-Datei geschrieben

```
ggsave("dia.svg", width=12, height=12, units="cm", dpi=600)
```

und anschließend von dort in den Text importiert. So entsteht im Vergleich zu den **ggplot2** - Abbildungen in früheren Abschnitten eine bessere Qualität.

8.3.3.3 Einfaches Streudiagramm mit Konfidenzzone

Um das Diagramm um eine 95% - Konfidenzzone zu erweitern, „addieren“ wir eine Schicht mit einem Glättungs-Geom. Die zuständige Funktion **geom_smooth()** beherrscht unterschiedliche Glättungsverfahren, z.B. die lokal optimierte Anpassung (Aufruf ohne Argument) und die Anpassung einer linearen Funktion (Aufruf mit dem Wert **"lm"** für das Argument **method**):¹



Sollen die Markierungspunkte *über* der (partiell transparenten) Konfidenzzone liegen und nicht überlagert werden (siehe Negativbeispiel im rechten Diagramm), ist auf die richtige Reihenfolge der **geom**-Aufrufe zu achten. z.B.:

```
> sd + geom_smooth()+geom_point(colour="red",size=3)+labs(x="Größe",y="Gewicht")
```

8.3.3.4 Gruppiertes Streudiagramm

Um ein *gruppiertes* Streudiagramm zu erzielen, legen wir per **ggplot()** - Aufruf ein neues Plot-Objekt an und machen dabei die Markierungsfarbe von der Variablen **geschlecht** abhängig. Dazu wird per **aes()** - Aufruf das ästhetische Attribut **colour** an die Variable **geschlecht** gebunden:

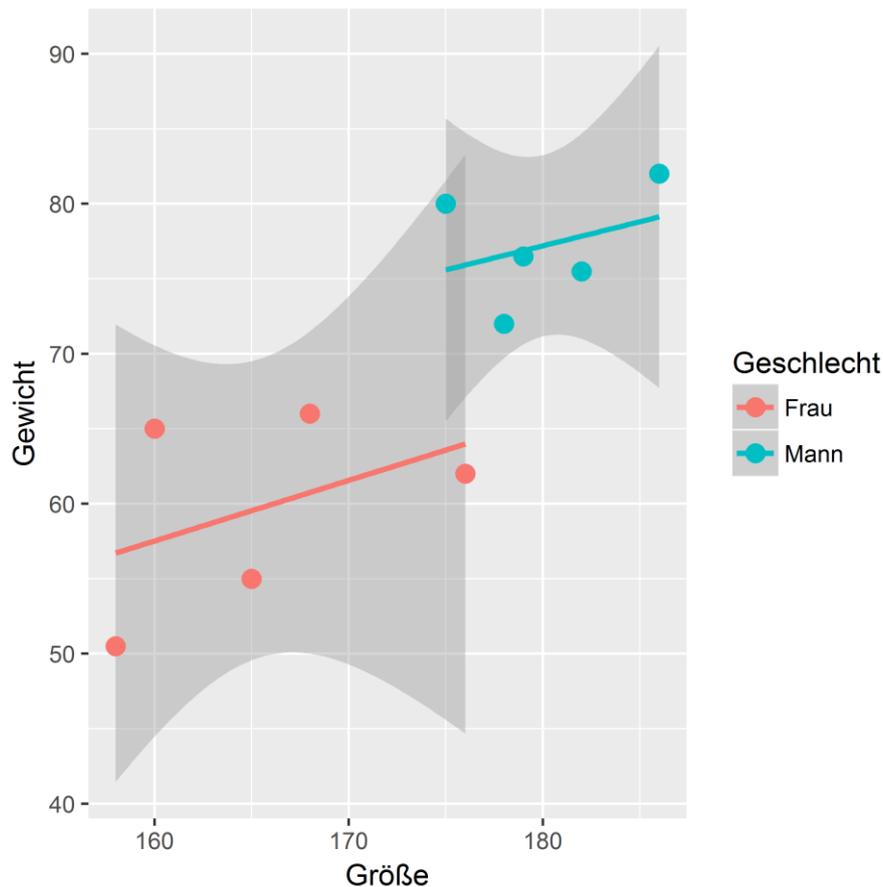
```
> gsd <- ggplot(ggg, aes(größe, gewicht, colour=geschlecht))
```

Mit dem sukzessiven Diagrammaufbau

```
> gsd + geom_smooth(method="lm") + geom_point(size=3)
+ labs(x="Größe", y="Gewicht", colour="Geschlecht")
```

¹ Zur Übernahme via Zwischenablage in Microsoft Word wurde das Bitmap-Format benutzt, weil beim Metafile-Transfer die Konfidenzzone aufgrund der Transparenzdarstellung verloren geht.

erhalten wir das folgende Ergebnis, wobei auch die Anpassungsfunktion für das Glättungs-Geom auf die Gruppierung reagiert:



Um die Legendenbeschriftung zu ändern, wurde im **labs()** - Aufruf dem Argument (bzw. dem ästhetischen Attribut) **colour** die gewünschte Zeichenfolge zugewiesen.

Im aktuellen Fall erscheint es *nicht* sinnvoll, auf eine Legende zu verzichten. Trotzdem soll demonstriert werden, wie dies mit einer „additiv“ ergänzten Theme-Modifikation (vgl. Abschnitt 8.3.1.7.2) möglich ist:

```
> gsd + geom_smooth(method="lm") + geom_point(size=3)
+ labs(x="Größe", y="Gewicht") + theme(legend.position="none")
```

Wie man einzelne Geome aus der Legende fernhält, wurde in Abschnitt 8.3.1.4 beschrieben. Im Beispiel kann man auch durch eine doppelte Verzichtserklärung (**show.legend=FALSE**) für das komplette Verschwinden der Legende sorgen:

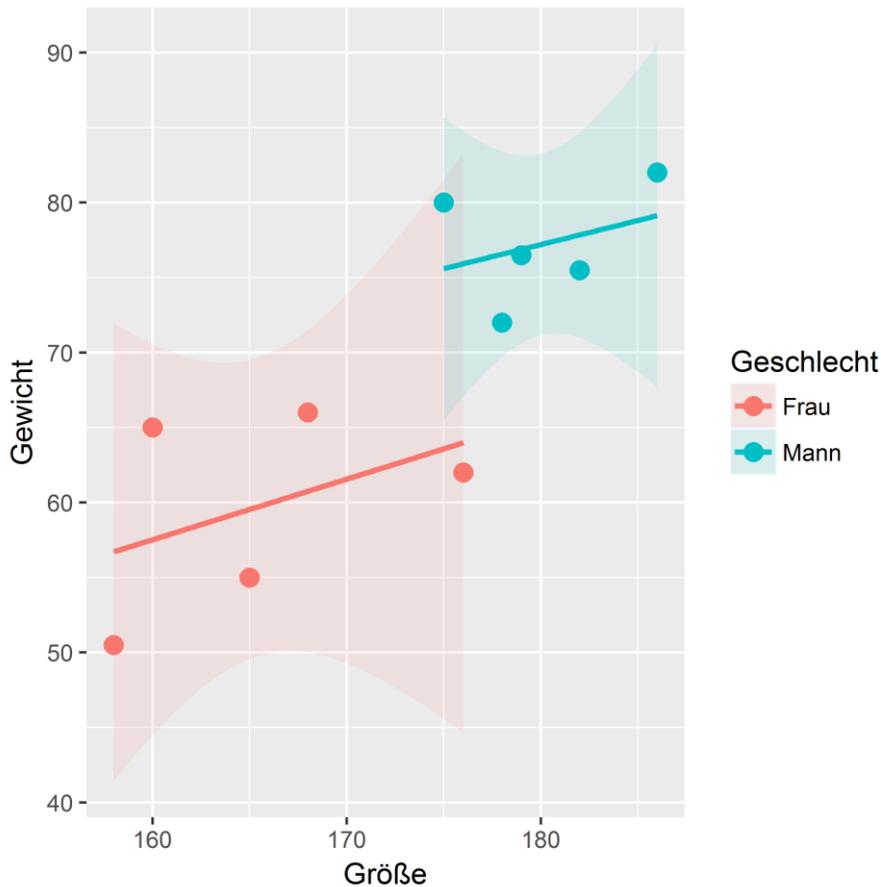
```
> gsd + geom_smooth(method="lm", show.legend=FALSE)
+ geom_point(size=3, show.legend=FALSE) + labs(x="Größe", y="Gewicht")
```

Im folgenden Kommando nach einem Vorschlag von Field (2012, S. 141) werden die Konfidenzzonen Geschlechts-abhängig gefärbt und außerdem mit einem stärkeren Transparenzgrad versehen:¹

¹ Weil auch das für die Linien- und Markierungsfarbe zuständige und per Voreinstellung auf die Variable **geschlecht** abgebildete ästhetische Attribut **colour** einen Legendenbeitrag bewirkt, muss mit der Skalenfunktion **scale_colour_hue("Geschlecht")** sein Legendentitel mit dem Legendenbeitrag das **fill**-Attributs harmonisiert werden, weil sonst *zwei* Legenden erscheinen.

```
> gsd + geom_smooth(method="lm", aes(fill=geschlecht), alpha=0.1) + geom_point(size=3)
+   scale_fill_hue("Geschlecht") + scale_colour_hue("Geschlecht")
+   labs(x="Größe", y="Gewicht")
```

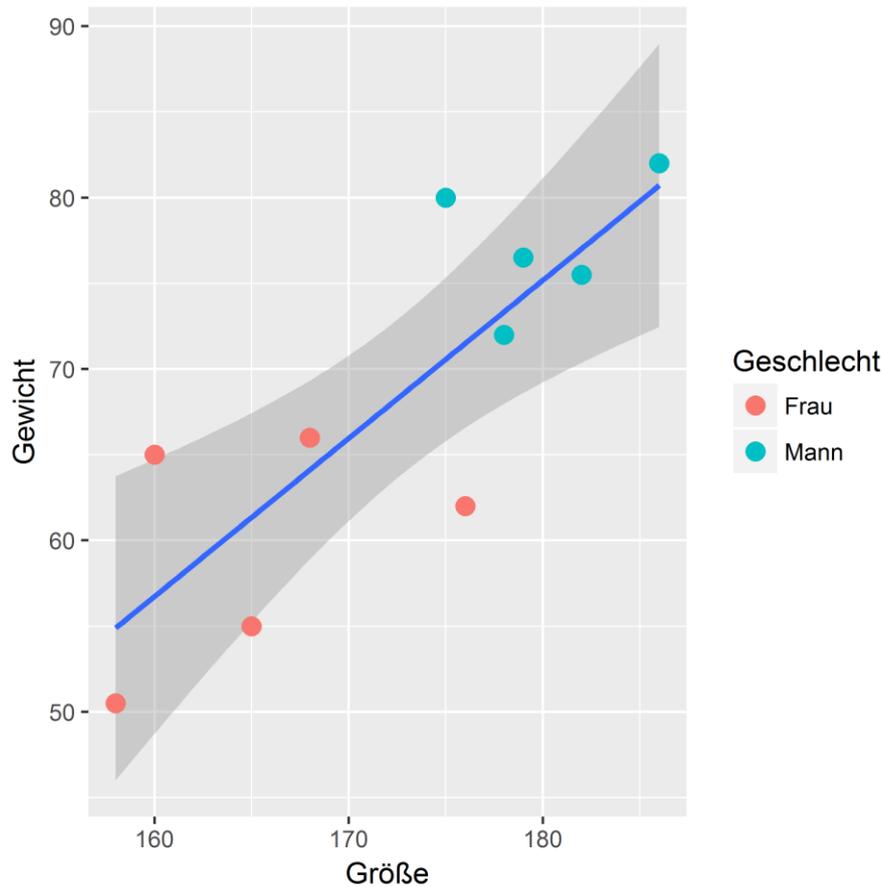
Ergebnis:



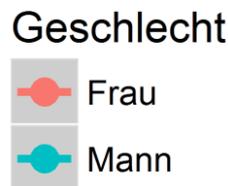
Wenn sich die gruppenspezifischen Anpassungsfunktionen kaum unterscheiden, kommt eine *gemeinsame* Anpassungsfunktion in Frage, wobei es aber in der Regel trotzdem sinnvoll ist, die Gruppen farblich zu unterscheiden. Sobald auf Plot-Ebene dem ästhetischen Attribut **colour** eine Variable zugewiesen ist, besteht eine Gruppierung, die sich auch auf das Glättungs-Geom auswirkt. Um dies zu verhindern, muss im **geom_smooth()** - Aufruf das ästhetische Attribut **group** auf den konstanten Wert 1 abgebildet werden:

```
> gsd <- ggplot(ggg, aes(größe, gewicht, colour=geschlecht))
> gsd + geom_smooth(method="lm", aes(group=1), show.legend=FALSE)
+   geom_point(size=3) + labs(x="Größe", y="Gewicht", colour="Geschlecht")
```

Hier ist das gewünschte Ergebnis:



Um die unpassende Legende



durch eine Alternative ohne gruppenspezifische Linienfarben zu ersetzen, wird im `geom_smooth()` - Aufruf mit dem Wert **FALSE** für das Argument **show.legend** verhindert, dass die Schicht in der Legende Berücksichtigung findet.

8.3.3.5 Schichtaufbau mit `qplot()` starten

Am grundsätzlichen Schichtaufbau eines `ggplot2`-Diagramms ändert sich übrigens nichts, wenn statt `ggplot()` die Funktion `qplot()` verwendet wird. Ein `qplot()` - Aufruf erstellt ein Plot-Objekt und ergänzt Schichten. Man kann sogar eine `qplot()` - Produktion als Ausgangsbasis für den weiteren Schichtaufbau mit `geom`-Funktionen verwenden, z.B.:

```
> qplot(gewicht, größe, data = ggg) + geom_smooth(method="lm")
```

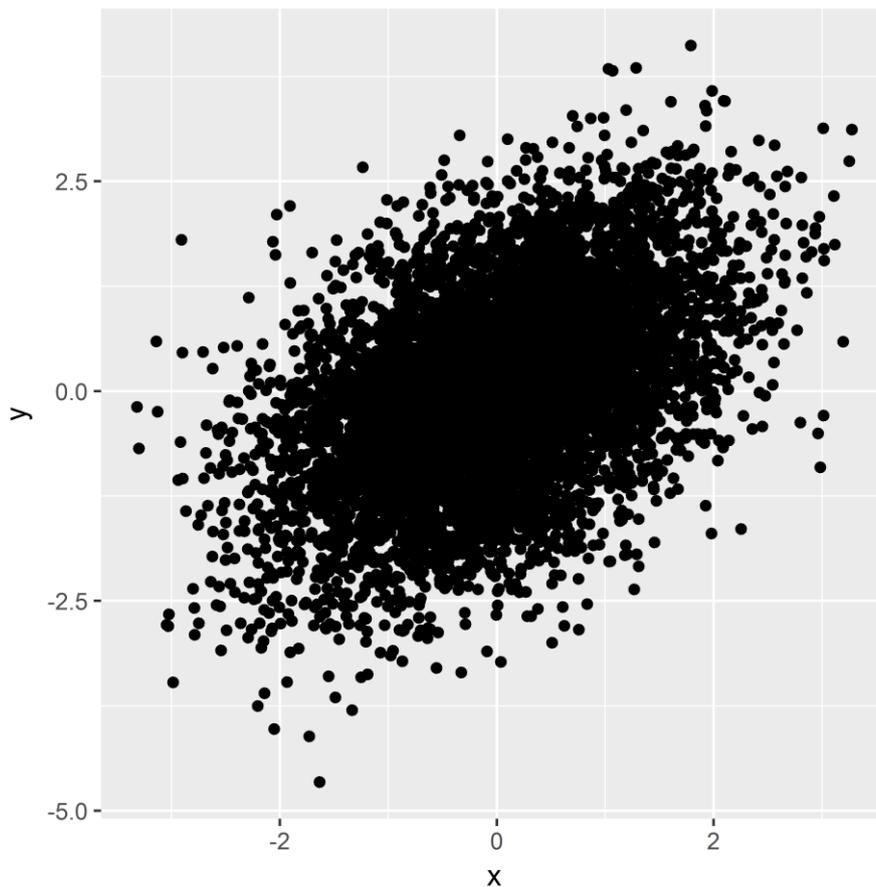
Im Beispiel landet allerdings die Glättungsschicht mit der transparenten Konfidenzzone *über* der Punkteschicht, was die Färbung einiger Punkte ändert. Die im Vergleich zu `qplot()` deutlich größere Flexibilität der Funktion `ggplot()` bei der Grafikproduktion zeigt sich u.a. in der perfekten Kontrolle über den Schichtaufbau.

8.3.3.6 Dichtedarstellung bei großen Stichproben

Wir haben schon in Abschnitt 8.2.5.2.6 im Zusammenhang mit der traditionellen **R**-Grafik Darstellungstechniken für Streudiagramme mit sehr vielen Fällen kennengelernt. Wie das folgende Beispiel mit simulierten Daten

```
> n <- 7000
> set.seed(12)
> x <- rnorm(n,0,1)
> res = rnorm(n,0,1)
> y <- 0.5* x + res
> df <- data.frame(x,y)
> ggplot(df, aes(x,y)) + geom_point()
```

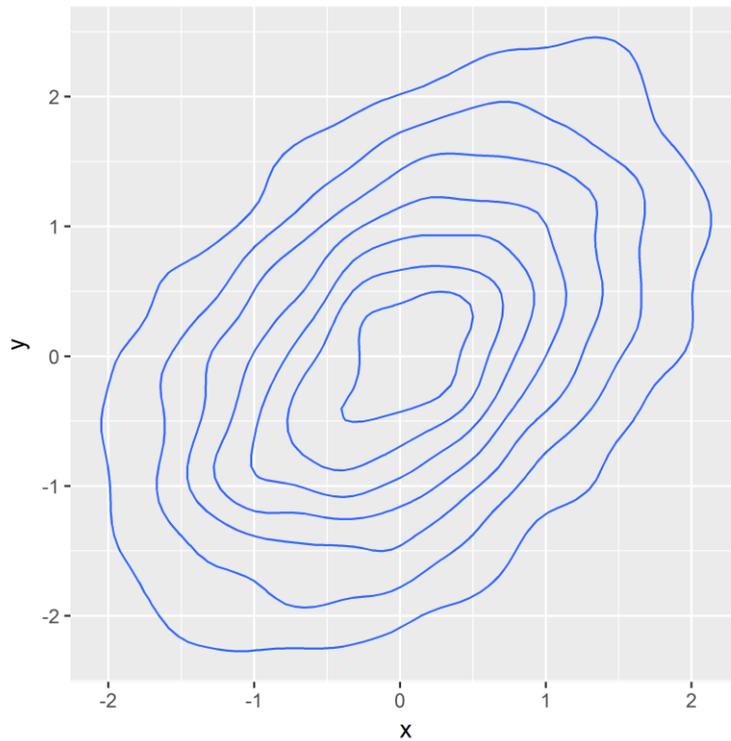
zeigt, ist die voreingestellte Darstellung von `geom_point()` nicht ideal:



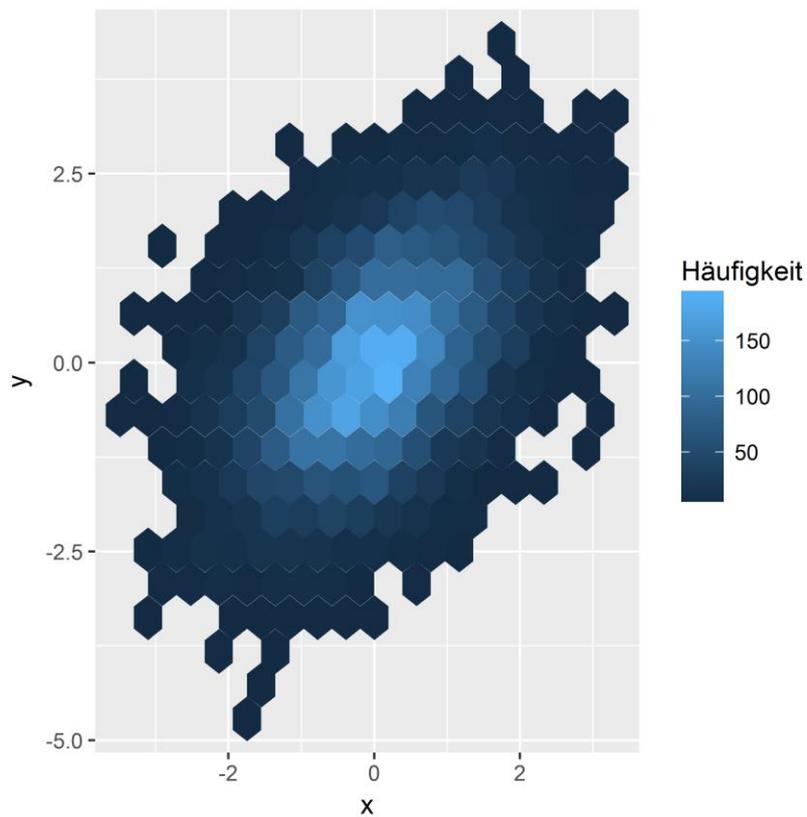
Von der Funktion `geom_density2d()`

```
> ggplot(df, aes(x,y)) + geom_density2d()
```

erhält man die **Dichtelinien** der Verteilung:



Eine in **ggplot2** (nach der Zusatzinstallation des Pakets **hexbin**) sowie auch in SPSS (via GPL) verfügbare Alternative ist das Streudiagramm mit **hexagonaler Gruppierung**. Dieses Exemplar



wurde durch folgende Syntax erzeugt:

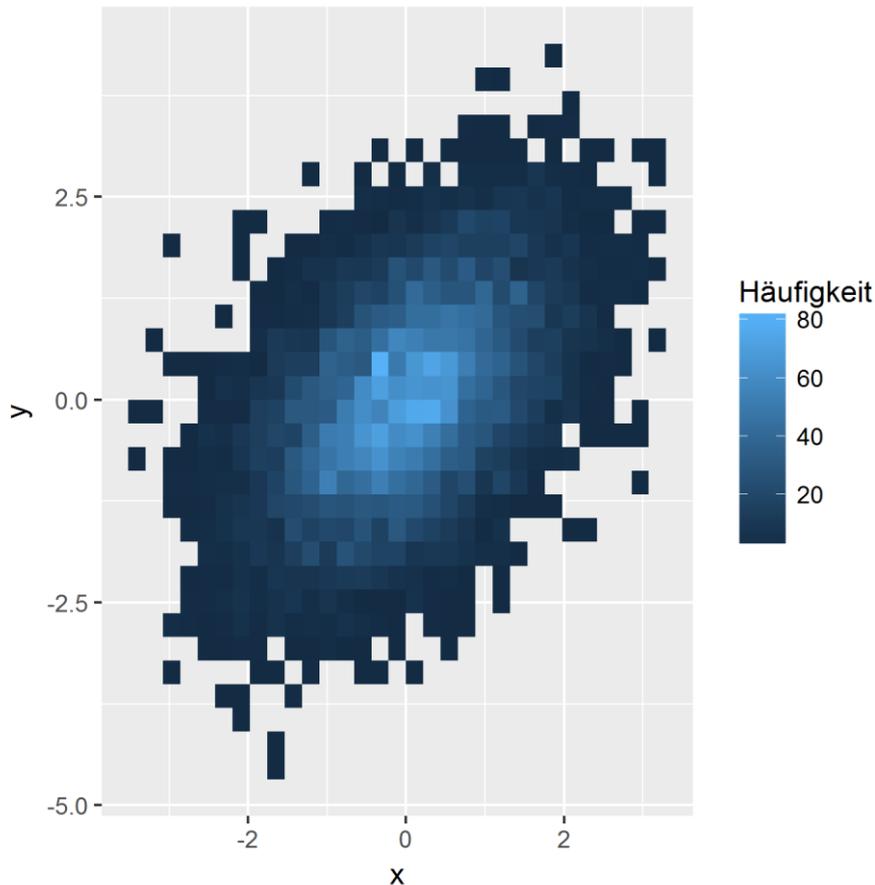
```
> library(hexbin)
> ggplot(df, aes(x,y))+stat_binhex(bins=17)+scale_fill_continuous(name='Häufigkeit')
```

Leider wird die Darstellungsqualität durch unregelmäßig auftretende horizontale Sägezahnmuster beeinträchtigt. Diese lassen sich über die Anzahl der Hexagone (`stat_binhex()` - Argument `bins` mit der Voreinstellung 30) beeinflussen, aber nicht zuverlässig beseitigen.

Gibt man sich durch mit Quadraten anstelle von Hexagonen zufrieden und verwendet die Funktion `geom_bin2d()` anstelle der Funktion `stat_binhex()`,

```
> ggplot(df, aes(x, y)) + geom_bin2d() + scale_fill_continuous(name='Häufigkeit')
```

dann resultiert eine korrekte Darstellung:



8.3.4 Weitere Diagrammtypen

Zur Demonstration des Histogramms verwenden wir die SPSS-Datendatei `kfa.RData`, die sich an dem im Vorwort vereinbarten Ort befindet. Aus dem folgenden Funktionsaufruf resultiert die Datentabelle `kfa`:

```
> load("kfa.RData")
```

8.3.4.1 Histogramm und Dichteschätzung

In der Datentabelle `kfa` befindet sich u.a. die Variable `aergo` mit dem auf einer Skala von 0 bis 10 gemessenen Ärger von 31 Probanden über einen verpassten Flug. Auf dem Weg zu einem Histogramm für `aergo` erstellen wir mit der Funktion `ggplot()` ein Plot-Objekt und vereinbaren per `aes()` - Funktion die darzustellende Variable:

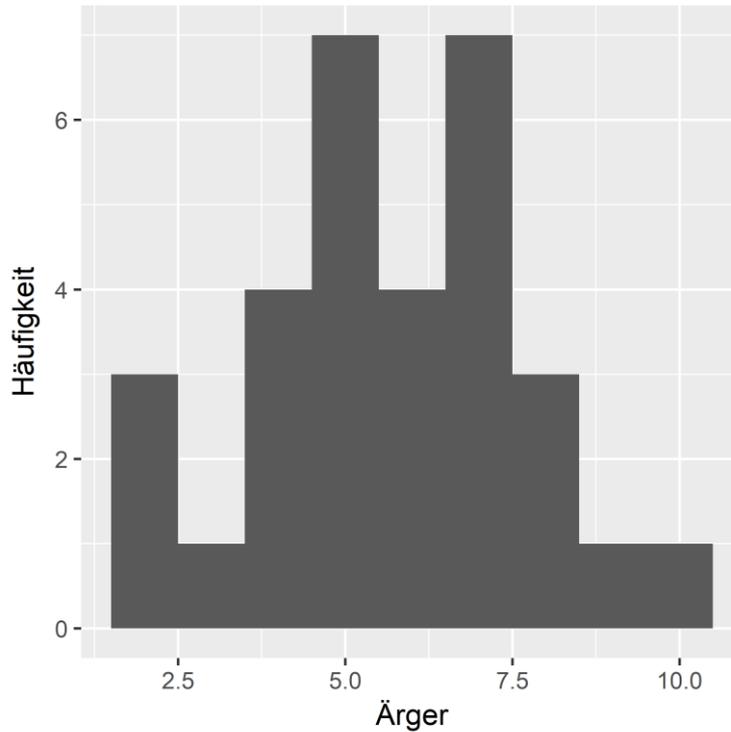
```
> histo <- ggplot(kfa, aes(aergo))
```

Wir ergänzen die Schicht mit dem Histogramm und sorgen im zuständigen `geom_histogram()` - Funktionsaufruf über das Argument `binwidth`, das an Transformationsfunktion `stat_bin()` durchgereicht wird

(vgl. Abschnitt 8.3.1.3), für eine passende Intervallbreite. Außerdem werden die Achsenbeschriftungen per `labs()` - Funktion modifiziert:

```
> histo + geom_histogram(binwidth = 1) + labs(x="Ärger", y="Häufigkeit")
```

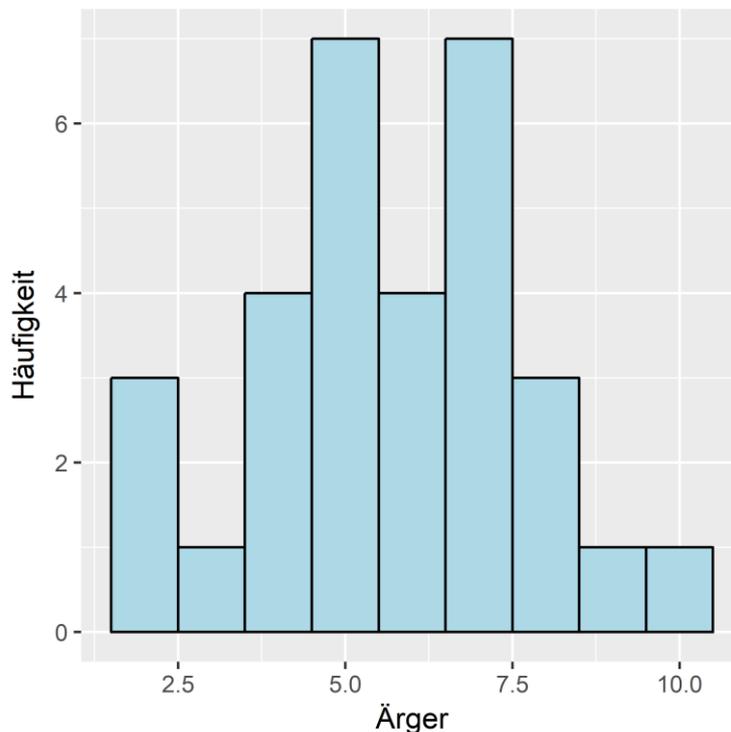
Im Ergebnis stört die triste Balkenfarbe:



Mit dem folgenden Kommando

```
> histo+geom_histogram(binwidth=1,colour="black",fill="lightblue")+labs(x="Ärger",y="Häufigkeit")
```

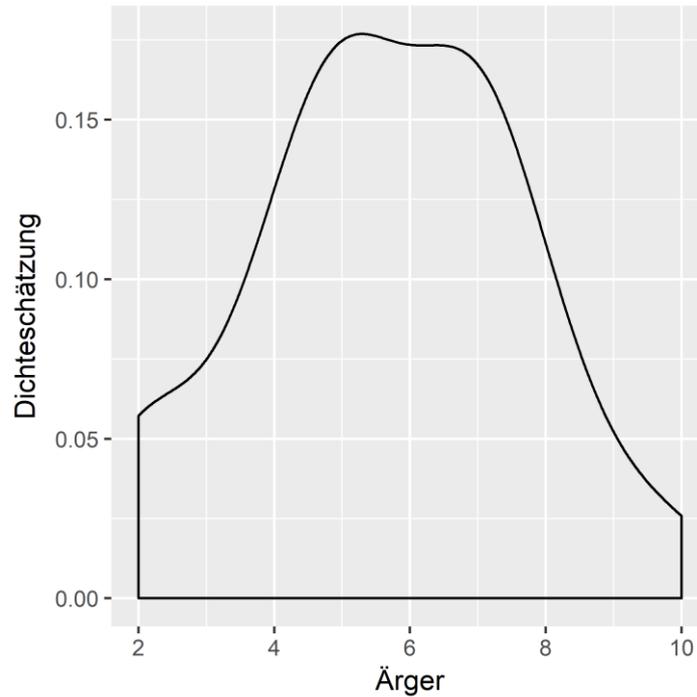
gelingt eine freundlichere Farbgestaltung (schwarz umrandete hellblaue Balken):



Über die Funktion `geom_density()` kann man die univariate Verteilung einer metrischen Variablen durch eine geschätzte Dichtefunktion beschreiben lassen.

```
> dichte <- ggplot(kfa, aes(aergo))  
> dichte + geom_density() + labs(x="Ärger", y="Dichteschätzung")
```

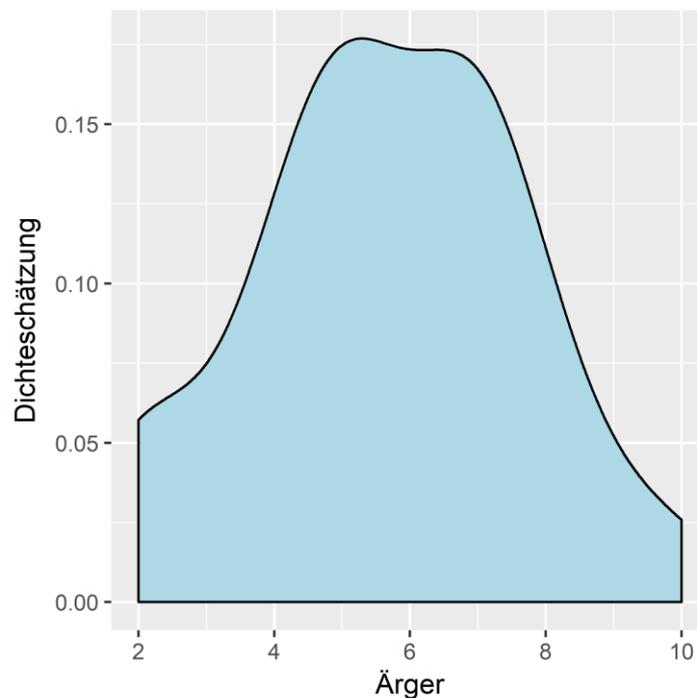
Im Beispiel resultiert eine „Mütze“:



Wird im Aufruf der Geom-Funktion das ästhetische Attribut `fill` auf einen konstanten Farbwert gesetzt,

```
> dichte + geom_density(fill="lightblue") + labs(x="Ärger", y="Dichteschätzung")
```

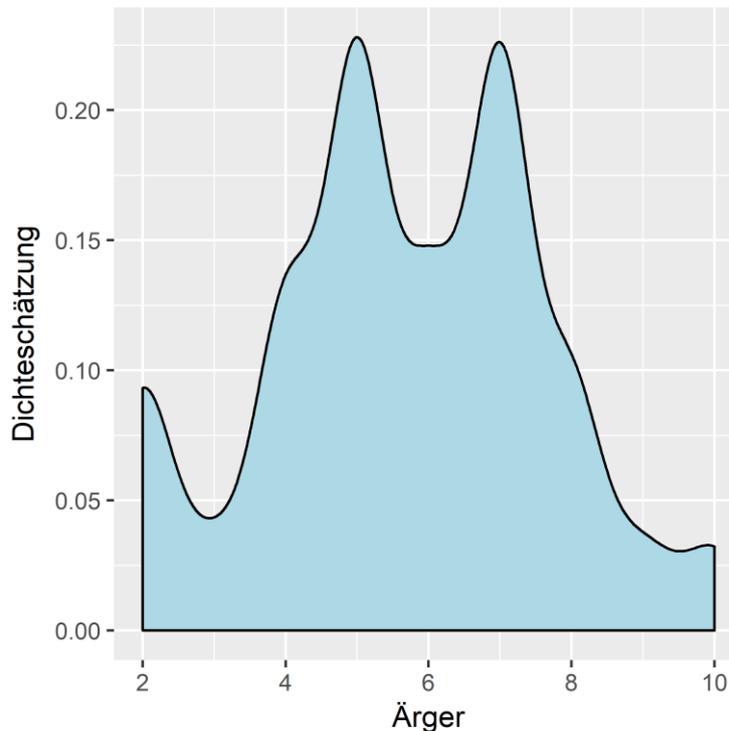
geht es etwas bunter zu:



Zur Steuerung des Glättungsgrads besitzt die Funktion `geom_density()` das Argument `adjust` mit dem Voreinstellungswert 1. Mit dem alternativen Wert 0,5

```
> dichte+geom_density(adjust=0.5,fill="lightblue")+labs(x="Ärger",y="Dichteschätzung")
```

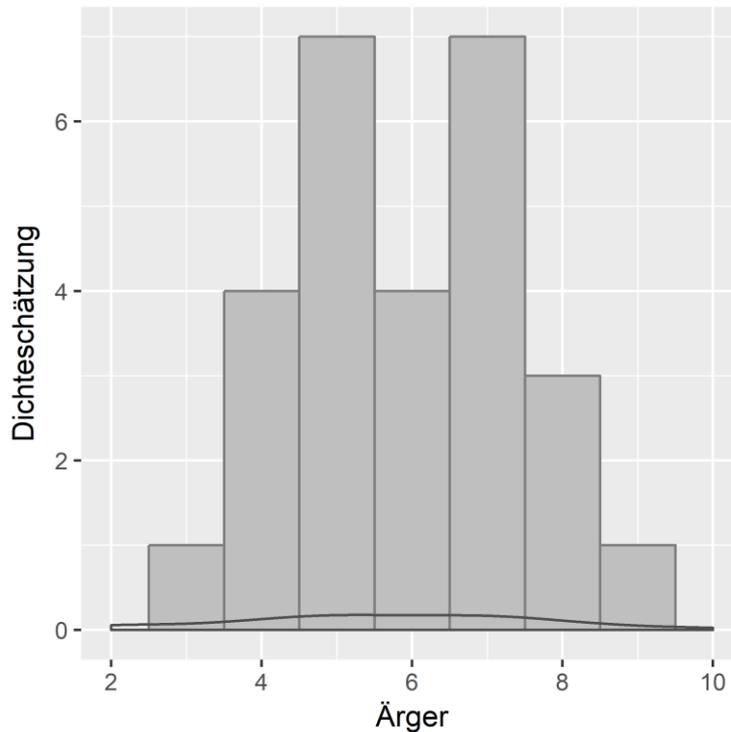
wird im Beispiel aus der Mütze eine Tiersilhouette:



Im folgenden Beispiel basierend auf einer Anregung aus Chang (2013, S. 125f) wird dem Histogramm eine Dichteschätzung überlagert. Von den Variablen, welche die bei `geom_histogram()` voreingestellte Transformation `stat_bin()` produziert, verwendet das Histogramm-Geom per Voreinstellung die Variable `count` mit den absoluten Häufigkeiten der Intervalle. Deren Werte übertreffen die Y-Werte von `geom_density()` (mit dem vorgeschriebenen Integrationsergebnis 1) bei weitem, so dass der erste Versuch

```
> histodichte <- ggplot(kfa, aes(aergo))
> histodichte + geom_histogram(binwidth=1,fill="gray75",colour="gray50")
+   geom_density(colour="gray30")+labs(x="Ärger",y="Dichteschätzung")+xlim(2,10)
```

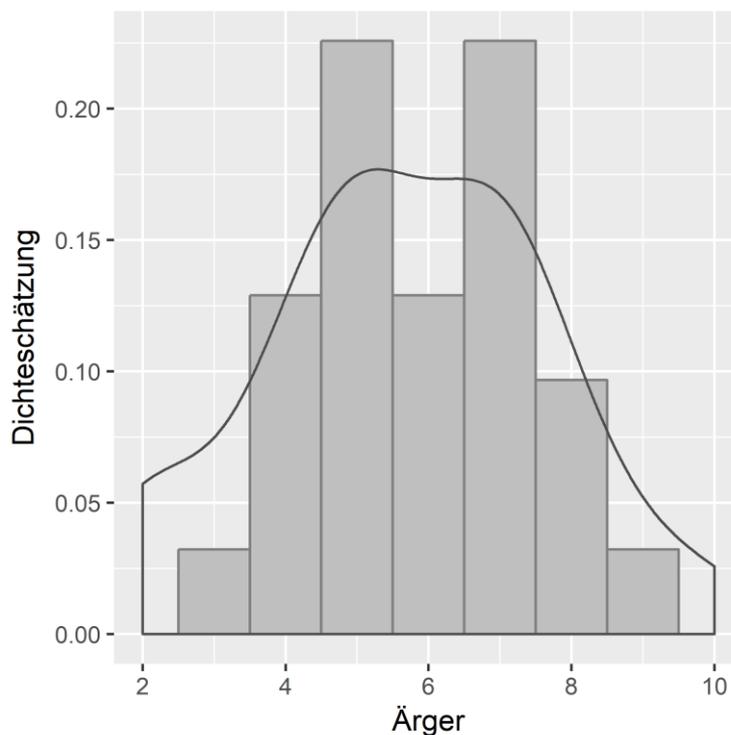
zu einem unbrauchbaren Ergebnis führt:



Zur Lösung des Problems wird dem `geom_histogram()` - Argument **mapping** per `aes()` - Aufruf eine Abbildung der Y-Achse auf die Ergebnisvariable **density** der involvierten statistischen Transformation `stat_bin()` zugewiesen (zur Syntax vgl. Abschnitt 8.3.1.3):

```
> histodichte <- ggplot(kfa, aes(aergo))
> histodichte+geom_histogram(aes(y=..density..),binwidth=1,fill="gray75",colour="gray50")
+ geom_density(colour="gray30")+labs(x="Ärger",y="Dichteschätzung")+xlim(2,10)
```

Durch Wahl passender Grauwerte kommen Histogramm und Dichtefunktion gut zur Geltung:



8.3.4.2 Boxplot

Zur Demonstration des Boxplots verwenden wir weiterhin die metrische Variable `aergo` in der Datentabelle `kfa` (siehe Abschnitt 8.3.4.1). Zunächst soll die `aergo`-Verteilung in der Gesamtstichprobe dargestellt werden. Wir erstellen mit der Funktion `ggplot()` ein Plot-Objekt und vereinbaren per `aes()` - Funktion die darzustellende Variable:

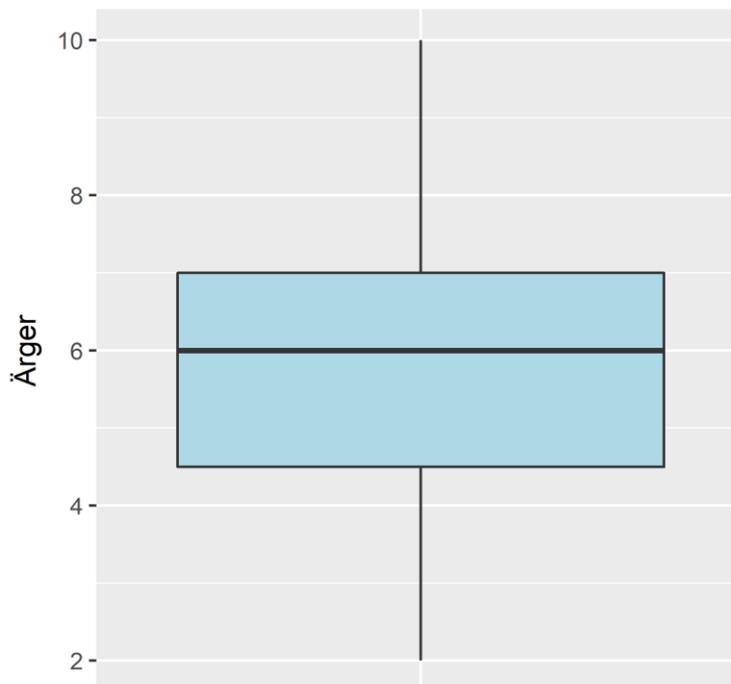
```
> box <- ggplot(kfa, aes(x=factor(0), y=aergo))
```

Das Boxplot-Geom erwartet X-Achsenvariable zur Aufteilung der Stichprobe und eine Y-Achsenvariable mit der zu beschreibenden Variablen. Weil wir die Gesamtstichprobe darstellen wollen, liefern wir als Wert für das `x`-Attribut über den Funktionsaufruf `factor(0)` einen Dummy-Faktor mit 0 Ausprägungen.

Wir ergänzen die Schicht mit dem Boxplot und sorgen im zuständigen `geom_boxplot()` - Funktionsaufruf über das Argument `fill` für eine angenehme Farbe. Mit der `labs()` - Funktion werden die Achsenbeschriftungen modifiziert, und der `theme()` - Aufruf unterdrückt die X-Achsenteilstriche sowie deren Beschriftungen:¹

```
> box + geom_boxplot(fill="lightblue") + labs(x="",y="Ärger")
+ theme(axis.ticks.x = element_blank(), axis.text.x = element_blank())
```

Das Ergebnis:

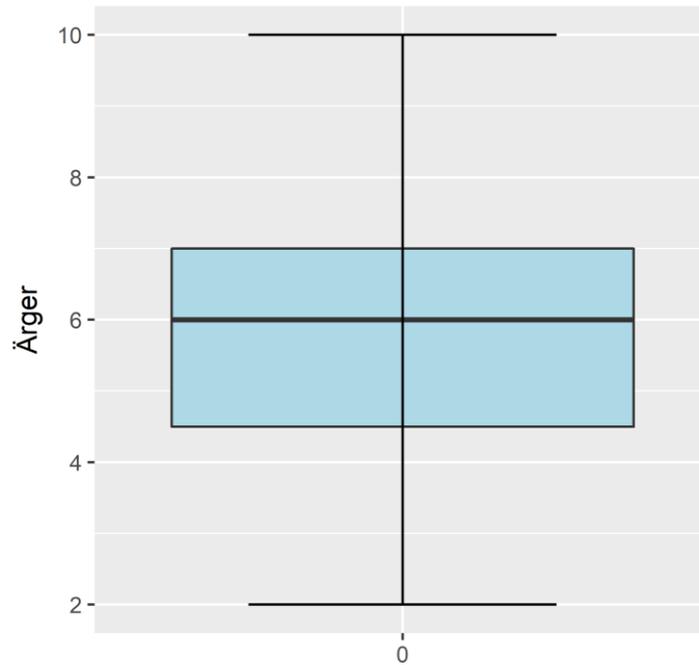


Die traditionelle **R**-Grafik stellt in einem Boxplot den kleinsten und den größten Wert in der Stichprobe, die keine Ausreißer sind, mit einer horizontalen Linie dar (siehe Abschnitt 8.2.5.3). Über eine durch die Funktion `stat_boxplot()`

```
> box <- ggplot(kfa, aes(x=factor(0), y=aergo))
> box + geom_boxplot(fill="lightblue")
+ stat_boxplot(geom = 'errorbar', width=0.5) + labs(x="",y="Ärger")
```

erstellte Zusatzschicht lässt sich diese Darstellung auch mit `ggplot2` erzielen:

¹ Anderenfalls würde ein überflüssiger Teilstrich mit der Beschriftung „0“ erscheinen.



Im Beispiel wird die Breite der horizontalen Linien durch das `stat_boxplot()` - Argument **width** festgelegt.

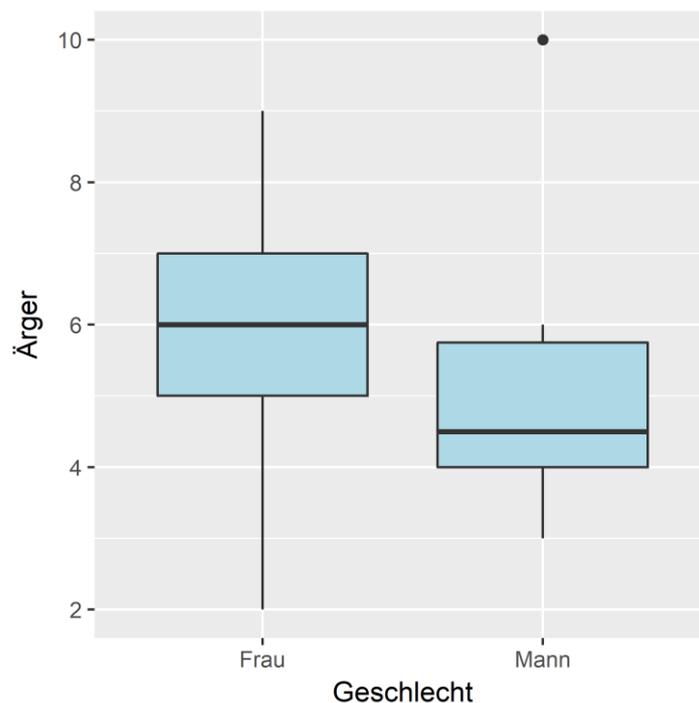
Sollen z.B. die Ärgerverteilungen bei Frauen und Männern gegenübergestellt werden, vereinbart man (z.B. bei der Erstellung des Plot-Objekts) eine passende Gruppierungsvariable, z.B.:

```
> box <- ggplot(kfa, aes(geschlecht, aergo))
```

Die Anweisung

```
> box + geom_boxplot(fill="lightblue") + labs(x="Geschlecht",y="Ärger")
```

liefert zwei nebeneinander stehende Boxplots, die einen Vergleich der Geschlechts-bedingten Verteilungen hinsichtlich Lage und Dispersion erlauben:



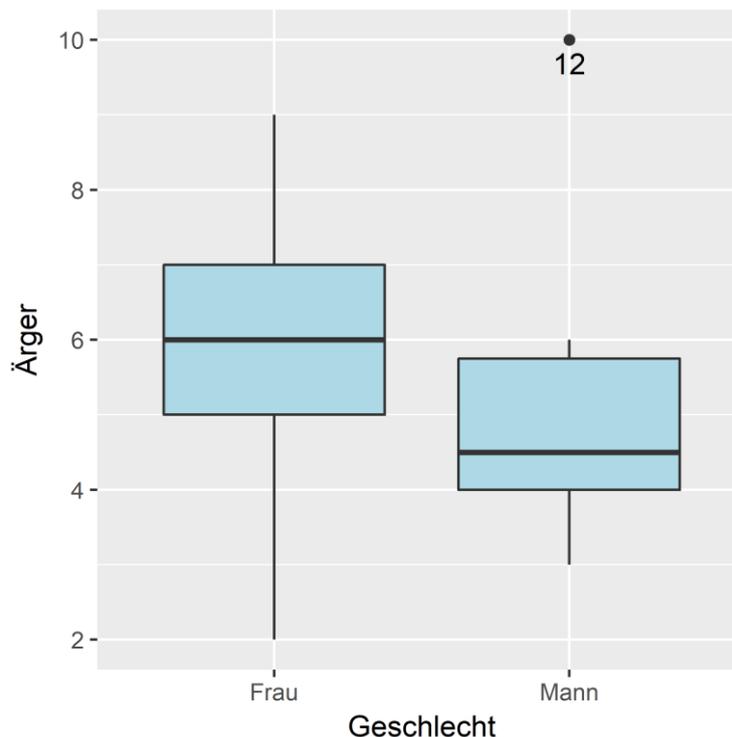
In der männlichen Teilstichprobe zeigt sich zudem ein Ausreißer mit dem Ärgerwert 10, der um mehr als 1,5 Boxbreiten vom 3. Quartil (dem oberen Rand der Box) entfernt ist. Um diesen Fall zu etikettieren, wird eine Variable vorbereitet, die bei Personen mit dem Ärgerwert 10 die Fallnummer (den Wert der Variablen `fnr`) enthält und ansonsten den Indikator für fehlende Werte:

```
> etikett <- kfa$fnr; etikett[kfa$aergo<10] <- NA
```

Nun wird eine Schicht mit einem `text`-Geom unter Verwendung der vorbereiteten Etikettierungsvariablen ergänzt:

```
> box + geom_boxplot(fill="lightblue") + labs(x="Geschlecht", y="Ärger")
+ geom_text(aes(label=etikett), vjust=1.6)
```

Mit der Eigenschaft `vjust` wird der Text in vertikaler Richtung vom Datenpunkt weg bewegt, um eine Überlagerung zu verhindern:



8.3.4.3 Balkendiagramme

Mit einem einfachen Balkendiagramm kann man darstellen:

- der Verteilung einer diskreten (z.B. kategorialen) Variablen,
- die Mittelwerte oder andere statistische Zusammenfassungen einer metrischen Variablen für die Ausprägungen einer kategorialen Variablen,
- die Mittelwerte von mehreren metrischen Variablen.

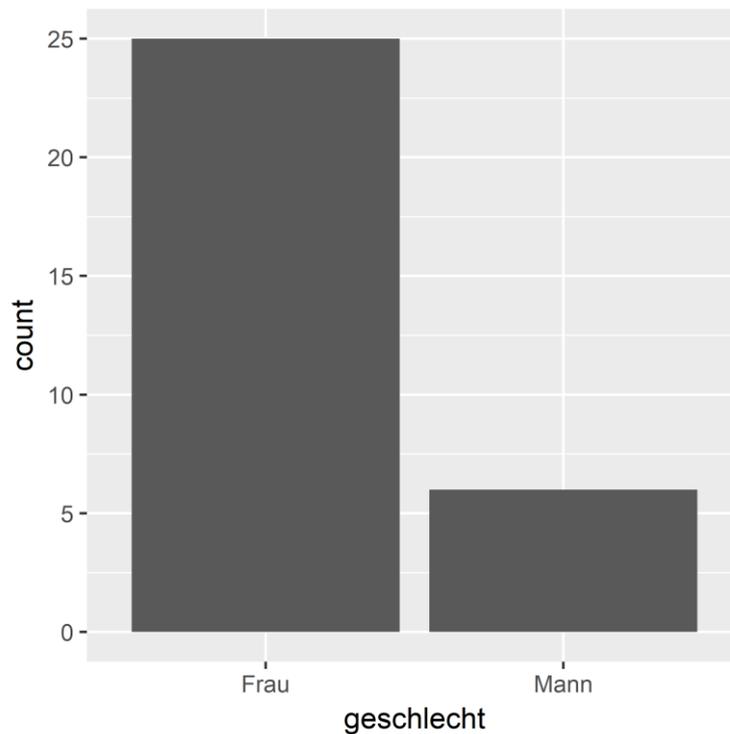
Zur Demonstration der verschiedenen Balkendiagramme verwenden wir weiterhin die Datentabelle `kfa`.

8.3.4.3.1 Diskrete Verteilungen

Über die Funktion `geom_bar()` lässt sich mit sehr wenig Aufwand

```
> ggplot(kfa, aes(x=geschlecht)) + geom_bar()
```

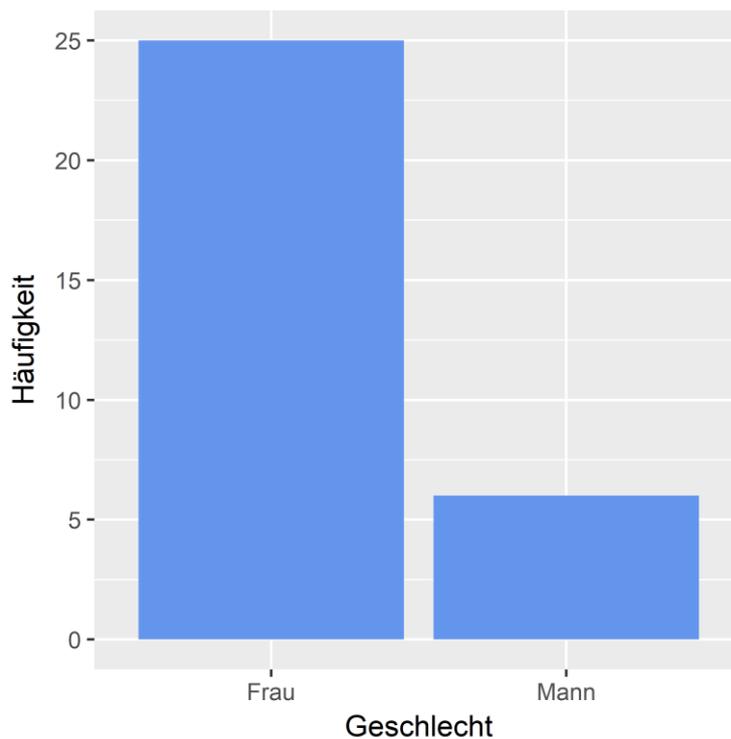
ein Balkendiagramm zur Anzeige der absoluten Häufigkeiten für die Kategorien einer diskret verteilten Variablen erstellen, z.B.:



Ein kleiner Zusatzaufwand

```
> ggplot(kfa, aes(x=geschlecht)) + geom_bar(fill="cornflowerblue")  
+ labs(x="Geschlecht", y="Häufigkeit")
```

erlaubt die wahlfreie Färbung und Achsenbeschriftung, z.B.:



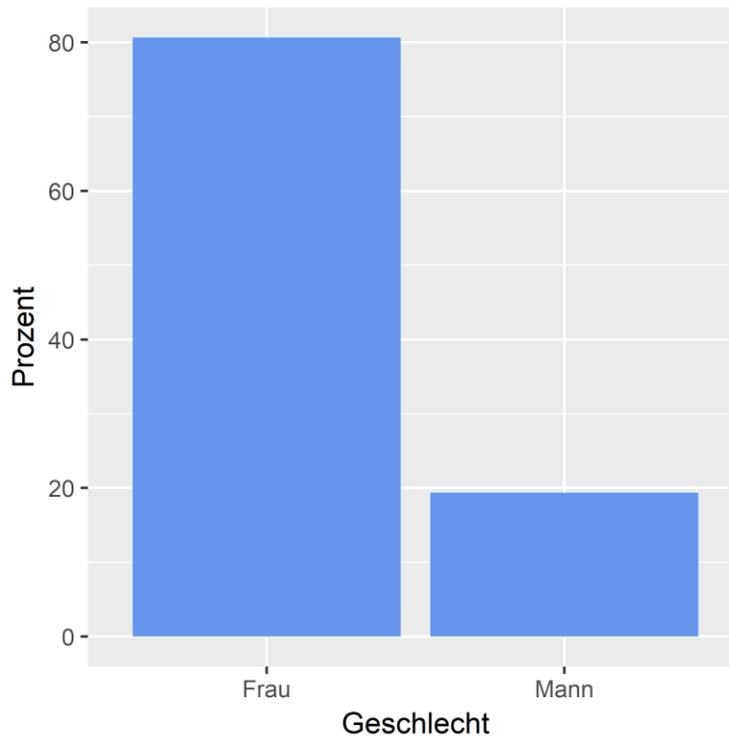
Indem das ästhetische Attribut **y** auf die folgende Funktion

```
(..count..)*100/sum(..count..)
```

der Transformations-Ergebnisvariablen **count** abgebildet wird,

```
> ggplot(kfa,aes(x=geschlecht,y=..count..*100/sum(..count..)))
+   geom_bar(fill="cornflowerblue")
+   labs(x="Geschlecht",y="Prozent")
```

erhält man die relativen Häufigkeiten (in Prozent):



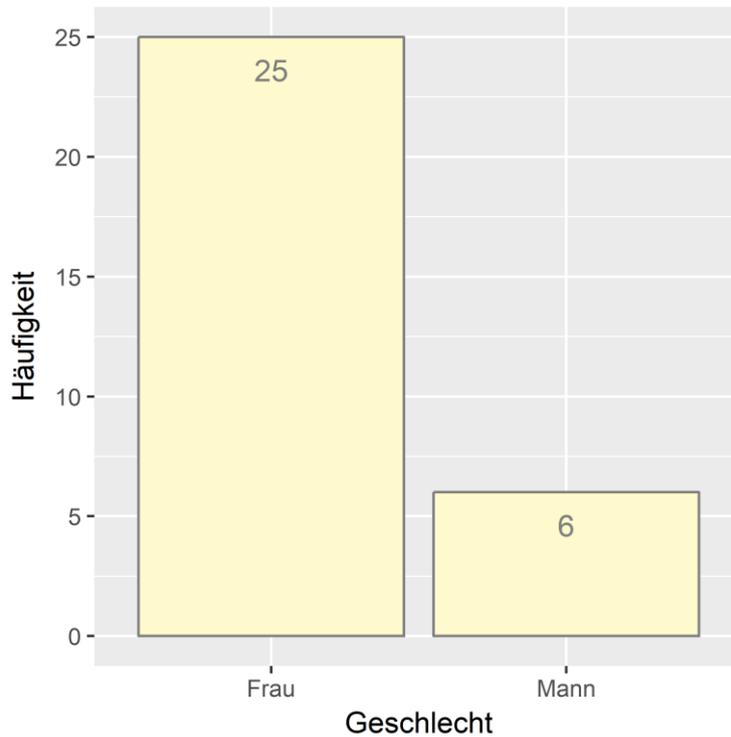
Über das **text**-Geom können die Balken mit den absoluten Häufigkeiten beschriftet werden:

```
> ggplot(kfa,aes(x=geschlecht)) + geom_bar(fill="lemonchiffon",colour="gray50")
+   labs(x="Geschlecht",y="Häufigkeit")
+   geom_text(stat="count",mapping=aes(label=..count..),vjust=2,size=4,colour="gray50")
```

Die bei **geom_text()** voreingestellte identische Transformation wird durch **count** ersetzt. Die somit verfügbare Transformationsergebnisvariable **..count..** wird mit dem ästhetischen Attribut **label** verbunden, das die auszugebenden Texte festlegt. Andere ästhetische Attribute von **geom_text()** erhalten einen festen Wert:

- **vjust** legt eine vertikale Verschiebung der Texte relativ zu den (durch **count** festgelegten) **y**-Positionen fest.
- **size** gibt die Textgröße an.
- **colour** bestimmt die Textfarbe.

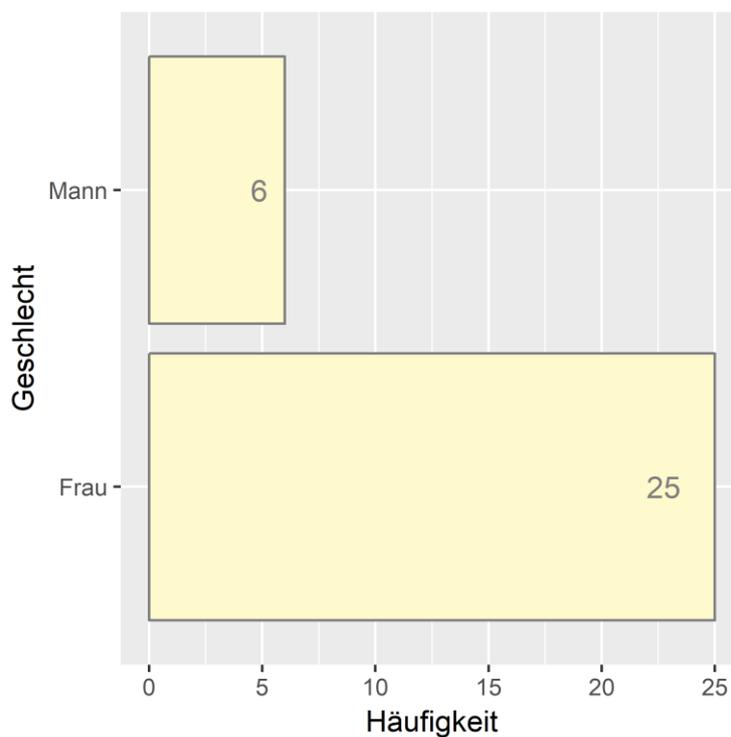
Auf der Schicht mit dem **text**-Geom erscheinen die gewünschten Etiketten mit den gewählten Attributen:



Um *waagerechte* Balken zu erhalten, muss man lediglich auf ein fertiges Balkendiagramm die Funktion **coord_flip()** anwenden. Im aktuellen Beispiel ist es allerdings ratsam, zusätzlich für die Textpositionen eine andere Anpassung zu wählen (**hjust** statt **vjust**):

```
> ggplot(kfa, aes(x=geschlecht)) + geom_bar(fill="lemonchiffon", colour="gray50")  
+ labs(x="Geschlecht", y="Häufigkeit")  
+ geom_text(stat="count", mapping=aes(label=..count..), hjust=2, size=4, colour="gray50")  
+ coord_flip()
```

Das Ergebnis:

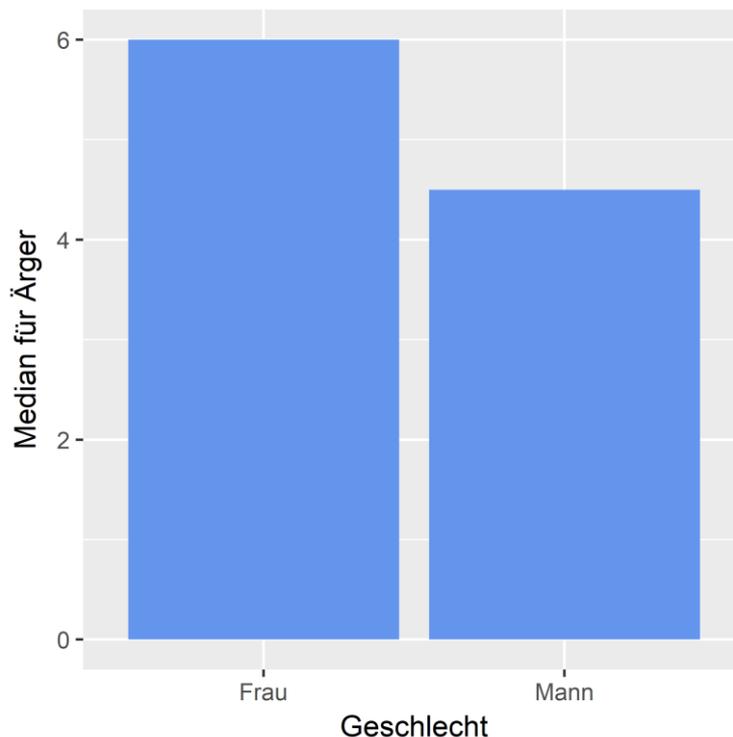


8.3.4.3.2 Statistische Kennwerte einer Variablen für mehrere Gruppen vergleichen

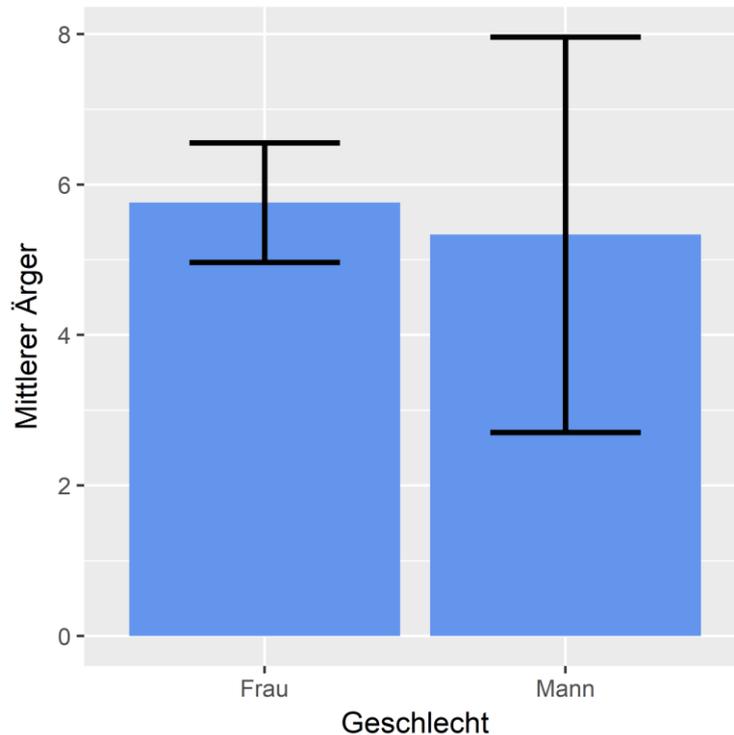
Kombiniert man `geom_bar()` mit der statistischen Transformation `stat_summary()` (an Stelle der Voreinstellung `stat_bin()`), kann man statistische Kennwerte einer abhängigen Variablen (z.B. Mittelwert, Median, Maximum) für die Ausprägungen einer diskreten Variablen darstellen. Das ästhetische Attribut `y` lässt man von einer individuellen Zusammenfassungsfunktion liefern, die über das Argument `fun.y` bestimmt wird, z.B.:

```
> ggplot(kfa, aes(x=geschlecht, y=aergo))  
+   geom_bar(stat="summary", fun.y=median, fill="cornflowerblue")  
+   labs(x="Geschlecht", y="Median für Ärger")
```

Im Beispiel zeigt das resultierende Balkendiagramm die Mediane der Variablen `aergo` für Frauen und Männer:



Im nächsten Diagramm ersetzen wir den Median durch das arithmetische Mittel und ergänzen eine von der Funktion `geom_errorbar()` erstellte Schicht mit den normalverteilungsbasierten 95% - Vertrauensintervallen der beiden Gruppen:



Für die Funktion `geom_errorbar()` wird das Paket **Hmisc** benötigt, das also nötigenfalls installiert und geladen werden muss:

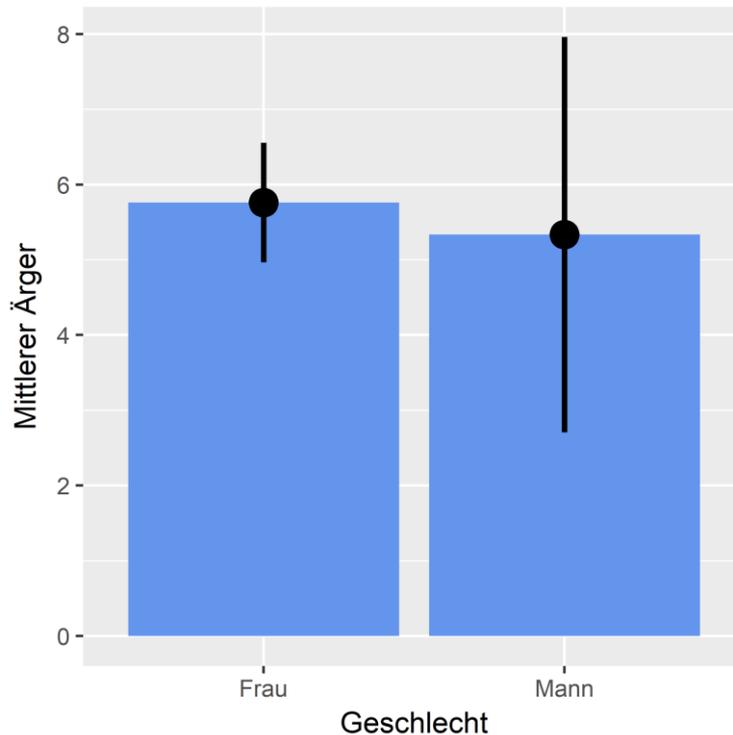
```
> install.packages("Hmisc")
> library(Hmisc)
```

Auch bei `geom_errorbar()` kommt die **summary**-Transformation zum Einsatz:

```
> ggplot(kfa, aes(x=geschlecht, y=aergo))
+   geom_bar(stat="summary", fun.y=mean, fill="cornflowerblue")
+   geom_errorbar(stat="summary", fun.data=mean_cl_normal, width=0.5, size=1)
+   labs(x="Geschlecht", y="Mittlerer Ärger")
```

Dem Argument **fun.data** wird die Zusammenfassungsfunktion `mean_cl_normal()`, zugewiesen, deren Vektor-wertige Rückgabe die Vertrauensschranken enthält, welche die ästhetischen Attribute **ymin** und **ymax** des **errorbar**-Geoms mit Werten versorgen. Über die **errorbar**-Attribute **width** und **size** beeinflusst man die Breite und die Linienstärke der Dispersionsindikatoren.

Um Dispersionsindikatoren *ohne* horizontale Begrenzungen zu erzeugen, verwendet man statt `geom_errorbar` die Alternative `geom_pointrange`:



8.3.4.3.3 Statistische Kennwerte von mehreren Variablen vergleichen

Während ein Balkendiagramm zur Darstellung der Effekte eines *Gruppierungsfaktors* in **ggplot2** leicht zu erstellen ist (siehe Abschnitt 8.3.4.3.2), muss man bei einem *Messwiederholungsfaktor* etwas mehr Aufwand investieren. Die Datentabelle `kfa` enthält in den Variablen `aergo` bzw. `aergm` die Ärgereinschätzungen der Probanden zu zwei Varianten eines imaginierten Schadensfalls:

- Ein Flug wurde deutlich verpasst.
- Ein Flug wurde knapp verpasst.

In der ersten Situationsvariante ist eine kontrafaktische (also positive) Alternative zum Schadensfall weit entfernt, in der zweiten Variante ist sie hingegen nah und steigert den Ärger.

Weil **ggplot2** für ein Balkendiagramm mit den beiden Variablen einen Faktor für die X-Achse benötigt, konvertieren wir die Daten mit Hilfe der Funktion `melt()` aus dem Paket **reshape2** in das Langformat und erstellen einen Faktor namens `kfa` mit den Stufen `aergo` und `aergm`:

```
> library(reshape2)
> kfa.long <- melt(kfa, id.vars = "fnr",
+   measure.vars = c("aergo", "aergm"), variable.name = "kfa")
> levels(kfa.long$kfa) <- c("Ohne", "Mit")
```

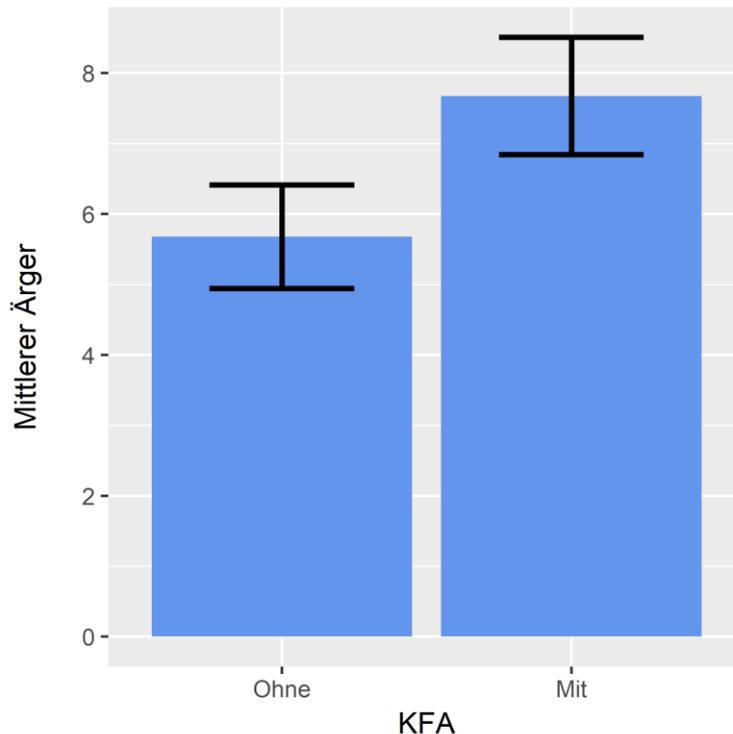
Wir erhalten die folgende Datentabelle `kfa.long` mit dem Faktor `kfa` und dem numerischen Vektor `value`:

```
> kfa.long
  fnr kfa value
  1   1 Ohne    5
  2   2 Ohne    5
  .   .   .
  61  30 Mit   10
  62  31 Mit    9
```

Um dazu ein Balkendiagramm mit Fehlerbalken zu erstellen, kann die **ggplot2**-Syntax aus Abschnitt 8.3.4.3.2 weitgehend übernommen werden:

```
> ggplot(kfa.long, aes(x=kfa, y=value))
+   geom_bar(stat="summary", fun.y=mean, fill="cornflowerblue")
+   geom_errorbar(stat="summary", fun.data=mean_cl_normal, width=0.5, size=1)
+   labs(x="KFA", y="Mittlerer Ärger")
```

Der nicht unerhebliche Aufwand wird durch eine gute Darstellungsqualität belohnt:



Für die Funktion **geom_errorbar()** wird das Paket **Hmisc** benötigt (siehe Abschnitt 8.3.4.3.2).

Mit SPSS ist ein weitgehend identisches Diagramm mit wenigen Mausklicks zu erstellen. Allerdings bietet die **ggplot2** - Lösung mehr Flexibilität, die z.B. zur Erstellung von korrigierten Fehlerbalken genutzt werden könnte: Bei der beschriebenen Fehlerbalkenerstellung bleibt unberücksichtigt, dass ein Messwiederholungsfaktor dargestellt wird. Folglich harmonisiert der wesentlich auf den Fehlerbalken basierende optische Eindruck eventuell schlecht mit der inferenzstatistischen Beurteilung durch den t-Test für abhängige Stichproben, der die interindividuellen Unterschiede aus den Fehlervarianzen eliminiert. Field (2012, S. 361ff) schlägt daher vor, aus den Messwerten die Personeneffekte zu entfernen und die Fehlerbalken aus diesen adjustierten Messwerten zu erstellen.

8.3.4.3.4 Gruppierete Balken

Bislang wurden einfache Balkendiagramme vorgestellt, die entweder eine univariate Verteilung oder den Effekt *eines* (gruppierten oder messwiederholten) Faktors auf eine abhängige Variable darstellen. Nun erstellen wir ein *gruppiertes* Balkendiagramm, das die kombinierten Effekte von zwei Faktoren auf ein metrisches Kriterium zeigt. Wie in Abschnitt 8.3.4.3.3 zu sehen war, werden Messwiederholungsfaktoren durch den Wechsel zum Langformat auf den Gruppierungsfall zurückgespielt. Im aktuellen Abschnitt kommt eine Langvariante der Datentabelle *kfa* zum Einsatz, die auch den Faktor *geschlecht* aus der Ausgangstabelle übernimmt (siehe **melt()** - Argument **id.vars**):

```
> library(reshape2)
> kfa.long <- melt(kfa, id.vars = c("fnr", "geschlecht"),
+   measure.vars = c("aergo", "aergm"), variable.name = "kfa")
> kfa.long$kfa <- factor(kfa.long$kfa, labels=c("Ohne", "Mit"))
```

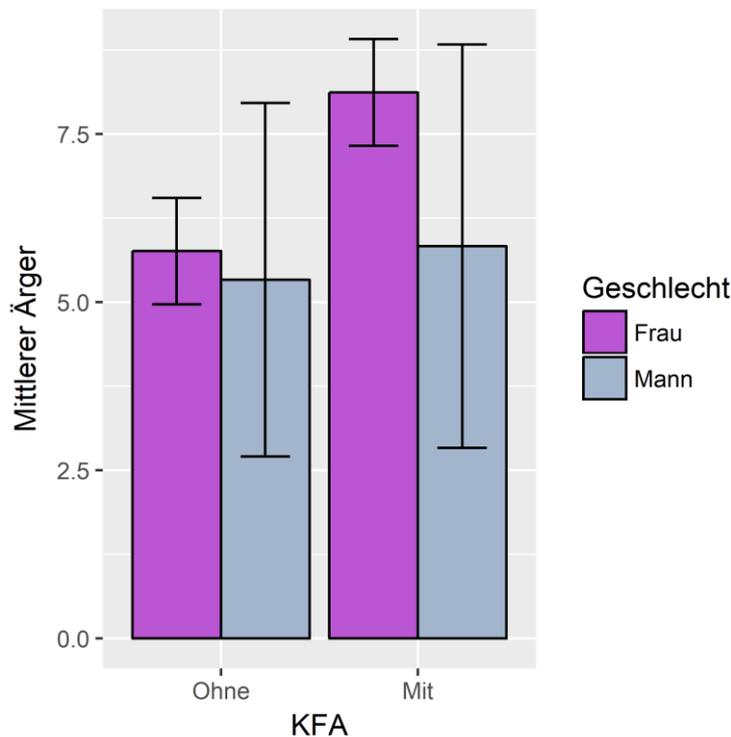
Wir erhalten die folgende Datentabelle `kfa.long` mit den Faktoren `geschlecht` und `kfa` sowie dem numerischen Vektor `value`:

```
> kfa.long
  fnr geschlecht kfa value
1    1         Frau Ohne    5
2    2         Frau Ohne    5
. . .
61  30         Frau Mit   10
62  31         Frau Mit    9
```

Damit anschließend die Funktion `geom_errorbar()` zur Erstellung von Fehlerindikatoren verfügbar ist, muss das Paket **Hmisc** nötigenfalls installiert und geladen werden:

```
> install.packages("Hmisc")
> library(Hmisc)
```

Um das folgende Ergebnis zu erzielen,



wird der Faktor `geschlecht` auf das ästhetische Attribut `fill` abgebildet:

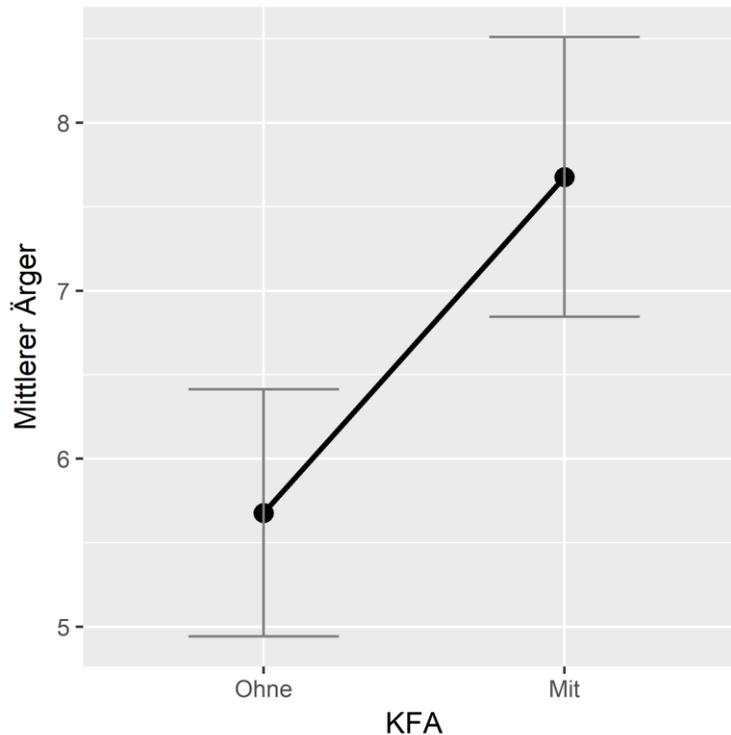
```
> plot <- ggplot(kfa.long, aes(kfa, value, fill=geschlecht))
> plot <- plot + geom_bar(stat="summary", fun.y=mean,
+   position=position_dodge(), colour="black")
> plot <- plot + geom_errorbar(stat="summary", fun.data=mean_cl_normal,
+   position=position_dodge(width=0.9), width=0.5)
> plot <- plot + labs(x="KFA", y="Mittlerer Ärger", fill="Geschlecht")
> plot + scale_fill_manual(values=c("mediumorchid", "lightsteelblue3"))
```

An den beiden X-Achsenpositionen (definiert durch den Faktor `kfa`) sind jeweils *zwei* Balken auszugeben (für Frauen und Männer). Ohne Maßnahme zur Positionsanpassung (vgl. Abschnitt 8.3.1.5) würden

die Balken übereinander gestapelt. Mit der Positionsanpassungsmethode `position_dodge()` wird stattdessen die Ausgabe von gruppierten Balken erreicht. Auf der Schicht mit den Fehlerbalken benötigt die Funktion `position_dodge()` das Argument `width`, weil Balken und Fehlerbalken unterschiedlich breit sind. Mit der Funktion `scale_fill_manual()` wird für eine individuelle Farbauswahl gesorgt.

8.3.4.4 Liniendiagramme

Wir erstellen eine Linienvariante des in Abschnitt 8.3.4.3.3 vorgestellten Balkendiagramms zum Effekt des Messwiederholungsfaktors KFA auf den Ärger der Probanden:



Im Schichtaufbau werden zunächst mit `geom_point()` zwei Punkte zur Darstellung der beiden Mittelwerte ausgegeben, wobei die statistische Transformation `summary` mit der Zusammenfassungsfunktion `fun.y` zum Einsatz kommt:

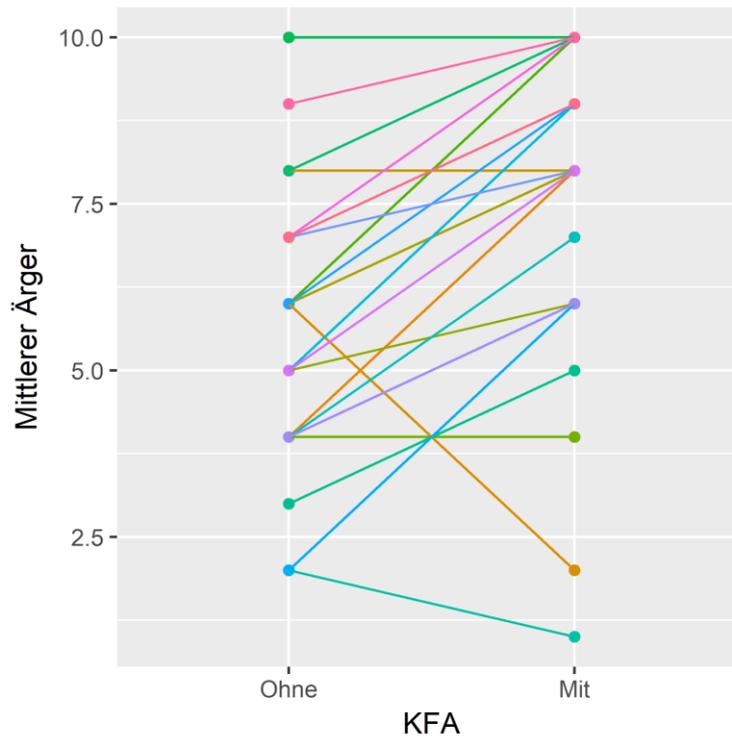
```
> plot <- ggplot(kfa.long, aes(x=kfa, y=value)) + geom_point(stat="summary",
+   fun.y=mean, size=3)
> plot <- plot + geom_line(stat="summary", fun.y=mean, aes(group=1), size=1)
> plot <- plot + geom_errorbar(stat="summary", fun.data=mean_cl_normal,
+   width=0.5, colour="gray50")
plot + labs(x="KFA", y="Mittlerer Ärger")
```

Dieselbe Zusammenfassungsfunktion wird auch im Aufruf von `geom_line()` verwendet, wo ansonsten die Abbildung des ästhetischen Attributs `group` auf den Wert 1 zu beachten ist. Damit wird die Gruppierung aufgrund des X-Achsen-Faktors `kfa` aufgehoben, sodass die beiden aus der Zusammenfassung hervorgehenden Punkte zu *einer* Gruppe gehören. Diese Maßnahme ist erforderlich, weil `geom_line()` zu jeder Gruppe eine Linie erstellt, die alle enthaltenen Punkte verbindet. Für die Funktion `geom_errorbar()` wird das Paket **Hmisc** benötigt (siehe Abschnitt 8.3.4.3.2).

Um die Veränderung zwischen den beiden Bedingungen (ohne bzw. mit KFA) für jede *einzelne* Person zu betrachten, legen wir die Variable `fnr` mit der Fallnummer als farbdefinierend und gruppenbildend fest, nachdem sie in einen Faktor gewandelt worden ist:

```
> plot <- ggplot(kfa.long, aes(kfa, value, colour=factor(fnr)))  
> plot <- plot + geom_point(show.legend=FALSE)  
> plot <- plot + geom_line(aes(group=fnr), show.legend=FALSE)  
> plot + labs(x="KFA", y="Mittlerer Ärger")
```

Es resultiert ein „Mikadodiagramm“, wobei die Legenden zum **point**- bzw. **line**-Geom durch den Wert **FALSE** für das Argument **show.legend** unterdrückt werden:



Vorschläge zu weiteren Liniendiagrammen finden sich z.B. bei Chang (2013, S. 49ff) und Field (2012, S. 155ff).

9 Statistische Datenanalyse mit R

In diesem Abschnitt wird keinesfalls versucht, den enormen **R**-Funktionsumfang zur statistischen Datenanalyse zu beschreiben. Das (leider bei weitem noch nicht erreichte) Ziel besteht in einer Sammlung nützlicher Werkzeuge für Personen, die ihre statistischen Analysen überwiegend mit SPSS erledigen wollen. In Abschnitt 9.1 werden einfache Funktionen zur univariaten Verteilungsbeschreibung vorgestellt, die sich potentiell zur Verwendung in **R**-Skripten eignen. Abschnitt 9.2 behandelt die **R**-Syntax zur Modellformulierung, die in zahlreichen Auswertungsfunktionen Verwendung findet. Im Abschnitt 9.3 werden exemplarisch einige **R**-Funktionen beschrieben, die Lücken im Statistikangebot von SPSS schließen. Sie sind über das Formulieren von **R**-Anweisungen zu nutzen, während man bei den in Abschnitt 3 vorgestellten SPSS-Erweiterungskommandos von **R** profitiert, ohne seine Syntax anwenden zu müssen.

9.1 Einfache univariate Verteilungsbeschreibung

In diesem Abschnitt werden einfache **R**-Funktionen zur Beschreibung von univariaten Verteilungen vorgestellt, die sich potentiell zur Verwendung in eigenen **R**-Skripten eignen. Es geht *nicht* darum, die in **R** zahlreich vorhandenen Schätz- und Testmethoden für univariate Verteilungen vorzustellen.

9.1.1 Univariate Verteilungsbeschreibung für metrische Variablen

9.1.1.1 Kompakte Verteilungsbeschreibung

Von der Funktion `summary()` erhält man für die Stichprobe in einem Vektor das Minimum, das Maximum, das 1., 2. und 3. Quartil sowie den Mittelwert. Im folgenden Beispiel werden die se Statistiken für eine per `rnorm()` (siehe Abschnitt 7.4.1.1) generierte Zufallsstichprobe aus einer normalverteilten Population ermittelt:

```
> summary(rnorm(100, 3, 2))
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-2.649  2.126   3.357   3.336   4.515   8.401
```

9.1.1.2 Statistische Funktionen für numerische Vektoren

Statistische Funktionen für numerische Vektoren:

- **min()**, **max()**
Die Funktionen `min()` bzw. `max()` liefert das kleinste bzw. größte Element, z.B.:

```
> min(c(1,2,3,5))
[1] 1
> max(c(1,2,3,5))
[1] 5
```
- **sum()**
Liefert die Summe der Elemente, z.B.:

```
> sum(c(1,2,3,5))
[1] 11
```
- **mean()**
Liefert das arithmetische Mittel der Elemente, z.B.:

```
> mean(c(1,2,3,5))
[1] 2.75
```

- **median()**

Liefert den Median der Elemente, z.B.:

```
> median(c(1,2,3,5))
[1] 2.5
```

- **var(), sd()**

Diese Funktionen liefern die Varianz bzw. Standardabweichung einer Stichprobe, z.B.:

```
> var(c(1,2,3,5,4,5,2,3,1,4,5,7,8,2,3,4,2,3,8))
[1] 4.508772
> sd(c(1,2,3,5,4,5,2,3,1,4,5,7,8,2,3,4,2,3,8))
[1] 2.123387
```

- **quantile()**

Liefert per Voreinstellung die Quartile, z.B.:

```
> quantile(c(1,2,3,5,4,5,2,3,1,4,5,7,8,2,3,4,2,3,8))
 0%  25%  50%  75% 100%
  1   2   3   5   8
```

Über einen Parametervektor mit Wahrscheinlichkeiten sind auch andere Quantile verfügbar, z.B.:

```
> quantile(c(1,2,3,5,4,5,2,3,1,4,5,7,8,2,3,4,2,3,8), c(1/3, 2/3))
33.33333% 66.66667%
      3      4
```

Als Rückgabe erhält man einen Vektor mit benannten Elementen (vgl. Abschnitt 5.3.4.2.6).

- **IQR()**

Liefert den Interquartilsabstand (Abstand vom 1. bis zum 3. Quartil), z.B.:

```
> IQR(c(1,2,3,5,4,5,2,3,1,4,5,7,8,2,3,4,2,3,8))
[1] 3
```

Enthält ein Vektor mindestens einen fehlenden Wert (**NA** oder **NaN**), dann liefern die genannten Funktionen das Ergebnis **NA** oder **NaN**, z.B.:

```
> a <- c(1, NA, 3)
> mean(a)
[1] NA
```

Soll stattdessen aus den vorhandenen Argumenten ein Ergebnis ermittelt werden, ist das Argument **na.rm** auf den Wert **TRUE** zu setzen, z.B.:

```
> mean(a, na.rm=TRUE)
[1] 2
```

9.1.2 Absolute und relative Häufigkeiten für kategoriale Variablen ausgeben

Die Funktion **table()** liefert für einen Vektor oder Faktor die absoluten Häufigkeiten der auftretenden Werte, z.B.:

```
> daten <- c(1,1,1,2,2,2,2,3,3,3,3,4,4,5,5,5,5)
> (tabelle <- table(daten))
daten
 1 2 3 4 5
 3 4 4 2 4
```

Zu den relativen Häufigkeiten verhilft die Funktion **prop.table()**, die eine Tabelle mit absoluten Häufigkeiten als Argument benötigt, z.B.:

```
> prop.table(tabelle)
daten
      1      2      3      4      5
0.1764706 0.2352941 0.2352941 0.1176471 0.2352941
```

9.2 Modellformulierung

Ein univariates Modell (z.B. für die Funktionen **lm()**, **glm()** oder **rlm()**) verwendet folgende Syntax:

$$\text{Kriterium} \sim \text{Design}$$

Das zu modellierende Kriterium kann sein:

- eine Variable
- ein Ausdruck, z.B.
 $\log(y) \sim x$

Das Design enthält mindestens einen Term. Treten mehrere Terme auf, sind diese durch den Operator + zu verbinden, z.B.:

$$y \sim x1 + x2$$

Erlaubte Terme:

- Variablen (z.B. Vektoren oder Faktoren aus einer Datentabelle)
- Ausdrücke
- Die Zahl **1** steht für den konstanten Term eines linearen Modells und muss nicht angegeben werden, weil der konstante Term implizit enthalten ist.
- Wechselwirkungsterme bestehend aus mindestens zwei durch den Operator **:** verbundene Variablen, z.B.:

$$y \sim a + b + a:b$$

Formuliert man einen Wechselwirkungsterm mit dem Operator *****, dann sind alle Terme niedrigerer Ordnung einbezogen, z.B.:

Kurzschreibweise	enthaltene Terme
a*b	a, b, a:b
a*b*c	a, b, c, a:b, a:c, b:c, a:b:c

Über den Minus-Operator kann ein Term aus dem Design entfernt werden, was gelegentlich mit dem implizit enthaltenen konstanten Term eines Modells (symbolisiert durch die Zahl **1**) geschieht. Im folgenden Beispiel mit den numerischen Vektoren **x** und **y** resultiert das Modell einer bivariaten linearen Regression mit einer durch den Ursprung des Koordinatensystems verlaufenden Regressionsgeraden:

$$y \sim x - 1$$

Im Design kann an Stelle einer Variablen auch ein Ausdruck stehen, z.B.:

$$y \sim \log(x)$$

Werden im Ausdruck Operatoren benötigt, die in der Modellformulierung eine alternative Bedeutung haben, muss per **I()** - Funktion die arithmetische Bedeutung der Operatoren wiederhergestellt werden, z.B.:

$$\text{lm}(y \sim \text{I}(x1*x2))$$

Weitere Details zur Modellformulierung finden sich z.B. bei Venables et al. (2014, Abschnitt 11.1) und Wollschläger (2010, S. 163ff).

9.3 Lücken im SPSS-Statistikangebot füllen

In diesem Abschnitt werden weitere statistische Methoden behandelt, die in SPSS fehlen und die mit **R** leicht zu realisieren sind. Viele Ergänzungen der SPSS-Funktionalität konnten schon in Abschnitt 3 vorgestellt werden, weil sie über Erweiterungsbundles ohne **R**-Kenntnisse zu nutzen sind. In Abschnitt 4 wurde erläutert, wie man **R**-Programme im SPSS-Syntaxfenster erstellt und darin auf SPSS-Variablen zugreift. In Abschnitt 5 wurde **R** als Programmierumgebung vorgestellt, so dass wir nun die Angebote im **R**-Universum nahezu uneingeschränkt nutzen können.

Es fällt allerdings schwer, Lücken von allgemeinem Interesse im SPSS-Statistikangebot zu finden, die noch nicht durch eine bequem via SPSS-Menü verfügbare Erweiterung geschlossen worden sind. Solche Erweiterungen sind oft **R**-basiert, aber ohne direkten Kontakt mit **R** zu verwenden. Anschließend werden zwei aktuell nur via **R**-Syntax verfügbare Analysemethoden vorgestellt:

- Gerichtete t-Tests und einseitige Vertrauensintervalle
- Soziale Netzwerkanalyse mit dem R-Paket **igraph**

9.3.1 Gerichtete t-Tests und einseitige Vertrauensintervalle

Bei t-Tests liefert SPSS grundsätzlich nur den p-Wert zum *zweiseitigen* Test und dazu passend das zweiseitige Konfidenzintervall. Wenn ein *gerichtetes* Testproblem vorliegt, kann der benötigte einseitige p-Wert leicht berechnet werden (durch Halbieren des zweiseitigen p-Werts). Das zugehörige einseitige Vertrauensintervall zu bestimmen, ist etwas umständlich. Diese Mühe kann man sich ersparen durch Verwendung der **R**-Funktion **t.test()**, die über das Argument **alternative** mit den Werten

- **"two.sided"** (Voreinstellung),
- **"greater"** oder
- **"less"**

eine Testausrichtung entgegennimmt und bei einseitiger Testung auch ein passendes einseitiges Vertrauensintervall berechnet. Wir rechnen als Beispiel einen Einstichproben - t-Test zum Hypothesenpaar

$$H_0: \mu \leq 0 \text{ versus } H_1: \mu > 0$$

und verwenden dazu eine Zufallsstichprobe ($N = 50$) aus einer normalverteilten Population mit dem Mittelwert 0,2 und der Varianz 1:

```
> sam <- rnorm(50, 0.2, 1)
> t.test(sam)
> t.test(sam, alternative="greater")
```

Beim ungerichteten t-Test als Folge des ersten **t.test()** - Aufrufs resultiert ein p-Level von 0,063, und die Nullhypothese muss beibehalten werden. Das zweiseitige Konfidenzintervall (zum voreingestellten Niveau von 0,95) enthält dementsprechend den Wert Null:

One Sample t-test

```
data: sam
t = 1.9022, df = 49, p-value = 0.06303
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.01514014  0.55151178
sample estimates:
mean of x
0.2681858
```

Beim gerichteten t-Test als Folge des zweiten `t.test()` - Aufrufs resultiert ein p-Level kleiner als 0,05, so dass die einseitige Nullhypothese verworfen werden kann. Das einseitige Konfidenzintervall liegt dementsprechend komplett rechts von der 0:

One Sample t-test

```
data: sam
t = 1.9022, df = 49, p-value = 0.03152
alternative hypothesis: true mean is greater than 0
95 percent confidence interval:
 0.03181230      Inf
sample estimates:
mean of x
0.2681858
```

9.3.2 Soziale Netzwerkanalyse mit igraph

Die statistische Datenanalyse geht in der Regel von einer Tabelle (Matrix) aus mit ...

- Beobachtungseinheiten (Fällen) in den Zeilen
- und Merkmalen (Eigenschaften) der Beobachtungseinheiten in den Spalten.

In den letzten Jahren finden (soziale) Netzwerke als Forschungsgegenstand zunehmend Interesse, z.B.

- Soziale Kontakte in einer Kindergartengruppe
- Handelsbeziehungen zwischen Staaten
- Wechselseitiges Zitieren in einer Gruppe von Wissenschaftlern

Die Akteure im Netzwerk bezeichnet man oft als *Knoten* (engl.: *nodes*) oder *Ecken* (engl. *vertices*). Ihre Beziehungen werden durch Linien zwischen den Knoten dargestellt, die oft als *Kanten* (engl.: *edges*) bezeichnet werden. In diesem Abschnitt kann natürlich keine gründliche Beschreibung der Begriffe und Techniken der sozialen Netzwerkanalyse angestrebt werden, die eine umfängliche Monographie erfordert (siehe z.B. Yang et al. 2016).

Eine typische Datenbasis zur Beschreibung und Analyse von sozialen Netzwerken enthält:

- Eine Tabelle zur Beschreibung der Akteure (engl. Bezeichnung: *nodelist*)
Diese Tabelle besitzt das eingangs beschriebene Standardformat mit den Akteuren als Zeilen. Oft dient die einzige Spalte zur Identifikation der Akteure. Es können aber in weiteren Spalten auch zusätzliche Eigenschaften der Akteure vorhanden sein.
- Eine Beschreibung der Beziehungen zwischen den Akteuren
Dabei kommen folgende Formate in Frage:
 - Eine quadratische Matrix mit (gerichteten) Beziehungsstärken in den Zellen (engl. Bezeichnung: *adjacency matrix*)
 - Eine Tabelle mit den Beziehungen als Zeilen (engl. Bezeichnung: *edgelist*)
Oft dienen die erste Spalte zur Identifikation des Senders und die zweite Spalte zur Identifikation des Empfängers. In weiteren Spalten können optional zusätzliche Eigenschaften der Beziehungen vorhanden sein (z.B. Stärke, Typ).

Anschließend wird versucht, in Anlehnung an Griffin (2017) einen ersten Eindruck von der Darstellung und Analyse sozialer Netzwerke mit dem von Gábor Csárdi erstellten **R**-Paket **igraph** (Version 1.0.1) zu vermitteln. Dabei bleiben zahlreiche, von **igraph** unterstützte Netzwerkoptionen unberücksichtigt, z.B.

- gerichtete Verbindungen
- Netzwerke mit verschiedenen Knotentypen (z.B. Informationsquellen und Informationskonsumenten, siehe Beispiel in Ognyanova 2016)

Eine ausführliche Beschreibung des sehr leistungsfähigen **igraph**-Pakets bietet z.B. das Buch von Kolaczyk & Csárdi (2014).

9.3.2.1 Visualisierung

Wir laden die von Griffin (2017) in der Datei **Pan.csv** bereit gestellten Daten zum Interaktionsverhalten von 23 Schimpansen herunter:

https://raw.githubusercontent.com/rgriff23/Primate_social_networks/master/data/Pan.csv

Das Paket **igraph** wird nötigenfalls installiert und geladen:

```
> install.packages("igraph")
> library("igraph")
```

Wir lesen die Beziehungsmatrix in **Pan.csv** mit der Funktion **read.csv()** und sorgen durch einen umrahmenden **data.frame()** - Aufruf mit **row.names** - Argument dafür, dass die Einträge der ersten Spalte als Zeilennamen interpretiert werden:

```
> relations <- data.frame(read.csv("Pan.csv"), row.names=1)
```

Mit Hilfe der **igraph**-Funktion **graph_from_adjacency_matrix()** erstellen wir aus der Beziehungsmatrix ein **igraph**-Objekt mit ungerichteten, gewichteten Kanten:

```
> igr.obj <- graph_from_adjacency_matrix(as.matrix(relations), mode="undirected",
+                                       weighted=TRUE)
```

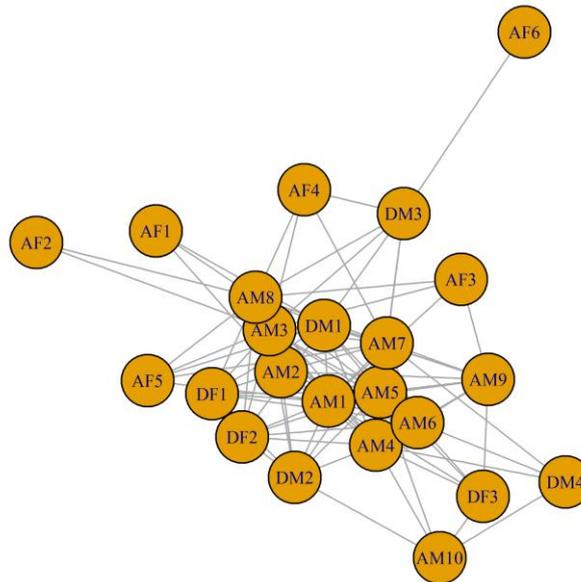
Dabei werden die Knoten mit der Identifikation als der einzigen Eigenschaft automatisch aus der Beziehungsmatrix extrahiert.

Mit dem folgenden **plot()** - Aufruf

```
> plot(igr.obj, main="Interaktionen zwischen Schimpansen (Griffin 2017)",
+       vertex.size=20, vertex.label.cex=0.7)
```

ist eine erste Darstellung des Netzwerks unter Verwendung des voreingestellten Netzwerklayouts schnell erstellt:

Interaktionen zwischen Schimpansen (Griffin 2017)



Zur Gestaltung des Diagramms werden im **plot()** - Aufruf folgende Argumente verwendet:

- **main**
Legt den Titel fest
- **vertex.size**
Beeinflusst die Größe der Knoten
- **vertex.label.cex**
Durch einen positiven Skalierungsfaktor kann die Größe der Knotenbeschriftungen im Vergleich zum Standard (Wert = 1) verändert werden.

In den Knoten befinden sich die Fallbeschriftungen, wobei im Beispiel der zweite Buchstabe das Geschlecht des Tieres angibt.

Um die verhaltensbiologische Interpretation des Netzwerks zu erleichtern, werden anschließend nach einem Vorschlag von Griffin (2017) folgende Verbesserungen vorgenommen:

- Das Geschlecht der Schimpansen wird farblich dargestellt (pink für weibliche und hellblau für männliche Tiere).
- Die Interaktionsintensität zwischen zwei Tieren wird durch die Breite der verbindenden Kante dargestellt.

Zunächst erstellen wir einen numerischen Vektor, der für weibliche Tiere eine 1 und für männliche Tiere eine 2 enthält:

```
> vert.col <- (substr(vertex_attr(igr.obj)$name,2,2)=="F")+1
```

Durch die **igraph**-Funktion **vertex_attr()** erhalten wir eine Liste mit den Knotenattributen. Darin befindet sich der Zeichenfolgenvektor **name** mit den Fallbeschriftungen. Mit der Funktion **substr()** aus dem **R**-Basispaket wird aus den Fallbeschriftungen das 2. Zeichen extrahiert, um es per Identitätsoperator mit „F“ zu vergleichen.

Im folgenden **plot()** - Aufruf¹

¹ Eine Dokumentation der im **plot()** - Aufruf unterstützten Parameter ist hier zu finden:

<http://igraph.org/r/doc/plot.common.html>

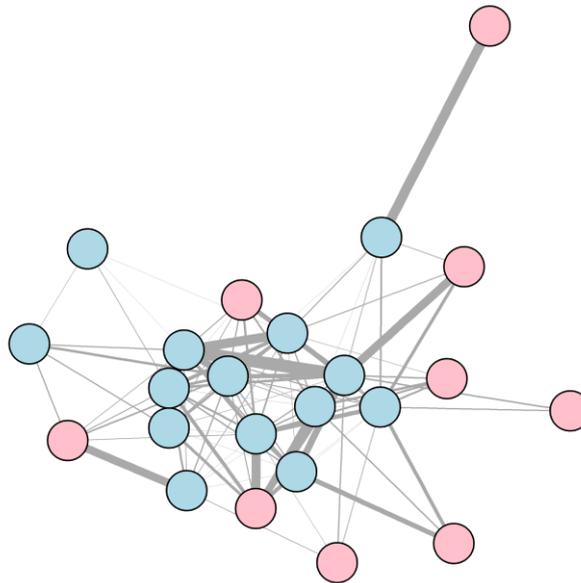
```
> plot(igr.obj,  
+      main="Interaktionen zwischen Schimpansen (Griffin 2017)",  
+      vertex.color=c("lightblue","pink")[vert.col],  
+      edge.width=edge_attr(igr.obj)$weight*0.001,  
+      vertex.label=NA)
```

erhält das zur Gestaltung der Knotenfarben dienende Argument **vertex.color** als Wert einen Vektor mit geschlechtsabhängigen Farbnamen (**lightblue** bei männlichen und **pink** bei weiblichen Tieren). Außerdem wird dem zur Gestaltung der Kantenbreiten dienenden Argument **edge.width** der numerische Vektor **weight** zugewiesen. Er stammt aus der Liste mit den Kantenattributen, welche die **igraph**-Funktion **edge_attr()** aus dem **igraph**-Objekt extrahiert. Schließlich werden der Übersichtlichkeit halber die Knotenbeschriftungen mit dem Wert **NA** für das Argument **vertex.label** abgeschaltet.

Die Anordnung der Knoten ist durch die Beziehungsmatrix keinesfalls determiniert und kann über den **plot()** - Parameter **layout** beeinflusst werden. Sein voreingestellter Wert **layout_nicely** sorgt u.a. dafür, dass die Knoten gleichmäßig verteilt werden und möglichst wenige Überschneidungen der Kanten auftreten. Weil der Platzierungs-Algorithmus mit einer zufälligen Ausgangslage startet, erhält man für ein Netzwerk bei jedem **plot()** - Aufruf ein anderes Ergebnis.

In der patriarchalisch organisierten Schimpansengruppe wird das Zentrum des Netzwerks dominiert von intensiv miteinander kommunizierenden männlichen Tieren. Die weiblichen Tiere befinden sich am Rand des Netzwerks und interagieren bevorzugt mit einzelnen männlichen Tieren:

Interaktionen zwischen Schimpansen (Griffin 2017)



9.3.2.2 Analyse

Zur Analyse von sozialen Netzwerken wurden Indizes für einzelne Knoten und gesamte Netzwerke definiert.

Einige Zentralitätsindizes auf Individualebene sind:

- **Grad der Zentralität**

Der Zentralitätsgrad eines Knotens ist definiert durch die Zahl seiner Verbindungen mit anderen Knoten. In **igraph** erhält man die Zentralitätsgrade der Knoten über die Funktion **degree()**. Für die in Griffin (2017) beschriebene Schimpansengruppe resultieren folgende Werte:

```
> (degree(igr.obj))
  AM1  AM2  AM3  AM4  AM5  AM6  AM7  AM8  AM9  AM10  DM1  DM2  DM3  DM4  AF1
    15   11   16   14   15   12   14   15    9    5   13    9    7    4    3
  AF2  AF3  AF4  AF5  AF6  DF1  DF2  DF3
    2    4    4    5    1    8    8    6
```

- **Gewichteter Grad bzw. Stärke der Zentralität**

Die Zentralitätsstärke eines Knotens ist definiert durch die Summe seiner Verbindungsgewichte. In **igraph** erhält man die Zentralitätsstärken der Knoten über die Funktion **strength()**. Für die Schimpansengruppe resultieren folgende Werte:

```
> (strength(igr.obj))
  AM1  AM2  AM3  AM4  AM5  AM6  AM7  AM8  AM9  AM10  DM1  DM2
22107 16182 36067 28120 19576 15237 17102 15334  9251  4823 10381  7895
  DM3  DM4  AF1  AF2  AF3  AF4  AF5  AF6  DF1  DF2  DF3
 9658  1478  6497  1909  2188  8698  5800  6190 21085  9213  9101
```

- **Verbindungs-zentralität**

Die Verbindungs-zentralität (engl.: *betweenness*) eines Knotens K ist definiert durch die Zahl der Knotenpaare, deren *kürzeste* Verbindung durch K verläuft. In **igraph** erhält man die Verbindungs-zentralitäten der Knoten über die Funktion **betweenness()**. Für die Schimpansengruppe resultieren folgende Werte:

```
> (betweenness(igr.obj))
  AM1  AM2  AM3  AM4  AM5  AM6  AM7  AM8  AM9  AM10  DM1  DM2  DM3  DM4  AF1
    37    4    1   19    0   17   34   65   57    0  138    2   54   21    0
  AF2  AF3  AF4  AF5  AF6  DF1  DF2  DF3
    0    0    0    0    0    0    0    0
```

Von den Indizes auf Netzwerkebene soll nur die *Dichte* angesprochen werden, die als Anteil der vorhandenen Verbindungen bezogen auf die Gesamtheit der möglichen Verbindungen definiert ist. Bei einem Netzwerk mit 5 Knoten sind nach der Paarbildungsregel maximal

$$\frac{5(5-1)}{2} = 10$$

Kanten möglich.

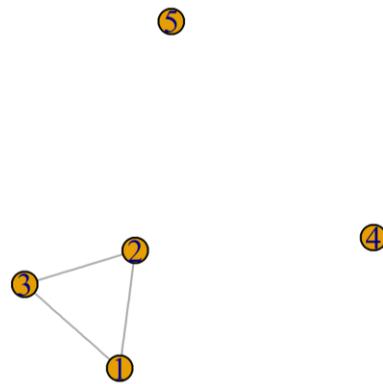
Anschließend werden kleine Netzwerke als Beispiele betrachtet, die wir per **graph()** - Funktion erstellen, statt eine Beziehungsmatrix einzulesen. Im ersten Beispiel definieren wir ein ungerichtetes Netzwerk mit 5 Knoten und 3 Kanten:

```
> n1 <- graph(edges=c(1,2, 2,3, 3,1), n=5, directed=FALSE)
```

Was der **graph()** - Aufruf bewirkt, ist per **plot()** - Aufruf

```
> plot(n1)
```

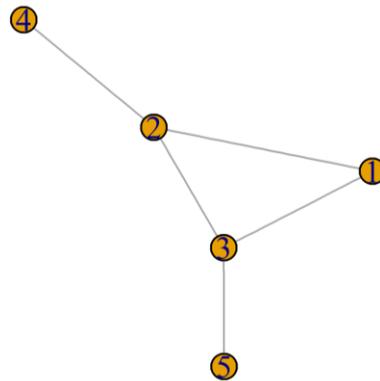
leicht zu visualisieren:



Dieses Netzwerk hat nach obiger Definition den Dichtewert 0,3, weil 3 von 10 möglichen Kanten vorhanden sind. Den bei einem größeren Netzwerk nicht ganz so leicht zu ermittelnden Dichtewert liefert die **igraph()** - Funktion **edge_density()**, z.B.:

```
> edge_density(n1)
[1] 0.3
```

Für das folgende Netzwerk



resultiert eine höhere Dichte:

```
> n2 <- graph(edges=c(1,2, 2,3, 3,1, 4,2, 5,3), n=5, directed=FALSE)
> edge_density(n2)
[1] 0.5
```

10 Mathematische Probleme mit R bewältigen

Mit **R** kann man nicht nur statistische Probleme lösen, sondern kann auch Aufgaben aus anderen Disziplinen der Mathematik bewältigen.

10.1 Mengenlehre

Ein Vektor kann eine Menge repräsentieren. Dann ...

- spielt die Reihenfolge der Elemente keine Rolle,
- werden Dubletten ignoriert.

Um Dubletten explizit zu entfernen, verwendet man die Funktion **unique()**, z.B.:

```
> set1 <- c(1,2,2,3,4,5)
> (uset <- unique(set1))
[1] 1 2 3 4 5
> length(unique(set1))
[1] 5
```

Den Durchschnitt zweier Mengen liefert die Funktion **intersect()**, z.B.:

```
> set2 <- c(3,4,5,6,7)
> intersect(set1, set2)
[1] 3 4 5
```

Um zwei Mengen zu vereinigen, verwendet man die Funktion **union()**, z.B.:

```
> union(set1, set2)
[1] 1 2 3 4 5 6 7
```

Subtrahiert man über die Funktion **setdiff()** von einer Menge **set1** eine andere Menge **set2**, dann resultiert als Differenz eine Menge mit allen **set1**-Elementen, die sich *nicht* in **set2** befinden, z.B.:

```
> setdiff(set1, set2)
[1] 1 2
```

Ob sich ein bestimmtes Element in einer Menge befindet, prüft man mit der Funktion **is.element()**, z.B.:

```
> is.element(2, set1)
[1] TRUE
```

Wird für einen Vektor die elementweise Existenzprüfung vorgenommen, resultiert ein Vektor vom Typ **logical**, z.B.:

```
> is.element(c(0,1,2), set1)
[1] FALSE TRUE TRUE
```

Alternativ zur Funktion **is.element()** kann mit demselben Ergebnis der Operator **%in%** verwendet werden, z.B.:

```
> c(0,1,2) %in% set1
[1] FALSE TRUE TRUE
```

Um für eine Menge zu prüfen, ob sie **Teilmenge** einer anderen Menge ist, wendet man die Funktion **all()** auf die Rückgabe der Funktion **is.element()** bzw. des Operators **%in%** an, z.B.:

```
> all(c(0,1,2) %in% c(0:4))
[1] TRUE
```

Weitere Mengenoperationen mit **R** werden in Wollschläger (2010, S. 40ff) beschrieben.

10.2 Lineare Algebra

Um das folgende Gleichungssystem zu lösen,

$$\begin{bmatrix} 5 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

kann man die **R**-Funktion `solve()` verwenden, die als ersten Parameter die Koeffizientenmatrix und als zweiten Parameter den Zielvektor (oder die Zielmatrix) erwartet.

```
> xmat <- matrix(c(1, 2, 2, 3), 2)
> xmat
      [,1] [,2]
[1,]    1    2
[2,]    2    3
> b <- c(5, 8)
> solve(xmat, b)
[1] 1 2
```

Lässt man den zweiten Parameter weg, wird dort die Einheitsmatrix angenommen, so dass die Inverse der Koeffizientenmatrix als Lösung verlangt wird:

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

Im Beispiel:

```
> ix <- solve(xmat)
> ix
      [,1] [,2]
[1,]   -3    2
[2,]    2   -1
```

Um das Ergebnis zu prüfen, multiplizieren wir die Matrizen `xmat` und `ix`, wobei das Operatorzeichen `%*%` zu verwenden ist:

```
> xmat %*% ix
      [,1] [,2]
[1,]    1    0
[2,]    0    1
```

Literatur

- Baltes-Götz, B. (2016). *Lineare Regressionsanalyse mit SPSS*. Online-Dokument: <http://www.uni-trier.de/index.php?id=22489>
- Baltes-Götz, B. (2017). *Mediator- und Moderatoranalyse mit SPSS und PROCESS*. Online-Dokument: <http://www.uni-trier.de/index.php?id=22528>
- Bellio, R. & Ventura, L. (2005). *An Introduction to Robust Estimation with R Functions*. Online-Dokument (abgerufen: 28.12.2015): <http://www.dst.unive.it/rsr/BelVenTutorial.pdf>
- Bliese, P.D. (2013). *Multilevel Modeling in R 2.5. A Brief Introduction to R, the multilevel package and the nlme package*. Online-Dokument (abgerufen: 28.12.2015): http://cran.r-project.org/doc/contrib/Bliese_Multilevel.pdf
- Chambers, J. (1998). *Programming with Data. A Guide to the S Language*. New York: Springer.
- Chang, W. (2013). *R Graphics Cookbook* (2nd ed.). Beijing: O'Reilly.
- Cohen, R. (2016). *Extension Bundles from IBM SPSS*. Online-Dokument (abgerufen: 01.12.2015): <https://www.ibm.com/developerworks/community/files/form/anonymous/api/library/b5bb8a42-04d2-4503-93bb-dc45d7a145c2/document/1d445b76-d706-44a6-85bf-9e3738d223a4/media/Extension%20Bundles%20from%20IBM%20SPSS.pdf>
- Field, A. (2012). *Discovering Statistics Using R*. London: SAGE Publications.
- Fox, J. (2002). *Robust Regression*. Online-Dokument (abgerufen am 21.01.2011): <http://cran.r-project.org/doc/contrib/Fox-Companion/appendix-robust-regression.pdf>.
- Fox, J. (2005). The R Commander: A Basic Statistics Graphical User Interface to R. *Journal of Statistical Software*, 14(9), 1-42.
- Fox, J. & Weisberg, S. (2011). *An R Companion to Applied Regression* (2nd ed.). Thousand Oaks: SAGE Publications.
- Fox, J. & Bouchet-Valat, M. (2015). *Getting Started With the R Commander*. Online-Dokument (abgerufen: 11.12.2016): <http://socserv.mcmaster.ca/jfox/Misc/Rcmdr/Getting-Started-with-the-Rcmdr.pdf>
- Gordon, R.A. (2010). *Regression Analysis for the Social Sciences*. New York: Routledge.
- Griffin, R. H. (2017). *Primate social networks. Basics of visualization and analysis with igraph*. Online-Dokument (abgerufen: 27.08.2017): <https://rgriff23.github.io/2017/04/26/primate-social-networks-in-igraph.html>
- Hain, J. (2011). *Statistik mit R*. RRZN-Handbuch.
- IBM Corp. (2016a). *Programming and Data Management for SPSS Statistics 24*. Online-Dokument https://www.ibm.com/developerworks/community/groups/service/html/communityview?communityUuid=ab16c38e-2f7b-4912-a47e-85682d124d32#fullpageWidgetId=W7de6c01e7f8d_483f_820f_d9d2108adc61
- IBM Corp. (2016b). *IBM SPSS Statistics 24 Core-System Benutzerhandbuch*. Online-Dokument: ftp://public.dhe.ibm.com/software/analytics/spss/documentation/statistics/24.0/de/client/Manuals/IBM_SPSS_Statistics_Core_System_User_Guide.pdf
- Kolaczyk, E.D. & Csárdi, G. (2014). *Statistical Analysis of Network Data with R*. New York: Springer.
- Ligges, U. (2007). *Programmieren mit R* (2. Aufl.). Berlin: Springer.
- Muenchen, R.A. (2011). *R for SAS and SPSS Users* (2nd ed.). New York: Springer.
- Murrell, P. (2011). *R Graphics* (2nd ed.). Boca Raton, FL: Chapman & Hall/CRC.

- Ognyanova, K. (2016). *Network Analysis and Visualization with R and igraph*. Online-Dokument (abgerufen am 27.08.2017): <http://kateto.net/networks-r-igraph>
- R Development Core Team (2016). *R Language Definition. Version 3.2.5*. Verfügbar über die Hilfe zum RGui.
- Ryan, T.S. (1997). *Modern Regression Methods*. New York: Wiley.
- IDRE UCLA (2014a). *R Data Analysis Examples: Robust Regression*. Online-Dokument (abgerufen am 30.11.2014): <http://www.ats.ucla.edu/stat/r/dae/rreg.htm>
- IDRE UCLA (2014b). *R Data Analysis Examples. Tobit Models*. Online-Dokument (abgerufen am 30.11.2014): <http://www.ats.ucla.edu/stat/R/dae/tobit.htm>
- Venables, W.N., Smith, D. M. & R Core Team. (2016). *An Introduction to R. Version 3.2.5*. Verfügbar über die Hilfe zum RGui.
- Weisberg, S. (1985). *Applied linear regression* (2nd ed.). New York: Wiley.
- Wickham, H. (2009). *ggplot2: Elegant graphics for data analysis*. New York: Springer.
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis* (2nd ed.). New York: Springer.
- Wilkinson, L. (2005). *The Grammar of Graphics* (2nd ed.). New York: Springer.
- Wollschläger, D. (2010). *Grundlagen der Datenanalyse mit R*. Heidelberg: Springer.
- Yang, S., Keller, F.B. & Zheng, L. (2016). *Social Network Analysis*. Los Angeles (CA): Sage.

dim..... 64
dim()..... 62, 67
dir()..... 39
dnorm()..... 108
Dreipunktargument..... 51
Dynamische Typisierung..... 55

E

Echo..... 81
edge_attr()..... 194
edge_density()..... 196
edit()..... 84
element_blank()..... 160
element_line()..... 160
element_rect()..... 159
element_text()..... 160
EMF..... 108, 161
encoding..... 94
EndDataStep()..... 37
Enhanced Windows Metafile..... 108
Entladen von Paketen..... 47
Erweiterungsbundle..... 10, 12
Erweiterungskommandos..... 9
example()..... 49
Excel..... 86
expand
 ggplot2..... 146
expression()..... 114, 137, 145, 154

F

facet_grid()..... 157
facet_wrap()..... 157
Facettierung..... 157
factor()..... 59
factorMode..... 31
Faktor..... 59
Fallauswahl..... 104
Farbe
 Hintergrundfarbe..... 115
 Vordergrundfarbe..... 115
 Zeichenfarbe..... 115
Farben..... 115
Farbskalen
 ggplot2..... 147
Fehlende Werte..... 32, 70, 188
file.choose()..... 99
fix()..... 84
Fleiss-Kappa..... 44
foreign..... 97
for-Schleife..... 79
FUN..... 63
fun.data..... 181
Funktionen..... 51
Funktionen definieren..... 82
Funktionsplots..... 136

G

gc()..... 31, 54
Generische Funktionen..... 81
Geom..... 139
geom_bar()..... 176
geom_bin2d()..... 169
geom_boxplot()..... 174
geom_density()..... 171

geom_density2d()..... 167
geom_errorbar()..... 180
geom_histogram()..... 169
geom_pointrange..... 181
geom_smooth()..... 143, 163
getwd()..... 39
ggplot2..... 107, 138
ggsave()..... 161
Gleichungssystem..... 198
GPL..... 107
Grad der Zentralität..... 195
Gradienten..... 147
Grafikfenster..... 108
Grafikparameter..... 114
Graphics Production Language..... 107
graphics.off()..... 112
grep()..... 104
Gruppiertes Balkendiagramm..... 183
Gruppiertes Streudiagramm..... 125, 163
Guides..... 144

H

help()..... 48, 92
help.search()..... 49
Hexagonale Gruppierung..... 168
High Level - Grafikfunktionen..... 134
Hintergrundfarbe..... 115
hist()..... 102, 103, 134
Histogramm
 ggplot2..... 169
 traditionelle R-Grafik..... 134

I

I() 67, 189
identify()..... 127
if-Anweisung..... 79
if-else - Anweisung..... 79
ifelse()..... 100
igraph..... 191
Indexmatrix..... 75
Indexvektoren..... 73
install.packages()..... 45
install.views()..... 47
Interquartilsabstand..... 188
intersect()..... 197
Inverse Matrix..... 198
IQR()..... 188
irr (R-Paket)..... 44
is.na()..... 33, 71
is.nan()..... 33
is.typ()..... 55

J

Jitter
 ggplot2..... 154
 traditionelle R-Grafik..... 128

K

Kommentare..... 50

L

labels-Argument in scale-Funktionen..... 145

labs() 153, 162, 170
 Langformat 182
 lattice 107
 legend() 120, 126
 Legende
 ggplot2 151
 Traditionelle Grafik 120
 length() 53
 levels 60
 levels() 59
 lib.loc 43
 library() 33, 43
 limits-Argument in scale-Funktionen 145
 lines() 120, 124, 134, 135
 Liniendiagramm
 ggplot2 185
 traditionelle R-Grafik 122
 Linienstärke 116
 Linientyp 118
 list() 64
 lm() 123
 load() 41
 Logische Operatoren 76
 Logischer Indexvektor 74, 104
 lowess() 124
 ls() 53

M

MARGIN 63
 MASS (R-Paket) 18, 21
 MATRIX 91
 matrix() 61
 max() 187
 mean() 187
 median() 188
 Mehrebenenanalyse 105
 melt() 182
 memory.size() 54
 Mengenlehre 197
 merge() 105
 Messniveau 67
 Metafile-Format 108, 110, 161
 mfrow 135
 Microsoft Office 110
 Mikadodiagramm 186
 min() 187
 missingValueToNA 33
 mode() 55
 Modulo 75
 M-Schätzer 18
 Multiplikation von Matrizen 198

N

NA 32, 56
 na.omit() 71
 na.rm 72, 106
 name-Argumente in scale-Funktionen 145
 Namensargumente 51
 names() 57, 65, 68
 NaN 32, 71
 ncol 61
 ncol() 62, 67
 ncvTest 23
 Nominaler Faktor 59
 Normalverteilte Zufallszahlen 102

nrow 61
 nrow() 62, 67
 NULL 69, 98
 numeric() 56, 100

O

objects() 40
 options() 41
 opts() 159
 order() 58
 Ordinaler Faktor 60

P

Pakete 43
 aktualisieren 46
 entladen 47
 installieren 44
 laden 43
 zitieren 47
 par() 114
 paste() 134
 Permutation 106
 Persönliche Bibliothek 45
 persp() 137
 Pivot-Tabelle 35
 plot() 116
 Plot-Objekt 161
 Plot-Typen 117
 Polychorische und polyseriale Korrelation 23
 position_dodge() 185
 position_jitter() 155
 Positionsanpassungen 154
 Positionsargumente 51
 predict() 124
 print() 29, 52
 prop.table() 188
 Python-Essentials 11

Q

q() 54
 qplot() 139
 quantile() 188
 quit() 54

R

R Commander 86
 rank() 58
 rbind() 61
 rbinom() 103
 Rboxplot 132
 rchisq() 102
 Rcmdr 86
 RData 40
 read.csv2() 95
 read.delim() 93
 read.spss() 97
 read.table() 95
 recode() 100
 Recycling-Regel 77
 reencode 98
 remove() 31
 rep() 59
 require() 44

reshape2.....	182
RGB-Wert.....	115
RGui.....	38, 78
Rhistory.....	40
rlm().....	18, 21
rm().....	31, 41, 53
rnorm().....	102
Robuste Regression.....	18
row.names().....	68
rownames().....	62
Rprofile.site.....	42
R-Skripte.....	80

S

S	8
sample().....	106
save(),.....	41
save.image().....	41
scale_colour_brewer().....	148
scale_colour_gradient().....	147
scale_colour_gradient2().....	148
scale_colour_gradientn().....	148
scale_colour_hue().....	148
scale_colour_manual().....	148, 149
scale_fill_brewer().....	148
scale_fill_gradient().....	147
scale_fill_gradient2().....	148
scale_fill_gradientn().....	148
scale_fill_hue().....	148
scale_fill_manual().....	148, 185
scale_shape_manual().....	150
scale_size_continuous().....	156
scale_x_log10().....	147
scale-Funktionen.....	145
Schnitt von Mengen.....	197
search().....	43, 54, 69
Sekundärstichprobe.....	106
seq().....	56, 77
Sequenzoperator.....	30, 56, 77
Serienargumente.....	51
setdiff().....	197
setwd().....	39
show.legend.....	152
Skalen	
ggplot2.....	144
Skalierung	
Tobit-Regression.....	27
Skripte.....	80
smoothScatter().....	131
SNA.....	191
Solide Regression.....	21
solve().....	198
sort().....	58
source().....	81, 83
Soziale Netzwerkanalyse.....	191
Spaltendominanz.....	61
SPD-Dateien.....	16
spsspivottable.Display().....	35
Stärke Zentralität.....	195
stat_binhex().....	168
stat_boxplot().....	174
stat_sum().....	156
str().....	32, 62, 67
strength().....	195
Streudiagramm	
ggplot2.....	161
traditionelle R-Grafik.....	122

Struktur einer Matrix.....	62
subset().....	104
Suchpfad.....	69
sum().....	71, 187
summary().....	34, 81, 187
Symbole.....	117

T

t()	63
t.test().....	190
table().....	188
Task Views.....	47
Teilmenge.....	197
text().....	114
Textdateien.....	92
lesen.....	93
schreiben.....	96
theme().....	154, 159
theme_bw().....	158
theme_grey().....	158
theme_set().....	159
Themes.....	158
to.data.frame.....	97
Tobit-Regression.....	25
trans	
ggplot2.....	147
Transferargumente.....	51
Transparenz.....	110, 115
Transponieren.....	63
Trellis-Grafiken.....	107
t-Test.....	190

U

Umkodieren.....	99
union().....	197
unique().....	197
units.....	161
update.packages().....	46
update.views().....	47
use.value.labels.....	98

V

values.....	149
var().....	188
Variablennamen.....	50
Vektor.....	55
Verbindungszentralität.....	195
Vereinigung von Mengen.....	197
Vergleichsoperatoren.....	76
Verkettungsfunktion.....	52
vertex_attr().....	193
vjust.....	176
Vordergrundfarbe.....	115

W

Waagerechte Balken.....	179
Wertargumente.....	82
which().....	74, 77
Wiederholungsanweisung.....	79
win.metafile().....	108
windowsFonts().....	115
with().....	70
Workspace.....	40

write.csv2()	97
write.foreign().....	98
write.table()	96

X

xlab()	153
xlim()	153

Y

ylab()	153
ylim()	153

Z

Zeichenfarbe	115
Zitieren von Paketen	47
Zufallszahlen	
Binomialverteilung	103
Chiquadratverteilung	102
Normalverteilung	102
Zusammenfassungsfunktionen.....	180, 181
Zuweisungsoperatoren	78