

Statistical Computing

Einführung in

Günther Sawitzki
StatLab Heidelberg

17. März 2008

noch in Vorbereitung

E-mail address: `gs@statlab.uni-heidelberg.de`

URL: `http://www.statlab.uni-heidelberg.de/projects/r/`

Key words and phrases. statistical computing, S programming language, R programming, data analysis, exploratory statistics, residual diagnostics

17. März 2008.

Inhaltsverzeichnis

Einleitung	v
0.1. Was ist R?	v
0.2. Referenzen	vii
0.3. Umfang und Aufbau des Kurses	vii
0.4. Dank	viii
0.5. Literatur und weitere Hinweise	viii
Kapitel 1. Grundlagen	1-1
1.1. Programmierung: Konventionen	1-1
1.2. Erzeugung von Zufallszahlen und Mustern	1-4
1.2.1. Zufallszahlen	1-4
1.2.2. Muster	1-8
1.3. Fallstudie: Verteilungsdiagnostik	1-9
1.3.1. Erster Durchgang zu Beispiel 1.1: Verteilungsfunktion	1-11
1.3.2. Erster Durchgang zu Beispiel 1.2: Histogramm	1-14
1.3.3. Zweiter Durchgang zu Beispiel 1.1: Verteilungsfunktion	1-20
1.3.4. Zweiter Durchgang zu Beispiel 1.2: Histogramm	1-27
1.4. Momente und Quantile	1-33
1.5. Ergänzungen	1-37
1.5.1. Ergänzung: Zufallszahlen	1-37
1.5.2. Ergänzung: Grafische Vergleiche	1-37
1.5.3. Ergänzung: Grafik-Aufbereitung	1-42
1.5.4. Ergänzung: Funktionen	1-43
1.5.5. Ergänzung: Das Innere von R	1-46
1.5.6. Ergänzung: Pakete	1-47
1.6. Statistische Zusammenfassung	1-49
1.7. Literatur und weitere Hinweise:	1-50
Kapitel 2. Regression	2-1
2.1. Allgemeines Regressionsmodell	2-1
2.2. Lineares Model	2-2
2.2.1. Faktoren	2-4
2.2.2. Kleinste-Quadrate-Schätzung	2-5
2.2.3. Weitere Beispiele für lineare Modelle	2-14
2.2.4. Modellformeln	2-15
2.2.5. Gauß-Markov-Schätzer und Residuen	2-17
2.3. Streuungszerlegung und Varianzanalyse	2-18
2.4. Simultane Schätzung	2-25
2.4.1. Scheffé's Konfidenz-Bänder	2-25
2.4.2. Tukey's Konfidenz-Intervalle	2-27

2.5.	Nichtparametrische Regression	2-32
2.5.1.	Zwischenspiel: verallgemeinerte lineare Modelle	2-33
2.5.2.	Lokale Regression	2-34
2.6.	Ergänzungen	2-37
2.6.1.	Ergänzung: Diskretisierungen	2-37
2.6.2.	Ergänzung: Externe Daten	2-37
2.6.3.	Ergänzung: Software-Test	2-38
2.6.4.	R-Datentypen	2-39
2.6.5.	Klassen und polymorphe Funktionen	2-39
2.6.6.	Extraktor-Funktionen	2-40
2.7.	Statistische Zusammenfassung	2-41
2.8.	Literatur und weitere Hinweise:	2-41
Kapitel 3.	Vergleich von Verteilungen	3-1
3.1.	Shift/Skalenfamilien	3-3
3.2.	QQ-Plot, PP-Plot	3-5
3.2.1.	Kolmogorov Smirnov Tests	3-10
3.3.	Tests auf Shift	3-10
3.4.	Güte	3-17
3.4.1.	Theoretische Güte	3-17
3.4.2.	Simulation der Güte	3-21
3.4.3.	Quantilschätzung durch Simulation	3-24
3.5.	Qualitative Eigenschaften von Verteilungen	3-26
3.6.	Ergänzungen	3-27
3.7.	Statistische Zusammenfassung	3-29
3.8.	Literatur und weitere Hinweise:	3-29
Kapitel 4.	Dimensionen 1, 2, 3, ..., ∞	4-1
4.1.	Ergänzungen	4-1
4.2.	Dimensionen	4-4
4.3.	Selektionen	4-5
4.4.	Projektionen	4-6
4.4.1.	Randverteilungen und Scatterplot-Matrix	4-7
4.4.2.	Projection Pursuit	4-11
4.4.3.	Projektionen für dim 1, 2, 3, ..., 7	4-13
4.4.4.	Parallel-Koordinaten	4-13
4.5.	Schnitte, bedingte Verteilungen und Coplots	4-14
4.6.	Transformationen und Dimensionsreduktion	4-21
4.7.	Höhere Dimensionen	4-26
4.7.1.	Linearer Fall	4-26
4.7.2.	Nichtlinearer Fall	4-28
4.7.3.	“Curse of Dimension”	4-32
4.7.4.	Fallstudie	4-32
4.8.	Hohe Dimensionen	4-43
4.9.	Statistische Zusammenfassung	4-45
4.10.	Literatur und weitere Hinweise:	4-45
R	als Programmiersprache: Übersicht	A-1
A.1.	Hilfe und Information	A-1

A.2.	Namen und Suchpfade	A-1
A.3.	Anpassung	A-2
A.4.	Basis-Datentypen	A-3
A.5.	Ausgabe von Objekten	A-5
A.6.	Inspektion von Objekten	A-7
A.7.	Inspektion des Systems	A-9
A.8.	Komplexe Datentypen	A-11
A.9.	Zugriff auf Komponenten	A-13
A.10.	Tabellen-Transformationen	A-15
A.11.	Operatoren	A-17
A.12.	Funktionen	A-19
A.13.	Debugging und Profiling	A-21
A.14.	Kontrollstrukturen	A-23
A.15.	Verwaltung und Anpassung	A-25
A.16.	Ein- und Ausgabe in Ausgabeströme	A-27
A.17.	Externe Daten	A-29
A.18.	Libraries, Pakete	A-31
A.19.	Lineare Algebraoperatoren	A-33
A.20.	Modell-Beschreibungen	A-35
A.21.	Grafik-Funktionen	A-37
A.21.1.	high level Grafik	A-37
A.21.2.	low level Grafik	A-37
A.21.3.	Annotationen und Legenden	A-38
A.21.4.	Grafik-Parameter und Layout	A-39
A.22.	Einfache Statistische Funktionen	A-41
A.23.	Verteilungen, Zufallszahlen, Dichten. . .	A-43
A.24.	Verarbeitung von Ausdrücken	A-45
Literaturverzeichnis		Literatur-1
Index		Index-1

Einleitung

Diese Einführung in R ist als Arbeitsmaterial in einem Kompaktkurs oder zum Selbststudium gedacht. Der Kurs wendet sich an Studierende mit Grundkenntnissen in Stochastik. Begriffe wie Verteilungsfunktion, Quantil, Erwartungswert und Varianz und die damit verbundenen einfachen Eigenschaften werden vorausgesetzt. Ebenso sollten klassische Verteilungen wie Binomial-, Uniform- und Gaußverteilung sowie daraus abgeleitete Verteilungen und deren asymptotisches Verhalten bekannt sein. Kenntnisse in Statistik selbst werden nicht vorausgesetzt. Sie werden in diesem Kurs aber auch nicht vermittelt. Der Kurs konzentriert sich auf die “Computing”-Aspekte. Dabei werden statistische Betrachtungsweisen und Konzepte zwar eingeführt und diskutiert. Für eine eingehendere Diskussion wird aber auf die Statistik-Vorlesungen verwiesen.

Kenntnisse in der Rechnerbenutzung und zumindest oberflächliche Kenntnisse von Programmierkonzepten wie Variable, Schleifen und Funktionen werden vorausgesetzt. Weitergehende Kenntnisse werden nicht vorausgesetzt, aber auch nicht vermittelt. Der Kurs führt in die Benutzung von R als Anwender ein. Für eingehendere Diskussion der Computing-Aspekte wird auf die Arbeitsgemeinschaft “Computational Statistics” verwiesen.

`<http://www.statlab.uni-heidelberg.de/studinfo/compstat/>`

0.1. Was ist R?

R ist eine Programmiersprache, und auch der Name eines Software-Systems, das diese Sprache implementiert. Die Programmiersprache R ist eine für die Statistik und für stochastische Simulation entwickelte Programmiersprache, die mittlerweile zum Standard geworden ist. Genau genommen müsste man hier unterscheiden: Die Sprache heißt S, ihre Implementierung und das System heißen R. Die ursprünglichen Autoren von S sind John M. Chambers, R. A. Becker und A. R. Wilks, AT & T Bell Laboratories, Statistics Research Department. Die Sprache und ihre Entwicklung sind in einer Reihe von Büchern dokumentiert, nach ihrem Umschlag häufig als das weiße ([CH92]), blaue ([BCW88]) und grüne Buch ([Cha98]) bezeichnet.

Die AT & T-Implementierung von S war lange Zeit die “Referenz” für die Sprache S. Heute gibt es S als kommerzielles System S-Plus `<http://www.insightful.com/>` (basierend auf der AT & T-Implementierung) sowie als frei verfügbare Version R, auch “Gnu S” genannt¹ `<http://www.r-project.org/>`.

¹R heißt nur zufällig so, wie auch zufälligerweise die Vornamen der Originalautoren (Ross Ihaka & Robert Gentleman) mit R beginnen.



Mittlerweile hat sich R zur Referenz-Implementierung entwickelt. Wesentliche Präzisierungen, und - falls notwendig - auch Modifikationen der Sprache werden durch R definiert. Der Einfachheit halber sprechen wir hier und in den folgenden Kapiteln von der Sprache R, auch wenn es genauer heißen müsste: die Sprache S in der R-Implementierung.

R ist eine interpretierte Programmiersprache. Anweisungen in R werden unmittelbar ausgeführt. R beinhaltet neben den ursprünglichen Elementen von S eine Reihe von Erweiterungen, zum Teil um Entwicklungen in der Statistik angemessen zu berücksichtigen, zum Teil um experimentelle Möglichkeiten zu eröffnen. Parallel dazu gibt es Weiterentwicklungen der S-Sprache.

Die (2008) aktuelle Version von R ist R 2.x. Diese Version ist weitgehend kompatibel mit den Vorläuferversionen R 1.x. Die wesentlichen Veränderungen sind im Inneren des Systems. Für den Anfang gibt es praktisch keinen Unterschied zu R 1.x. Für den fortgeschrittenen Nutzer gibt es drei wesentliche Neuerungen:

- *Grafik*: Das Basis-Grafiksystem von R implementiert ein Modell, dass an der Vorstellung von Stift und Papier orientiert ist. Ein Grafik-Port (Papier) wird eröffnet und darauf werden Linien, Punkte/Symbole gezeichnet. Mit R 2.x gibt es zusätzlich ein zweites Grafiksystem, dass an einem Kamera/Objekt-Modell orientiert ist. Grafische Objekte in unterschiedlicher Lage und Richtung werden in einem visuellen Raum abgebildet.
- *Packages*: Das ursprüngliche System von R hat eine lineare Geschichte und einen einheitlichen Arbeitsraum. Mit R 2.x gibt es eine verbesserte Unterstützung von "Paketen", die in sich abgeschirmt werden können. Dazu dienen Sprachkonzepte wie "name spaces", aber auch unterstützende Werkzeuge.
- *Internationalisierung*: Die ursprüngliche Implementierung von R setzte Englisch als Sprache und ASCII als Zeichensatz voraus. Seit R 2.x gibt es umfassende Unterstützung für andere Sprachen und Zeichensätze. Dies ermöglicht es, "lokalisierte" Versionen zu erstellen. Derzeit ist man bei Kommandos, Ausgaben und Erklärungen jedoch noch auf Englisch angewiesen.

Zwei Aspekte sind in R nur unzureichend berücksichtigt: der interaktive Zugriff und die Einbettung in eine vernetzte Umgebung. Diese und weitere Aspekte sind Bestandteil von Omegahat - eines Versuchs, ein System der nächsten Generation zu entwickeln, das auf den Erfahrungen mit R aufbaut. Diese mehr experimentellen Arbeiten werden unter <http://www.omegahat.org/> bereitgestellt. Schon R bietet einfache Möglichkeiten, Prozeduren aus anderen Sprachen wie C und Fortran aufzurufen. Omegahat erweitert diese Möglichkeiten und bietet einen direkten Zugang zu Java, Perl . . .

Eine Java-basierte grafische Oberfläche ist als JGR unter <http://stats.math.uni-augsburg.de/software/> zugänglich. Dort findet sich als *iplots* auch eine Sammlung von interaktiven Displays für R.

Aktuelle Entwicklungen zu R finden sich in <http://r-forge.r-project.org/>. Zahlreiche hilfreiche Erweiterungen sind auch unter <http://www.bioconductor.org/> zu finden.

0.2. Referenzen

R ist für die praktische Arbeit in der Statistik entworfen. Nützlichkeit hat oft Vorrang vor prinzipiellen Design-Überlegungen. Als Folge ist eine systematische Einführung in R nicht einfach. Stattdessen wird ein verschlungener Pfad gewählt: Fallstudien und Beispiele, an die sich systematische Übersichten anschließen. Für die praktische Arbeit sollte auf das reichhaltige Online-Material zu R zugegriffen werden. Ein erster Zugriffspunkt sind dabei die “frequently asked questions” (FAQ) <http://www.cran.r-project.org/faqs.html>. “An Introduction to R” ([R D07a]) ist die “offizielle” Einführung. Diese Dokumentation und andere Manuale sind unter <http://www.cran.r-project.org/manuals.html> bereitgestellt.

R-Prozeduren sind zum Teil im Basis-System enthalten. Andere Prozeduren müssen aus Bibliotheken hinzugeladen werden. Eine Reihe von Bibliotheken ist in der Standard-Distribution von R enthalten und muss lediglich aktiviert werden. Die technischen Hinweise dazu sind jeweils angegeben. Speziellere Bibliotheken müssen evtl. hinzu geladen werden. Die erste Quelle dafür ist <http://www.cran.r-project.org/src/contrib/PACKAGES.html>.

Größere Unterschiede gibt es bei unterschiedlichen Versionen von S-Plus. S-Plus 4.x und S-Plus 2000 benutzen S Version 3 und sind weitestgehend mit R kompatibel. S-Plus 5 ist eine Implementierung von S Version 4 mit Änderungen, die eine Sonderbehandlung bei der Programmierung benötigen. Auf diese Besonderheiten wird hier nicht eingegangen. Informationen zu S-Plus findet man bei <http://www.insightful.com/>.

0.3. Umfang und Aufbau des Kurses

R beinhaltet in der Basis-Version mehr als 1500 Funktionen - zu viele, um sie in einem Kurs zu vermitteln, und zu viel, um sie sinnvollerweise zu lernen. Der Kurs kann nur dazu dienen, den Zugang zu R zu eröffnen.

Teilnehmerkreise können aus unterschiedlichem Hintergrund kommen und unterschiedliche Vorbedingungen mitbringen. Gerade für jüngere Schüler oder Studenten kann ein reiner Programmierkurs, der sich auf die technischen Grundlagen konzentriert, angemessen sein. Für diese Teilnehmer ist dieser Kurs nicht geeignet. Für Fortgeschrittene stellt sich eher die Frage nach einer sinnvollen Einordnung und nach dem Hintergrund. Hierauf zielt der vorliegende Kurs. Das “technische” Material bildet das Skelett. Daneben wird versucht, den Blick auf statistische Fragestellungen zu richten und das Interesse am Hintergrund zu wecken. Der Kurs soll Appetit auf die Substanz wecken, die eine fundierte statistische Vorlesung bieten kann.

Das hier bereitgestellte Material besteht zunächst aus einer thematisch geordneten Sammlung, in der anhand von Beispiel-Fragestellungen illustriert wird, wie ein erster Zugang mit R erfolgen kann. Hinzu kommt eine Zusammenstellung von Sprachbestandteilen und Funktionen, die als Orientierungshilfe für das umfangreiche in R enthaltene Informationsmaterial dient. Für die praktische Arbeit sind die Online-Hilfen und Manuale die erste Informationsquelle.

Der Kurs kann bei einer Auswahl der Aufgaben in etwa vier Tagen durchgeführt werden. Konzeptuell ist er eine viertägige Einführung in die Statistik mit den Themenbereichen

- Ein-Stichprobenanalyse und Verteilungen
- Regression
- Zwei- oder Mehr-Stichprobenanalysen
- Multivariate Analysen

Eine großzügigere Zeit für die Übungsaufgaben wird empfohlen (ein Halbtage zusätzlich für einführende Aufgaben, ein Halbtage zusätzlich für eine der Projektaufgaben). Mit dieser Zeit kann der Kurs als Block in einer Woche durchgeführt werden, wenn im Anschluss die Möglichkeit geschaffen wird, die aufgetreten Fragen zu beantworten und das geweckte Interesse am statistischen Hintergrund zu vertiefen.

Für ein anschließendes vertiefendes Selbststudium von R als Programmiersprache wird ([VR00]) empfohlen.

Beispiele und Eingaben im Text sind so formatiert, dass sie mit “Cut & Paste” übernommen und als Programmeingabe verwandt werden können. Deshalb sind bei Programmbeispielen im Text bisweilen Satzzeichen fortgelassen, und Eingabebeispiele werden ohne “Prompt” gezeigt. Einem Beispiel

Beispiel 0.1:

Eingabe

$1 + 2$

Ausgabe

3

entspricht auf dem Bildschirm etwa

```
> 1+2
[1] 3
>
```

wobei anstelle des Prompt-Zeichens ”>” je nach Konfiguration auch ein anderes Zeichen erscheinen kann.

0.4. Dank

Zu danken ist dem R core team für die Kommentare und Hinweise. Besonderen Dank an Friedrich Leisch vom R core team sowie an Antony Unwin, Univ. Augsburg.

0.5. Literatur und weitere Hinweise

[R D07a] R Development Core Team (2000-2007): An Introduction to R.
 Siehe: <http://www.r-project.org/manuals.html>.

- [**R D07d**] R Development Core Team (2000-2007): R Reference Manual.
Siehe: <http://www.r-project.org/manuals.html>.
- The Omega Development Group (2000): Omega.
Siehe: <http://www.omegahat.org/>.
- [**BCW88**] Becker, R.A.; Chambers, J.M.; Wilks, A.R. (1988): The New S Language.
New York: Chapman and Hall.
- [**CH92**] Chambers, J.M.; Hastie, T.J. (eds) (1992): Statistical Models in S. New York:
Chapman and Hall.
- [**Cle93**] Cleveland, W.F. (1993): Visualizing Data. Summit: Hobart Press.
- [**VR02**] Venables, W.N.; Ripley, B.D. (2002): Modern Applied Statistics with S.
Heidelberg:Springer.
Siehe: <http://www.stats.ox.ac.uk/pub/MASS4/>.
- [**VR00**] Venables, W.N.; Ripley, B.D. (2000): Programming in S. Heidelberg:Springer.
Siehe: <http://www.stats.ox.ac.uk/pub/MASS3/Sprog>.

KAPITEL 1

Grundlagen

1.1. Programmierung: Konventionen

Wie jede Programmiersprache hat R bestimmte Konventionen. Hier die ersten Grundregeln.

<i>R-Konventionen</i>	
Zahlen	<p>Dezimaltrenner ist ein Punkt. Zahlen können im Exponentialformat eingegeben werden; der Exponentialteil wird mit <i>E</i> eingeleitet. Zahlen können komplex sein. Der Imaginärteil wird mit <i>i</i> gekennzeichnet.</p> <p><i>Beispiel:</i> 1 2.3 3.4E5 6i+7.8</p> <p>Zahlen können auch die Werte <i>Inf</i>, <i>-Inf</i>, <i>NaN</i> für “not a number” und <i>NA</i> für “not available” = fehlend annehmen.</p> <p><i>Beispiel:</i> 1/0 ergibt <i>Inf</i>; 0/0 ergibt <i>NaN</i>.</p>
Zeichenketten	<p>Zeichenketten (Strings) werden zu Beginn und zu Ende durch " oder ' begrenzt.</p> <p><i>Beispiel:</i> "ABC" 'def' "gh'ij"</p>

Damit die folgenden Beispiele nicht zu simpel werden, greifen wir hier vor: in R ist *a:b* eine Sequenz von Zahlen von *a* bis höchstens *b* in Schritten von 1 bzw. -1.

<i>R-Konventionen</i>	
Objekte	<p>Die Datenbausteine in R sind Objekte. Objekte können Klassen zugeordnet werden.</p> <p><i>Beispiel:</i> Die Basis-Objekte in R sind Vektoren.</p>

(Fortsetzung)→

R-Konventionen (Fortsetzung)	
Namen	<p>R-Objekte können Namen haben. Dann kann anhand ihres Namens auf sie zugegriffen werden.</p> <p>Namen beginnen mit einem Buchstaben oder einem Punkt, gefolgt von einer Folge von Buchstaben, Ziffer, oder den Sonderzeichen <code>_</code> oder <code>.</code></p> <p>Beispiele: <code>x</code> <code>y_1</code></p> <p>Groß- und Kleinschreibung werden unterschieden.</p> <p>Beispiele: <code>Y87</code> <code>y87</code></p>
Zuweisungen	<p>Zuweisungen haben die Form</p> <p>Aufruf: <code>Name <- Wert</code> oder alternativ <code>Name = Wert</code>.</p> <p>Beispiel: <code>a <- 10</code> <code>x <- 1:10</code></p>
Abfragen	<p>Wird nur der Name eines Objekts eingegeben, so wird der Wert des Objekts ausgegeben.</p> <p>Beispiel: <code>x</code></p>
Indizes	<p>Auf Vektorkomponenten wird über Indizes zugegriffen. Die Index-Zählung beginnt mit 1.</p> <p>Beispiel: <code>x[3]</code></p> <p>Dabei können für die Indizes auch symbolische Namen oder Regeln verwandt werden.</p> <p>Beispiele: <code>x[-3]</code> <code>x[x^2 < 10]</code> <code>a[1]</code></p>

Hilfe und Inspektion	
Hilfe	<p>Dokumentation und Zusatzinformation für ein Objekt kann mit <code>help</code> angefordert werden.</p> <p>Aufruf: <code>help(Name)</code></p> <p>Beispiele: <code>help(exp)</code> <code>help(x)</code></p> <p>Alternative Form <code>?Name</code></p> <p>Beispiele: <code>?exp</code> <code>?x</code></p>

(Fortsetzung)→

Hilfe und Inspektion (Fortsetzung)	
Inspektion	<p><code>help()</code> kann nur vorbereitete Dokumentation bereitstellen. <code>str()</code> kann den aktuellen Zustand inspizieren und darstellen. <i>Aufruf:</i> <code>str(Object, ...)</code> <i>Beispiele:</i> <code>str(x)</code></p>

R-Konventionen	
Funktionen	<p>Funktionen in R werden aufgerufen in der Form <i>Aufruf:</i> <code>Name(Parameter ...)</code> <i>Beispiel:</i> <code>e_10 <- exp(10)</code></p> <p>Diese Konvention gilt selbst, wenn keine Parameter vorhanden sind. <i>Beispiel:</i> Um R zu verlassen ruft man eine "Quit"-Prozedur auf <code>q()</code>.</p> <p>Parameter werden sehr flexibel gehandhabt. Sie können Default-Werte haben, die benutzt werden, wenn kein expliziter Parameter angegeben ist. <i>Beispiele:</i> <code>log(x, base = exp(1))</code></p> <p>Funktionen können <i>polymorph</i> sein. Die aktuelle Funktion wird dann durch die Klasse der aktuellen Parameter bestimmt. <i>Beispiele:</i> <code>plot(x)</code> <code>plot(x, x^2)</code> <code>summary(x)</code></p>
Operatoren	<p>Für Vektoren wirken Operatoren auf jede Komponente der Vektoren. <i>Beispiel:</i> Für Vektoren y, z ist $y*z$ ein Vektor, der komponentenweise das Produkt enthält.</p> <p>Operatoren sind spezielle Funktionen. Sie können auch in Präfix-Form aufgerufen werden. <i>Beispiel:</i> <code>"+"(x, y)</code></p> <p>In Situationen, in denen die Operanden nicht gleiche Länge haben, wird der kürzere Operand zyklisch wiederholt. <i>Beispiel:</i> <code>(1:2)*(1:6)</code></p>

Wir beschäftigen uns im folgenden mit statistischen Methoden. Wir benutzen die Methoden zunächst in Simulationen, d.h. mit synthetischen Daten, deren Erzeugung wir weitgehend unter Kontrolle haben. Das erlaubt es uns, Erfahrung mit den Methoden zu gewinnen und sie kritisch zu beurteilen. Erst dann benutzen wir die Methoden zur Analyse von Daten.

1.2. Erzeugung von Zufallszahlen und Mustern

1.2.1. Zufallszahlen. Die Funktion `runif()` erlaubt die Erzeugung von uniform verteilten Zufallsvariablen. Mit `help(runif)` oder `?runif` erhalten wir Informationen, wie die Funktion benutzt werden kann:

help(runif)

Uniform

The Uniform Distribution

Description.

These functions provide information about the uniform distribution on the interval from `min` to `max`. `dunif` gives the density, `punif` gives the distribution function `qunif` gives the quantile function and `runif` generates random deviates.

Usage.

```
dunif(x, min=0, max=1, log = FALSE)
punif(q, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
qunif(p, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
runif(n, min=0, max=1)
```

Arguments.

<code>x,q</code>	vector of quantiles.
<code>p</code>	vector of probabilities.
<code>n</code>	number of observations. If <code>length(n) > 1</code> , the length is taken to be the number required.
<code>min,max</code>	lower and upper limits of the distribution.
<code>log, log.p</code>	logical; if TRUE, probabilities <code>p</code> are given as $\log(p)$.
<code>lower.tail</code>	logical; if TRUE (default), probabilities are $P[X \leq x]$, otherwise, $P[X > x]$.

Details.

If `min` or `max` are not specified they assume the default values of 0 and 1 respectively.

The uniform distribution has density

$$f(x) = \frac{1}{\max - \min}$$

for $\min \leq x \leq \max$.

For the case of $u := \min == \max$, the limit case of $X \equiv u$ is assumed.

References.

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

See Also.

.`Random.seed` about random number generation, `rnorm`, etc for other distributions.

Examples.

```
u <- runif(20)

## The following relations always hold :
punif(u) == u
dunif(u) == 1

var(runif(10000))#- ~ = 1/12 = .08333
```

Diese Hilfsinformation sagt uns: Als Parameter für `runif()` muss die Anzahl `n` der zu generierenden Zufallswerte angegeben werden. Als weitere Parameter für `runif()` können das Minimum und das Maximum des Wertebereichs angegeben werden. Geben wir keine weiteren Parameter an, so werden die Default-Werte `min = 0` und `max = 1` genommen. Z. B. `runif(100)` erzeugt einen Vektor mit 100 uniform verteilten Zufallsvariablen im Bereich (0, 1). Der Aufruf `runif(100, -10, 10)` erzeugt einen Vektor mit 100 uniform verteilten Zufallsvariablen im Bereich (-10, 10). Die zusätzlichen Parameter können in der definierten Reihenfolge angegeben werden, oder mithilfe der Namen spezifiziert werden. Bei Angabe des Namens kann die Reihenfolge frei gewählt werden. Anstelle von `runif(100, -10, 10)` kann also `runif(100, min = -10, max = 10)` oder `runif(100, max = 10, min = -10)` benutzt werden. Dabei können auch ausgewählt einzelne Parameter gesetzt werden. Wird zum Beispiel das Minimum nicht angegeben, so wird für das Minimum der Default-Wert eingesetzt: die Angabe von `runif(100, max = 10)` ist gleichwertig mit `runif(100, min = 0, max = 10)`. Der besseren Lesbarkeit halber geben wir oft die Namen von Parametern an, auch falls es nicht nötig ist.

Jeder Aufruf von `runif()` erzeugt 100 neue uniforme Zufallszahlen. Wir können diese speichern.

```
x <- runif(100)
```

erzeugt einen neuen Vektor von Zufallszahlen und weist ihn der Variablen `x` zu.

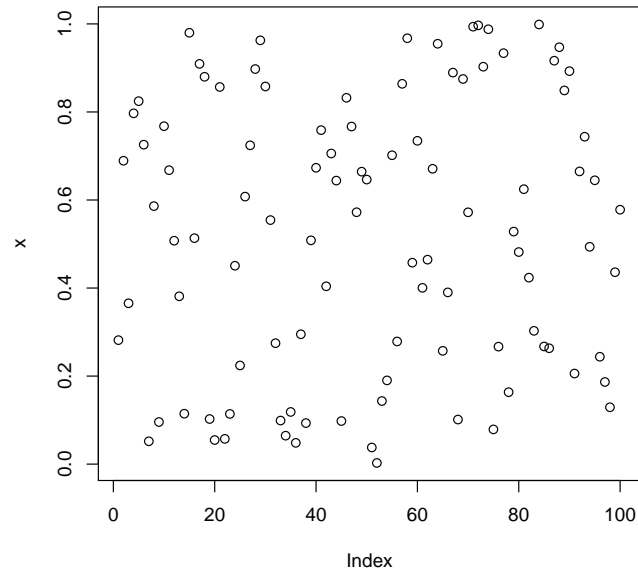
```
x
```

gibt jetzt dessen Werte aus; damit können wir die Resultate inspizieren. Eine grafische Darstellung, den *Serienplot* - einen Scatterplot der Einträge in `x` gegen den laufenden Index, erhalten wir mit

```
plot(x)
```

Beispiel 1.1:*Eingabe*

```
x <- runif(100)
plot(x)
```



Aufgabe 1.1	
	<p>Experimentieren Sie mit den Plots und <code>runif()</code>. Sind die Plots Bilder von Zufallszahlen?</p> <p>Genauer: Akzeptieren Sie die Plots als Bilder von 100 unabhängigen Realisationen von uniform auf $(0, 1)$ verteilten Zufallszahlen?</p> <p>Wiederholen Sie das Experiment und versuchen Sie, die Argumente, die für oder gegen die (uniforme) Zufälligkeit sprechen, möglichst genau zu notieren. Ihr Resumée?</p> <p>Gehen Sie die Argumente noch einmal durch und versuchen Sie, eine Prüfstrategie zu entwerfen, mit der Sie eine Folge von Zahlen auf (uniforme) Zufälligkeit überprüfen könnten. Versuchen Sie, diese Strategie möglichst genau zu formulieren.</p> <p style="text-align: right;">(Fortsetzung)→</p>

Aufgabe 1.1	(Fortsetzung)
	<p><i>Hinweis:</i> Sie können mehrere Abbildungen in einem Fenster halten. Mit</p> $\text{par}(\text{mfrow} = \text{c}(2, 3))$ <p>wird das Grafik-System so eingestellt, dass jeweils sechs Abbildungen zeilenweise als 2×3-Matrix angeordnet (2 Zeilen, 3 Spalten) gezeigt werden.</p> <p>Die Funktion par ist die zentrale Funktion, mit der die Grafik-Ausgabe parametrisiert wird. Weitere Information erhält man mit help(par).</p>

Wir lüften gleich das Geheimnis¹: die Zahlen sind nicht zufällig, sondern ganz deterministisch. Genauer: im Hintergrund von `runif()` wird eine deterministische Folge z_i generiert. Verschiedene Algorithmen stehen zur Verfügung. Informationen dazu erhält man mit `help(.Random.seed)`. Im einfachsten Fall, für lineare Kongruenzgeneratoren, werden aufeinanderfolgende Werte z_i, z_{i+1} sogar nur mit einer linearen Funktion generiert. Damit die Werte im kontrollierten Bereich bleiben, wird modulo einer oberen Grenze gerechnet, also

$$z_{i+1} = a z_i + b \quad \text{mod } M.$$

Die resultierenden Werte, die uns übergeben werden, sind umskaliert auf

$$\frac{z_i}{M} \cdot (\text{max} - \text{min}) + \text{min}.$$

Die dadurch definierte Folge kann regelmäßig sein und schnell zu periodischer Wiederholung führen. Bei geeigneter Wahl der Parameter, wie beim Beispiel in der Fußnote, kann sie jedoch zu einer sehr langen Periode (in der Größenordnung von M) führen und scheinbar zufällig sein. Die Zahlenfolge ist jedoch keine unabhängige Zufallsfolge, und die Verteilung ist auch nicht uniform auf (min, max) .

Selbst wenn man das Geheimnis kennt, ist es nur mit viel weiterem Wissen möglich nachzuweisen, dass die erzeugte Folge nicht den Gesetzen folgt, die für eine unabhängige Folge von identisch uniform verteilten Zufallszahlen gelten.

Zahlenfolgen, die den Anspruch erheben, sich wie zufällige Zahlen zu verhalten, nennen wir **Pseudo-Zufallszahlen**, wenn es wichtig ist, auf den Unterschied hinzuweisen. Wir benutzen diese Pseudo-Zufallszahlen, um uns geeignete Test-Datensätze zu generieren. Wir können damit untersuchen, wie sich statistische Verfahren unter nahezu bekannten Bedingungen verhalten. Dabei benutzen wir Pseudo-Zufallszahlen, als ob wir Zufallszahlen hätten.

Pseudo-Zufallszahlen sollten wir zum anderen als Herausforderung nehmen: Sind wir in der Lage, sie als nicht unabhängige Zufallszahlen zu erkennen? Wenn wir einen

¹... nur teilweise. Die benutzten Zufallsgeneratoren in R sind konfigurierbar und können wesentlich komplexer sein, als hier vorgestellt. Für unsere Diskussion reicht jedoch hier die Familie der linearen Kongruenzgeneratoren. Sie können deren Verhalten in anderen Programmiersystemen nachvollziehen. Die übliche Referenz ist dabei der "minimal standard generator" mit $x_{i+1} = (x_i \times 7^5) \text{ mod } 2^{31} - 1$.

Unterschied erkennen, werden wir versuchen, den Pseudo-Zufallszahlengenerator gegen einen besseren auszutauschen. Aber zunächst geht die Herausforderung an uns. Sind wir überhaupt in der Lage, z.B. eine mit einem linearen Generator erzeugte deterministische Folge als nicht zufällig zu erkennen? Falls nicht: welche intellektuellen Konsequenzen ziehen wir daraus?

1.2.2. Muster. Außer Pseudo-Zufallszahlen gibt es in R eine ganze Reihe von Möglichkeiten, regelmäßige Sequenzen zu generieren. Die in anderen Sprachen notwendigen Schleifen werden damit weitgehend vermieden. Hier eine erste Übersicht:

R Sequenzen	
:	Erzeugt Sequenz von <i>Anfang</i> bis höchstens <i>Ende</i> . <i>Aufruf:</i> <code>Anfang:Ende</code> <i>Beispiele:</i> <code>1:10</code> <code>10.1:1.2</code>
<code>c()</code>	“combine”. Kombiniert Argumente zu einem neuen Vektor. <i>Aufruf:</i> <code>c(..., recursive = FALSE)</code> <i>Beispiele:</i> <code>c(1, 2, 3)</code> <code>c(x, y)</code> Bezeichnen die Argumente zusammengesetzte Datentypen, so arbeitet die Funktion rekursiv absteigend in die Daten hinab, wenn sie mit <code>recursive = TRUE</code> aufgerufen wird.
<code>seq()</code>	Erzeugt allgemeine Sequenzen. <i>Aufruf:</i> Siehe <code>help(seq)</code>
<code>rep()</code>	Wiederholt Argument. <i>Aufruf:</i> <code>rep(x, times, ...)</code> <i>Beispiele:</i> <code>rep(x, 3)</code> <code>rep(1:3, c(2, 3, 1))</code>

Dabei steht “...” für eine variable Liste von Argumenten. Wir werden diese Notation noch häufiger benutzen.

Aufgabe 1.2	
	Generieren Sie mit $\text{plot}(\sin(1:100))$ einen Plot mit einer diskretisierten Sinusfunktion. (Falls Sie die Sinusfunktion nicht sofort erkennen, benutzen Sie <code>plot(sin(1:100), type = "l")</code> , um die Punkte zu verbinden. Benutzen Sie Ihre Strategie aus Aufgabe 1.1. Können Sie damit die Sinusfunktion als nicht zufällig erkennen?

Die Zahlenreihe eines Datensatzes, wie z.B. die Ausgabe eines Zufallszahlengenerators hilft selten, zugrunde liegende Strukturen zu erkennen. Nur wenig helfen einfache, unspezifische grafische Darstellungen wie der Serienplot. Selbst bei klaren Mustern sind diese Informationen selten aussagekräftig. Zielgerichtete Darstellungen sind nötig, um Verteilungseigenschaften zu untersuchen.

1.3. Fallstudie: Verteilungsdiagnostik

Wir brauchen genauere Strategien, um Strukturen zu erkennen oder deren Verletzung festzustellen. Wie diese Strategien aussehen können, skizzieren wir am Beispiel der Zufallszahlen. Wir konzentrieren uns hier auf die Verteilungseigenschaft. Angenommen, die Folge besteht aus unabhängigen Zufallszahlen mit einer gemeinsamen Verteilung. Wie überprüfen wir, ob dies die uniforme Verteilung ist? Wir ignorieren die mögliche Umskalierung auf (\min, \max) - dies ist ein technisches Detail, das die Fragestellung nicht wesentlich tangiert. Wir betrachten $\min = 0; \max = 1$.

Aus Realisierungen von Zufallsvariablen können Verteilungen nicht direkt abgelesen werden. Dies ist unser kritisches Problem. Wir brauchen Kennzeichnungen der Verteilungen, die wir empirisch überprüfen können. Wir können zwar Beobachtungen als Maße betrachten: Für n Beobachtungen X_1, \dots, x_n können wir formal die empirische Verteilung P_n definieren als das Maß $P_n = \sum (1/n)\delta_{x_i}$, wobei δ_{X_i} das Dirac-Maß an der Stelle X_i ist. Also

$$P_n(A) = \#\{i : X_i \in A\}/n.$$

Aber leider ist das empirische Maß P_n einer Beobachtungsreihe von unabhängigen Beobachtungen mit gemeinsamem Maß P im allgemeinen sehr von P verschieden. Einige Eigenschaften gehen unwiederbringlich verloren. Dazu gehören infinitesimale Eigenschaften: so ist z.B. P_n immer auf endlich viele Punkte konzentriert. Wir brauchen Konstrukte, die anhand von Realisierungen von Zufallsvariablen bestimmbar und mit den entsprechenden Konstrukten von theoretischen Verteilungen vergleichbar sind. Eine Strategie ist es, sich auf (empirisch handhabbare) Testmengen zu beschränken.

BEISPIEL 1.1. Verteilungsfunktion

Anstelle der Verteilung P betrachten wir ihre Verteilungsfunktion $F = F_P$ mit

$$F(x) = P(X \leq x).$$

Für eine empirische Verteilung P_n von n Beobachtungen X_1, \dots, X_n ist entsprechend die empirische Verteilungsfunktion

$$F_n(x) = \#\{i : X_i \leq x\}/n.$$

BEISPIEL 1.2. Histogramm

Wir wählen disjunkte Testmengen $A_j, j = 1, \dots, J$, die den Wertebereich von X überdecken. Für die uniforme Verteilung auf $(0, 1)$ können wir z.B. die Intervalle

$$A_j = \left(\frac{j-1}{J}, \frac{j}{J} \right]$$

als Testmengen wählen.

Anstelle der Verteilung P betrachten wir den Vektor $(P(A_j))_{j=1,\dots,J}$ bzw. den empirischen Vektor $(P_n(A_j))_{j=1,\dots,J}$.

Wir diskutieren diese Beispiele ausführlicher. Einige allgemeine Lehren können wir daraus ziehen. Wir machen mehrere Durchgänge, um von einem naiven Zugang zu einem entwickelten statistischen Ansatz zu kommen.

An dieser Stelle sei schon darauf hingewiesen, dass Histogramme kritisch von der Wahl der Testmengen abhängen. Insbesondere wenn Diskretisierungen in den Daten unglücklich mit der Wahl der Testmengen zusammentreffen, kann es zu sehr irreführenden Ergebnissen kommen. Eine Alternative zu Histogrammen ist es, die Daten zu glätten.

BEISPIEL 1.3. Glättung Wir ersetzen jeden Datenpunkt durch eine (lokale) Verteilung, d.h. wir verschmieren die Datenpunkte etwas. Wir benutzen dazu Gewichtsfunktionen. Diese Gewichtsfunktionen werden **Kerne** genannt und mit K bezeichnet. Wenn die Kerne integrierbar sind, normieren wir sie konventionell so, dass $\int K(x)dx = 1$. Einige übliche Kerne sind in Tabelle 1.9 aufgelistet und in Abb. 1.1 gezeigt. Wenn sie einen kompakten Träger haben, so ist als Träger das Intervall $[-1, 1]$ gewählt (Die R-Konvention ist es, die Kerne so zu standardisieren, dass sie die Standardabweichung 1 haben).

Kern	$K(x)$
Uniform	$1/2$
Dreieck	$1 - x $
Epanechnikov (quadratisch)	$3/4(1 - x^2)$
Biweight	$15/16(1 - x^2)^2$
Triweight	$35/32(1 - x^2)^3$
Gauß	$(2\pi)^{-1/2} \exp(-x^2/2)$

TABELLE 1.9. Einige übliche Kerne

Durch Verschiebung und Umskalierung definiert jeder Kern eine ganze Familie

$$\frac{1}{h}K\left(\frac{x - x_0}{h}\right).$$

Der Skalenfaktor h wird **Bandbreite** genannt. Der mit h skalierte Kern wird mit K_h bezeichnet:

$$K_h(x) = \frac{1}{h}K\left(\frac{x}{h}\right).$$

Die Funktion

$$x \mapsto \frac{1}{n} \sum_i K_h(x - X_i)$$

ergibt anstelle des Histogramms ein geglättetes Bild.

Näheres dazu findet man unter dem Stichwort **smoothing** in der Literatur.

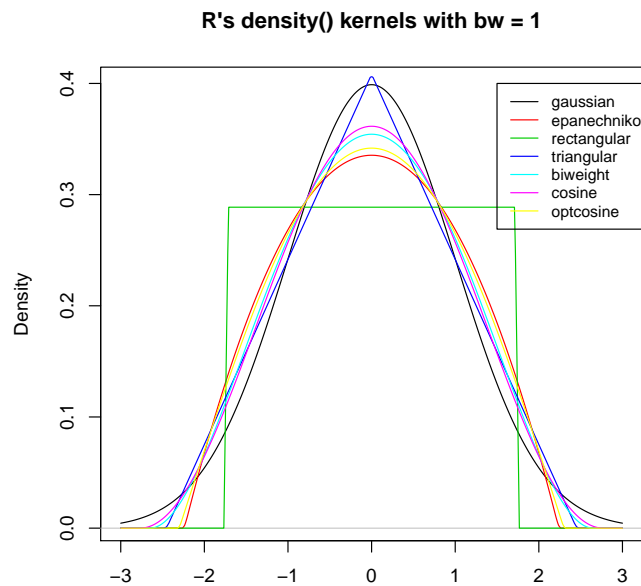


ABBILDUNG 1.1. Kerne in R

1.3.1. Erster Durchgang zu Beispiel 1.1: Verteilungsfunktion. Um zu prüfen, ob eine Zufallsfolge zu einer Verteilung mit Verteilungsfunktion F passt, vergleiche man F mit F_n . Im Fall der uniformen Verteilung auf $(0, 1)$ ist $F(x) = F_{unif}(x) = x$ für $0 \leq x \leq 1$. Der ganz naive Zugang berechnet die Funktionen F_n und F . Eine erste Überlegung sagt: F_n ist eine stückweise konstante Funktion mit Sprungstellen an den Beobachtungspunkten. Wir bekommen also ein vollständiges Bild von F_n , wenn wir F_n an den Beobachtungspunkten $X_i, i = 1..n$ auswerten. Ist $X_{(i)}$ die i . Ordnungsstatistik, so ist - bis auf Bindungen - $F_n(X_{(i)}) = i/n$. Wir vergleichen $F_n(X_{(i)})$ mit dem "Sollwert" $F(X_{(i)}) = X_{(i)}$. Eine R-Implementierung, mit Hilfsvariablen notiert:

```
n <- 100
x <- runif(n)
xsort <- sort(x)
i <- (1:n)
y <- i/n
plot(xsort, y)
```

Eine zusätzliche Gerade für die "Sollwerte" kann mit

```
abline(0, 1)
```

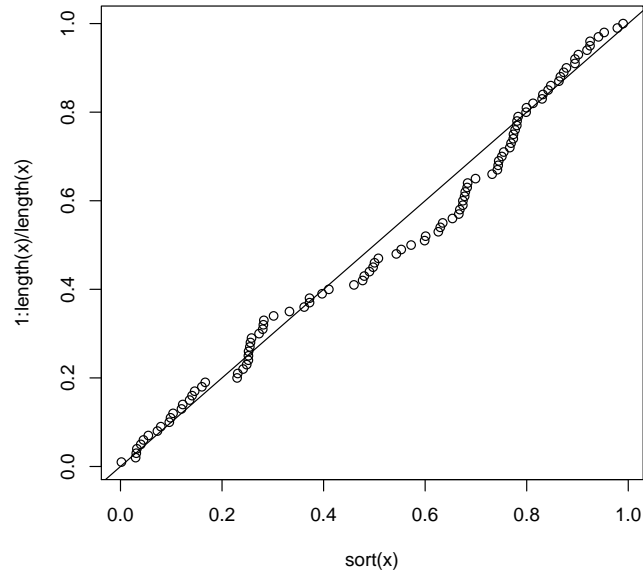
eingezeichnet werden.

Eine kompaktere Implementierung mit der Funktion `length()`:

Beispiel 1.2:

Eingabe

```
x <- runif(100)
plot(sort(x), 1:length(x)/length(x))
abline(0, 1)
```



R Funktionen	
<code>sort()</code>	Sortiert Vektor <i>Beispiel:</i> <code>sort(runif(100))</code>
<code>length()</code>	Länge eines Vektors <i>Beispiel:</i> <code>length(x)</code>
<code>abline()</code>	Fügt Linie in Plot hinzu <i>Beispiel:</i> <code>abline(a = 0, b = 2)</code>

Die Funktion `plot()` fügt defaultmäßig Beschriftungen hinzu. Damit die Grafik für sich aussagekräftig ist, wollen wir diese durch genauere Beschriftungen ersetzen. Dazu ersetzen wir die Default-Parameter von `plot()` durch unsere eigenen. Der Parameter `main` kontrolliert die Hauptüberschrift (Default: leer). Wir können diese zum Beispiel ersetzen wie in

```
plot(sort(x), (1:length(x))/length(x),
      main = "Empirische Verteilungsfunktion\n (X uniform)").
```

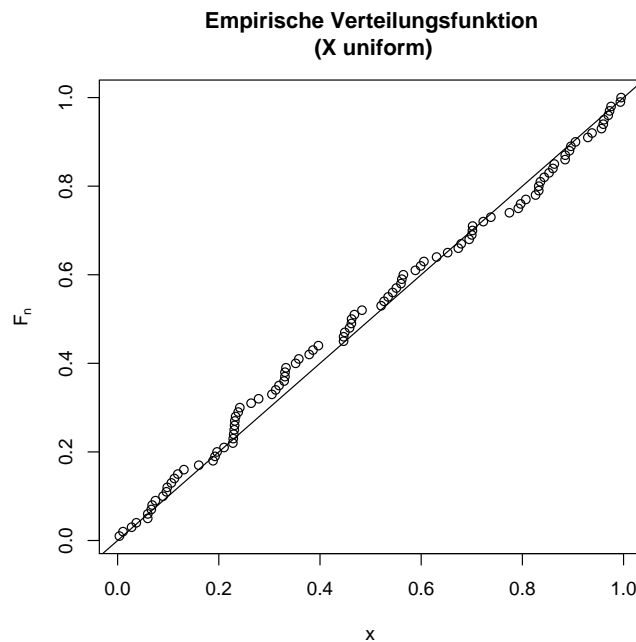

Mit `xlab` und `ylab` wird die Beschriftung der Achsen gesteuert. Über diese und weitere Parameter kann man Information mit `help(plot)` abfragen, werden dann aber weiter an `help(title)` verwiesen.

Die vertikale Achse gibt noch eine Herausforderung: mit `ylab = Fn(x)` als Parameter würden wir eine Beschriftung mit $F_n(x)$ erhalten. Die übliche Bezeichnung setzt aber den Stichprobenumfang als Index, also $F_n(x)$. Hier hilft eine versteckte Eigenschaft der Beschriftungsfunktionen: Wird als Parameter eine Zeichenkette übergeben, so wird sie ohne Umwandlung angezeigt. Wird als Parameter ein R-Ausdruck übergeben, so wird versucht, die mathematisch übliche Darstellung zu geben. Details findet man mit `help(plotmath)` und Beispiele mit `demo(plotmath)`. Die Umwandlung einer Zeichenkette in einen (unausgewerteten) R-Ausdruck geschieht mit `expression()`.

Beispiel 1.3:

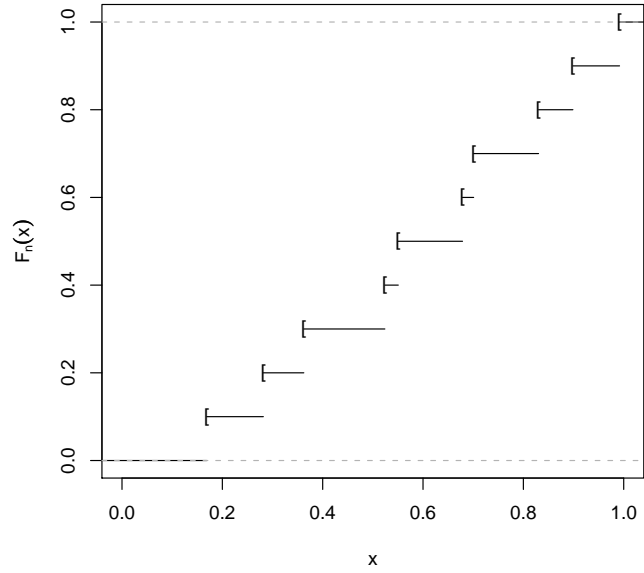
Eingabe

```
x <- runif(100)
plot(sort(x), (1:length(x))/length(x),
      xlab = "x", ylab = expression(F[n]),
      main = "Empirische Verteilungsfunktion\n (X uniform)"
)
abline(0, 1)
```



Dieses Beispiel dient hier nur zur Einführung. Es ist nicht notwendig, die Verteilungsfunktion selbst zu programmieren. In R gibt es z.B. die Klasse `ecdf` für die empirische Verteilungsfunktion. Wird die Funktion `plot()` auf ein Objekt der Klasse `ecdf` angewandt, so führt die "generische" Funktion `plot()` intern auf die spezielle

Funktion `plot.ecdf`, und diese zeichnet in der für Verteilungsfunktionen speziellen Weise. Wir können das Beispiel also abkürzen durch den Aufruf `plot(ecdf(runif(100)))`.

Aufgabe 1.3	
	<p>Ergänzen Sie den Aufruf <code>plot(ecdf(runif(10)))</code> durch weitere Parameter so, das die Ausgabe die folgende Form hat:</p> <div style="text-align: center;"> <p>Empirische Verteilungsfunktion (X uniform)</p>  </div>

Aufgabe 1.4	
	<p>Mit <code>rnorm()</code> generieren Sie gaußverteilte Zufallsvariablen. Versuchen Sie, anhand der Serienplots gaußverteilte Zufallsvariablen von uniform verteilten zu unterscheiden.</p> <p>Benutzen Sie dann die empirischen Verteilungsfunktionen. Können Sie damit gaußverteilte von uniform verteilten unterscheiden? Die Sinus-Serie von uniform verteilten? von gaußverteilten?</p> <p>Wie groß ist der benötigte Stichprobenumfang, um die Verteilungen verlässlich zu unterscheiden?</p>

1.3.2. Erster Durchgang zu Beispiel 1.2: Histogramm. Wir wählen Testmengen A_j , $j = 1, \dots, J$ im Wertebereich von X . Strategie: Um zu prüfen, ob eine Zufallsfolge zu einer Verteilung P gehört, vergleiche man den Vektor $(P(A_j))_{j=1, \dots, J}$

mit $(P_n(A_j))_{j=1,\dots,J}$. Für die uniforme Verteilung auf $(0, 1)$ können wir z.B. die Intervalle

$$A_j = \left(\frac{j-1}{J}, \frac{j}{J} \right]$$

als Testmengen wählen. Dann ist

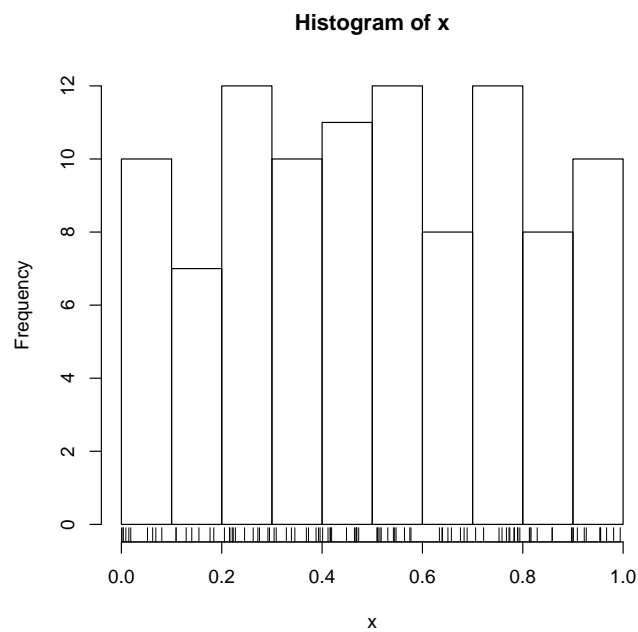
$$(P(A_j))_{j=1,\dots,J} = (1/J, \dots, 1/J)$$

der theoretische Vergleichsvektor zum Vektor der beobachteten relative Häufigkeiten $\frac{\#\{i: X_i \in A_j\}}{n}$ $j = 1, \dots, J$. Vorläufige Implementierung: wir benutzen hier gleich eine vorgefertigte Funktion, die Histogramme zeichnet. Als Seiteneffekt liefert sie uns die gewünschten Werte. Mit der Funktion `rug()` können wir die Originaldaten zusätzlich einblenden.

Beispiel 1.4:

Eingabe

```
x <- runif(100)
hist(x)
rug(x)
```



Zum Vergleich können wir einen Dichteschätzer überlagern. Da `density()` im Gegensatz zu `hist()` das Resultat nicht zeichnet, sondern ausdrückt, müssen wir die Grafik explizit anfordern. Damit die Skalen vergleichbar sind, fordern wir für das Histogramm mit dem Parameter `probability = TRUE` eine Wahrscheinlichkeitsdarstellung an.

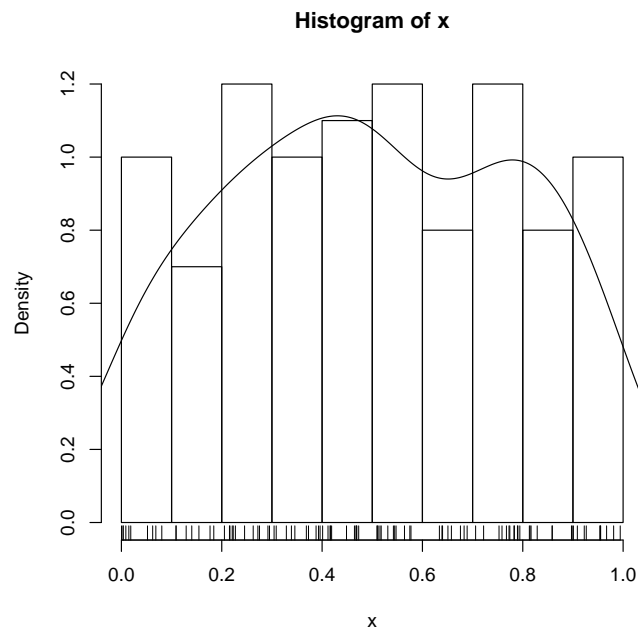
Beispiel 1.5:

```

hist(x, probability = TRUE)
rug(x)
lines(density(x))

```

Eingabe



Histogramm und Kern-Dichteschätzer haben jeweils ihre spezifischen Vorteile und Probleme. Histogramme leiden unter ihrer Diskretisierung, die mit einer Diskretisierung in den Daten unglücklich zusammen treffen kann. Kern-Dichteschätzer “verschmieren” die Daten, und können dadurch insbesondere am Rand des Datenbereichs zu unangemessenen Rand-Effekten führen.

Zurück zum Histogramm: Benutzen wir eine Zuweisung

```
xhist <- hist(x),
```

so wird die interne Information des Histogramm unter **xhist** gespeichert und kann mit

```
xhist
```

abgerufen werden. Sie ergibt z.B.

Beispiel 1.6:

```

----- Eingabe -----
x <- runif(100)
xhist <- hist(x)
xhist
----- Ausgabe -----
$breaks
[1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0

$count
[1] 8 13 8 12 12 7 6 11 10 13

$intensities
[1] 0.7999998 1.3000000 0.8000000 1.2000000 1.2000000 0.7000000
[7] 0.6000000 1.1000000 1.0000000 1.3000000

$density
[1] 0.7999998 1.3000000 0.8000000 1.2000000 1.2000000 0.7000000
[7] 0.6000000 1.1000000 1.0000000 1.3000000

$mids
[1] 0.05 0.15 0.25 0.35 0.45 0.55 0.65 0.75 0.85 0.95

$xname
[1] "x"

$equidist
[1] TRUE

attr(,"class")
[1] "histogram"

```

Counts gibt dabei die Besetzungszahlen der Histogrammzellen, d.h. die von uns gesuchte Anzahl. Die in *xhist* gespeicherte interne Information des Histogramms besteht aus fünf wesentlichen Komponenten - hier jeweils Vektoren. Diese Komponenten von *xhist* haben Namen und können mit Hilfe dieser Namen angesprochen werden. So gibt z.B.

```
xhist$count
```

den Vektor der Besetzungszahlen.

R <i>Datenstruktu-</i> <i>ren</i>	
--------------------------------------	--

(Fortsetzung)→

R Datenstrukturen (Fortsetzung)	
Vektoren	Komponenten eines Vektors werden über ihren Index angesprochen. Alle Elemente eines Vektors haben denselben Typ. <i>Beispiele:</i> <code>x</code> <code>x[10]</code>
Listen	Listen sind zusammengesetzte Datenstrukturen. Die Komponenten einer Liste haben Namen, über die sie angesprochen werden können. Teilkomponenten einer Liste können von unterschiedlichem Typ sein. <i>Beispiele:</i> <code>xhist</code> <code>xhist\$counts</code>

Weitere zusammengesetzte Datenstrukturen sind im Anhang (A.8) beschrieben.

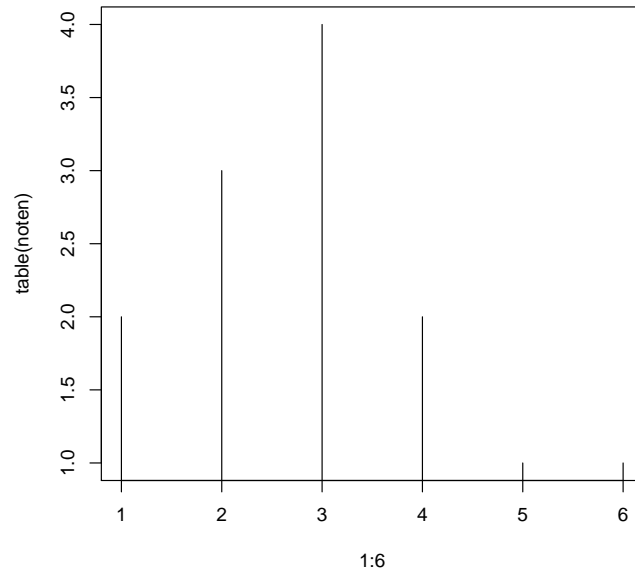
Die Wahl der Histogrammgrenzen erfolgt automatisch. Für die genaue Behandlung der Intervallgrenzen gibt es unterschiedliche Konventionen, deren Wahl durch Parameter von `hist()` gesteuert werden kann. Um unsere Testmengen zu benutzen, müssen wir die Aufrufstruktur von `hist()` erfragen.

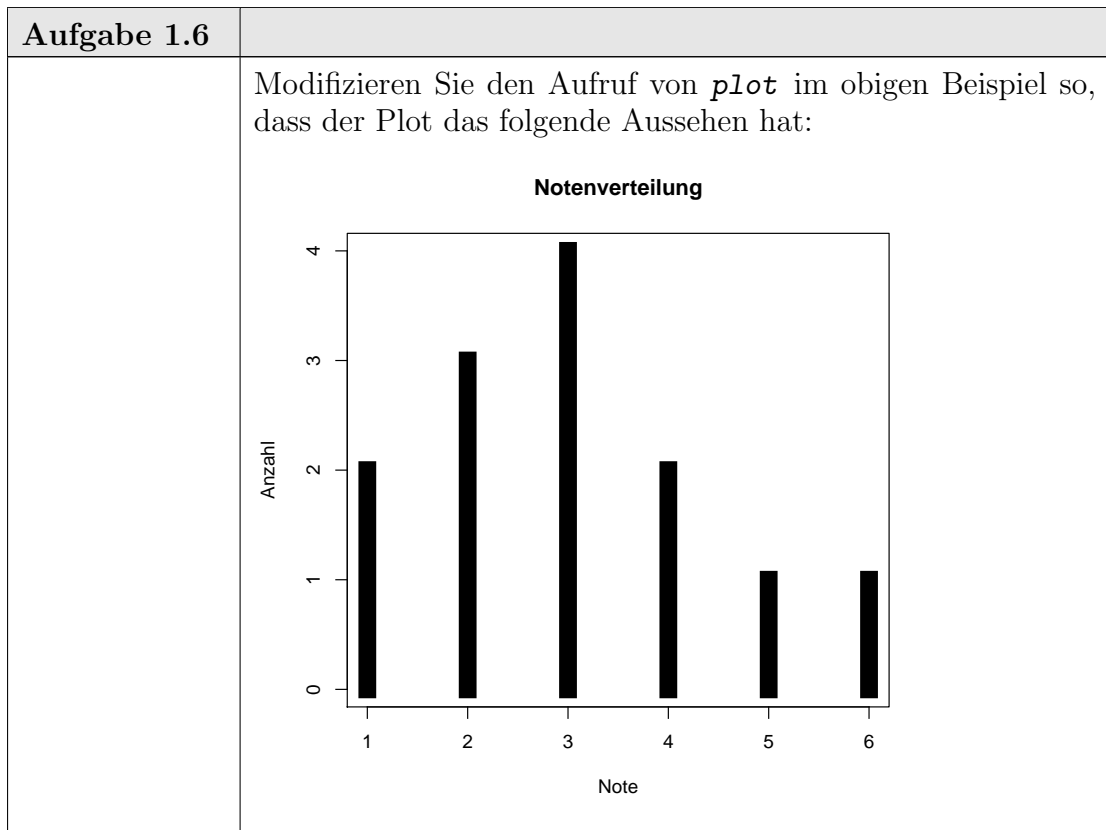
Aufgabe 1.5	
	<p>Generieren Sie zu <code>runif(100)</code> Histogramme mit 5, 10, 20, 50 gleichgroßen Zellen und ziehen Sie wiederholt Stichproben. Entsprechen die Bilder dem, was Sie von unabhängig uniform verteilten Zufallsvariablen erwarten? Versuchen Sie, ihre Beobachtungen möglichst genau zu notieren.</p> <p>Wiederholen Sie das Experiment mit zwei Zellen (0, 0.5], (0.5, 1).</p> <pre>hist(runif(100), breaks = c(0, 0.5, 1))</pre> <p>Wiederholen Sie das Experiment analog mit <code>rnorm(100)</code> und vergleichen Sie die Resultate von <code>runif(100)</code> und <code>rnorm(100)</code>.</p>

1.3.2.1. *Balkendiagramme.* Als Hinweis: wenn die Daten nicht quantitativ sind, sondern kategorial (durch Kategorien-Label bezeichnet, wie z.B. “sehr gut, gut, befriedigend, ...”, oder durch Kennziffern bezeichnet, wie z.B. “1, 2, 3, ...”), so ist ein Balkendiagramm eher geeignet. Einfache Balkendiagramme werden von `plot()` selbst durch den Parameter `type = h` unterstützt. Dazu müssen aus den Rohdaten die Häufigkeiten der einzelnen Stufen bestimmt werden. Dies kann mit der Funktion `table()` geschehen.

Beispiel 1.7:

```
noten <- c(2, 1, 3, 4, 2, 2, 3, 5, 1, 3, 4, 3, 6)
plot(1:6, table(noten), type = 'h')
```





1.3.3. Zweiter Durchgang zu Beispiel 1.1: Verteilungsfunktion. Wir machen jetzt einen Schritt von einem naiven Ansatz zu einer statistischen Betrachtung. Naiv haben wir für unabhängig identisch verteilte Variable (X_1, \dots, X_n) mit Verteilungsfunktion F angenommen, dass $i/n = F_n(X_{(i)}) \approx F(X_{(i)})$ und dies zur Überprüfung der Verteilungsannahme benutzt. Speziell für uniform $(0, 1)$ verteilte Variable ist diese naive Annahme: $i/n \approx X_{(i)} = F(X_{(i)})$.

Statistisch gesehen ist $X_{(i)}$ eine Zufallsvariable. Damit ist auch $F(X_{(i)})$ eine Zufallsvariable mit Werten in $[0, 1]$, und wir können die Verteilung dieser Zufallsvariablen untersuchen.

THEOREM 1.4. *Sind (X_1, \dots, X_n) unabhängig identisch verteilte Zufallsvariablen mit stetiger Verteilungsfunktion F , so ist $F(X_{(i)})$ verteilt nach der Beta-Verteilung $\beta(i, n - i + 1)$.*

BEWEIS. \rightarrow Wahrscheinlichkeitstheorie. Hinweis: Benutze

$$X_{(i)} \leq x_\alpha \Leftrightarrow (\#j : X_j \leq x_\alpha) \geq i.$$

Für stetige Verteilungen ist $(\#j : X_j \leq x_\alpha)$ binomialverteilt mit Parametern (n, α) . \square

KOROLLAR 1.5.

$$E(F(X_{(i)})) = i/(n + 1).$$

Aufgabe 1.7	
	Mit <code>help(rbeta)</code> erhalten Sie Informationen über die Funktionen, die für die Beta-Verteilungen bereitstehen. Plotten Sie die entsprechenden Dichten der Beta-Verteilungen für $n = 10, 50, 100$ und $i = n/4, n/2, 3n/4$. Benutzen Sie zum plotten die Funktion <code>curve()</code> . Zum Aufruf, siehe <code>help(curve)</code> .

Wir können also im statistischen Mittel für uniform auf $(0, 1)$ verteilte Variable nicht erwarten, dass $X_{(i)} \approx i/n$, sondern im Mittel erhalten wir $i/(n+1)$. Die “richtige Sollwertgerade” sollte also mit `abline(a = 0, b = n/n+1)` gezeichnet werden.

Aufgabe 1.8	
	Zeichnen Sie die Verteilungsfunktion mit der korrigierten Geraden.
*	Für die grafische Darstellung wird jeweils nur ein Plot benutzt. Ist der Erwartungswert von $X_{(i)}$ hier der richtige Vergleichsmaßstab? Gibt es Alternativen? Falls Sie Alternativen sehen: implementieren Sie diese.

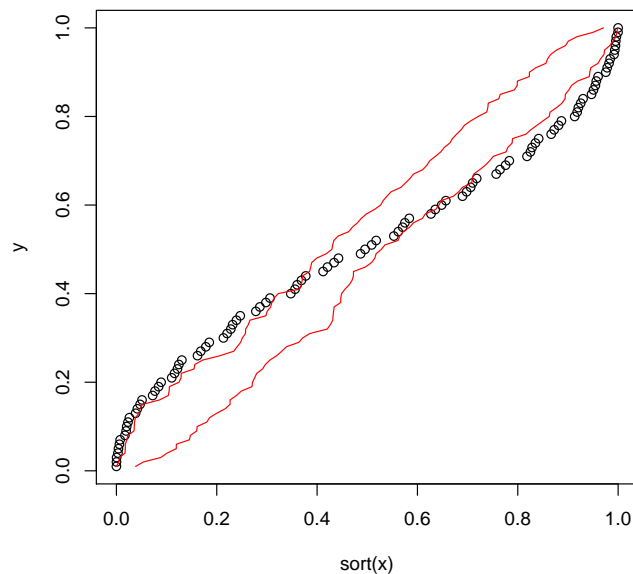
1.3.3.1. *Monte Carlo Konfidenzbänder.* Mit einer Simulation können wir uns auch ein Bild von der typischen Fluktuation verschaffen. Wir benutzen Zufallszahlen, um eine (kleine) Anzahl von Stichproben bekannter Verteilung zu generieren, und vergleichen die in Frage stehende Stichprobe mit den Simulationen. Dazu bilden wir für die Simulationen die Einhüllende, und prüfen, ob die Stichprobe innerhalb dieses Bereichs liegt. Ist x der in Frage stehende Vektor mit Länge n , so benutzen wir z.B. die folgende Programmidee:

Beispiel 1.8:

```

x <- (sin(1:100)+1)/2      Eingabe # demo example only
y <- (1:length(x))/length(x)
plot(sort(x), y)
nrsamples <- 19          # nr of simulations
samples <- matrix(data = runif(length(x)* nrsamples),
nrow = length(x), ncol = nrsamples)
samples <- apply(samples, 2, sort)
envelope <- t(apply(samples, 1, range))
lines(envelope[, 1], y, col = "red")
lines(envelope[, 2], y, col = "red")

```



Dieses Beispiel ist [VR02], entnommen, einer reichen Quelle von R-Beispielen. Für die Programmierung wird hier eine für R typische Strategie erkennbar. R ist eine interpretierte vektor-orientierte Sprache. Einzelne Interpretationsschritte sind zeitintensiv. Deshalb sind Operationen mit weniger, dafür komplexeren Schritten effektiver als Operationen aus mehreren elementaren Schritten.

- Operationen auf Vektorebene sind effektiver als Ketten einzelner elementare Operationen.
- Iterationen und Schleifen werden vermieden zugunsten strukturierter Vektor-Operationen.

Aufgabe 1.9	
	Benutzen Sie die <code>help()</code> -Funktion und kommentieren Sie das obige Beispiel Schritt für Schritt. Notieren Sie insbesondere die neu hinzugekommenen Funktionen.

R Iteratoren	
<code>apply()</code>	wendet eine Funktion auf die Zeilen oder Spalten einer Matrix an. <i>Beispiel:</i> <code>samples <- apply(samples, 2, sort)</code> sortiert spaltenweise.
<code>outer()</code>	erzeugt eine Matrix mit allen Paar-Kombinationen aus zwei Vektoren, und wendet eine Funktion auf jedes Paar an.

Wenn die Kurve für unsere Stichprobe die durch die Simulation gewonnenen Grenzen überschreitet, so widerspricht das der Hypothese, dass der Stichprobe und der Simulation das selbe Modell zugrunde liegt. Das hier skizzierte Verfahren heißt **Monte-Carlo-Test**. Die Idee dahinter ist von sehr allgemeiner Bedeutung.

Aufgabe 1.10	
*	Wieso 19? <i>Hinweis:</i> versuchen Sie, das Problem zunächst abstrakt und vereinfacht zu betrachten: sei T eine messbare Funktion und $X_0, X_1, \dots, X_{nrsamples}$ unabhängige Stichproben mit einer gemeinsamen stetigen Verteilungsfunktion. Berechnen Sie $P(T(X_0) > T(X_i))$ für alle $i > 0$. Formulieren Sie dann das obige Beispiel abstrakt. Spezialisieren Sie dann für $nrsamples = 19$.

Aufgabe 1.11	
*	Schätzen Sie die Überdeckungswahrscheinlichkeit des Monte-Carlo-Bands, in dem Sie wie folgt vorgehen: Generieren Sie zunächst analog zum obigen Beispiel ein Band. (Wie können Sie das Band zeichnen, ohne zuvor für eine spezielle Stichprobe einen Plot zu machen?) Ziehen Sie für eine zu wählende Anzahl sim (100? 1000? 999?) jeweils eine Stichprobe von uniform verteilten Zufallszahlen vom Stichprobenumfang 100. Zählen Sie aus, wie oft die empirische Verteilungsfunktion der Stichprobe innerhalb des Bands verläuft. Schätzen Sie hieraus die Überdeckungswahrscheinlichkeit. (Fortsetzung)→

Aufgabe 1.11	(Fortsetzung)
	Hinweis: <code>any()</code> kann benutzt werden, um für einen ganzen Vektor einen Vergleich zu machen.

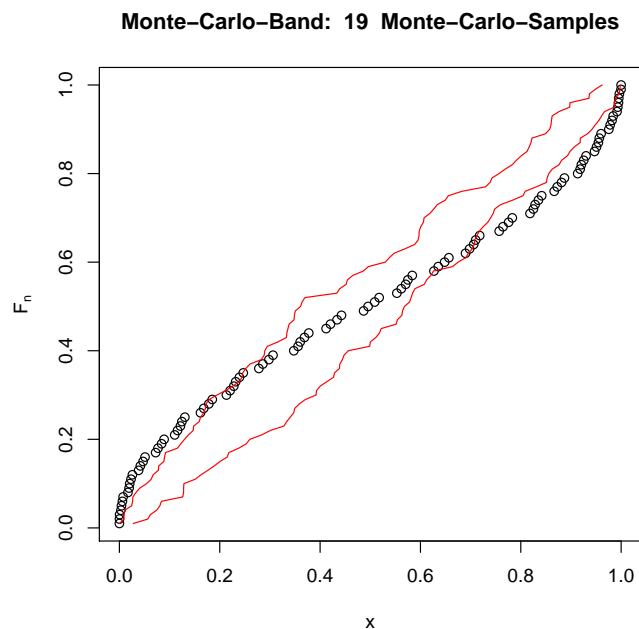
Wir wollen auch hier die Ausgabe noch überarbeiten, so dass der Plot genügend Information enthält. Bei der Beschriftung können wir zunächst analog zu Abschnitt 1.3.1 vorgehen. Die Anzahl `nrsamples` bedarf des Nachdenkens. Wenn wir nur eine feste Anzahl (z.B. 19) betrachten wollen, können wir diese wie gewohnt in die Beschriftung aufnehmen. Wenn das Programmfragment jedoch allgemeiner nutzbar sein soll, so müssten wir die jeweils gewählte Anzahl von Simulationen angeben. Dies von Hand zu tun ist eine Fehlerquelle, die vermieden werden kann. Die Funktion `bquote()` ermöglicht es, den jeweils aktuellen Wert zu erfragen oder im jeweiligen Kontext zu berechnen. Damit die Anzahl von Simulationen in die Überschrift übernommen werden kann, vertauschen wir die Anweisungen so dass die Anzahl der Simulationen vor dem Aufruf von `plot` festgelegt ist.

Beispiel 1.9:

```

x <- (sin(1:100)+1)/2                               Eingabe                               # demo example only
y <- (1:length(x))/length(x)
nrsamples <- 19                                     # nr of simulations
plot(sort(x), y,
      main = paste("Monte-Carlo-Band: ", bquote(.(nrsamples)), " Monte-Carlo-Samples"),
      xlab = 'x', ylab = expression(F[n]))
samples <- matrix(data = runif(length(x) * nrsamples),
                  nrow = length(x), ncol = nrsamples)
samples <- apply(samples, 2, sort)
envelope <- t(apply(samples, 1, range))
lines(envelope[, 1], y, col = "red")
lines(envelope[, 2], y, col = "red")

```



Für die Simulationen werden jeweils neue Monte-Carlo-Stichproben gezogen. Deshalb erhalten Sie bei jedem Aufruf unterschiedliche Monte-Carlo-Konfidenzbänder und die Bänder hier sind von denen im vorherigen Beispiel verschieden.

Für die praktische Arbeit kann es notwendig sein, die Verteilungsdiagnostik auf ein einfaches Entscheidungsproblem zu reduzieren, etwa um anhand von Tabellen oder Kontrollkarten zu entscheiden, ob eine Verteilung in einem hypothetischen Bereich liegt, oder eine Kenngröße anzugeben, die die Abweichung von einem gegebenen Modell charakterisiert. Wenn wir auf Tabellen oder einfache Zahlen zurückgreifen wollen, müssen wir uns weiter einschränken. Wir müssen die Information, die in den Funktionen (F_n, F) steckt, weiter reduzieren, wenn wir die Unterschiede numerisch zusammenfassen wollen. Eine Zusammenfassung ist etwa

$$\sup_x |F_n - F|(x).$$

Wenn wir diese Zusammenfassung als Kriterium benützen wollen, stehen wir wieder vor der Aufgabe, ihre Verteilung zu untersuchen.

THEOREM 1.6. (Kolmogorov, Smirnov) Für stetige Verteilungsfunktionen F ist die Verteilung von

$$\sup_x |F_n - F|(x)$$

unabhängig von F (jedoch abhängig von n).

BEWEIS. → Wahrscheinlichkeitstheorie. Z.B. [Gänßler & Stute, Lemma 3.3.8]. \square

THEOREM 1.7. (Kolmogorov): Für stetige Verteilungsfunktionen F und $n \rightarrow \infty$ hat

$$\sqrt{n} \sup |F_n - F|$$

asymptotisch die Verteilungsfunktion

$$F_{\text{Kolmogorov-Smirnov}}(y) = \sum_{m \in \mathbb{Z}} (-1)^m e^{-2m^2 y^2} \quad \text{für } y > 0.$$

BEWEIS. → Wahrscheinlichkeitstheorie. Z.B. [Gänßler & Stute, Formel (3.3.11)]. \square

Für die praktische Arbeit bedeutet dies: Für stetige Verteilungsfunktionen können wir eine Entscheidungsstrategie formulieren: wir entscheiden, dass die Beobachtung (X_1, \dots, X_n) nicht mit der Hypothese von unabhängig, identisch nach F verteilten Zufallsvariablen vereinbar ist, falls $\sup |F_n - F|$ zu groß ist:

$$\sup |F_n - F| > F_{\text{krit}} / \sqrt{n},$$

wobei F_{krit} aus der (von F unabhängigen) Verteilungsfunktion der Kolmogorov-Smirnov-Statistik zum Stichprobenumfang n entnommen wird. Wählen wir speziell das obere α -Quantil $F_{\text{krit}} = F_{\text{Kolmogorov-Smirnov}, 1-\alpha}$, so wissen wir, dass bei Zutreffen der Hypothese der Wert F_{krit} oder ein höherer Wert höchstens mit Wahrscheinlichkeit α erreicht wird. Damit können wir unsere Irrtumswahrscheinlichkeit für eine ungerechtfertigte Ablehnung der Hypothese kontrollieren.

Asymptotisch, für große n , können wir anstelle der Verteilungsfunktion die Kolmogorov-Approximation benutzen. Wenn die Modellverteilung F nicht stetig ist, sind weitere Überlegungen nötig.

Wir wollen uns hier auf die Programmierung konzentrieren und gehen nicht in die Details des Kolmogorov-Smirnov-Tests. Mit elementaren Mitteln können wir die Teststatistik $\sup_x |F_n - F|(x)$ für die uniforme Verteilung programmieren. Aus Monotoniegründen ist

$$\sup_x |F_n - F|(x) = \max_{X(i)} |F_n - F|X(i)$$

und für die uniforme Verteilung ist

$$\max_{X(i)} |F_n - F|X(i) = \max_i |i/n - X(i)|.$$

Damit gibt in R-Schreibweise der Ausdruck

$$\max(\text{abs}((1:\text{length}(x)) / \text{length}(x)) - \text{sort}(x)))$$

die für uns die gewünschte Statistik, wenn \mathbf{x} unser Datenvektor ist.

Diese Statistik (und viele weitere allgemein benutzte Statistiken) sind in der Regel schon programmiert, ebenso wie die zugehörigen Verteilungsfunktionen.²

Aufgabe 1.12	
	<p>Mit</p> <pre style="text-align: center;"><code>help(ks.test)</code></pre> <p>erhalten Sie die Information, wie die Funktion <code>ks.test</code> angewandt wird.</p> <p>Welche Resultate erwarten Sie, wenn Sie die folgenden Vektoren auf uniforme Verteilung testen:</p> <pre style="text-align: center;"><code>1:100</code> <code>runif(100)</code> <code>sin(1:100)</code> <code>rnorm(100)?</code></pre> <p>Führen Sie diese Tests durch (skalieren Sie dabei die Werte so, dass sie im Intervall $[0, 1]$ liegen, oder benutzen Sie eine uniforme Verteilung auf einem angepassten Intervall.) und diskutieren Sie die Resultate.</p>

1.3.4. Zweiter Durchgang zu Beispiel 1.2: Histogramm. Wie bei der Verteilungsfunktion machen wir einen Schritt in Richtung auf eine statistische Analyse. Der Einfachheit halber nehmen wir an, dass wir disjunkte Testmengen $A_j, j = 1, \dots, J$ gewählt haben, die den Wertebereich von X überdecken. Die Beobachtung (x_1, \dots, x_n) gibt dann Besetzungszahlen n_j

$$n_j = (\#i : X_i \in A_j).$$

Wenn $(X_i)_{i=1, \dots, n}$ unabhängig sind mit identischer Verteilung P , so ist $(n_j)_{j=1, \dots, J}$ ein Zufallsvektor mit Multinomialverteilung zu den Parametern $n, (p_j)_{j=1, \dots, J}$ mit $p_j = P(A_j)$. Für den Spezialfall $J = 2$ haben wir die Binomialverteilung. Da wir freie Wahl über die Testmengen A_j haben, können wir damit eine ganze Reihe von oft hilfreichen Spezialfällen abdecken, z.B.

Mediantest auf Symmetrie:

$$A_1 = \{x < x_{0.5}\} \quad A_2 = \{x \geq x_{0.5}\}$$

²Unterschiedliche Implementierungen können hier andere Aufrufstrukturen vorsehen. Der Kolmogorov-Smirnov-Test findet sich in `ks.test`.

Vor R Version 2.x gehörten diese jedoch nicht zum Basis-Umfang von R, sondern sind in speziellen Bibliotheken enthalten, die explizit hinzugeladen werden mussten. Die Bibliothek mit klassischen Tests in der R1.x-Implementierung heißt `ctest` und wird mit

```
library(ctest)
```

geladen.

Midrange-Test auf Konzentration:

$$A_1 = \{x_{0.25} \leq x < x_{0.75}\} \quad A_2 = \{x < x_{0.25} \text{ oder } x \geq x_{0.75}\}.$$

Für den allgemeinen Fall müssen wir jedoch die empirischen Besetzungszahlen n_j anhand der Multinomialverteilung beurteilen, und diese ist sehr unangenehm zu berechnen. Deshalb greift man oft auf Approximationen zurück. Auf Pearson geht folgende Approximation zurück:

LEMMA 1.8. (*Pearson*): Für $(p_j)_{j=1,\dots,J}$, $p_j > 0$ gilt im Limes $n \rightarrow \infty$ die Approximation

$$\begin{aligned} P_{mult}(n_1, \dots, n_j; n, p_1, \dots, p_j) &\approx \\ (2\pi n)^{-1/2} \left(\prod_{j=1,\dots,J} p_j \right)^{-1/2} &\cdot \\ \exp \left(-1/2 \sum_{j=1,\dots,J} \frac{(n_j - np_j)^2}{np_j} \right. & \\ -1/2 \sum_{j=1,\dots,J} \frac{n_j - np_j}{np_j} & \\ \left. + 1/6 \sum_{j=1,\dots,J} \frac{(n_j - np_j)^3}{(np_j)^2} + \dots \right). & \end{aligned}$$

BEWEIS. \rightarrow Wahrscheinlichkeitstheorie. Z.B. [JK70] p. 285. □

Der erste Term wird bestimmt von $\widehat{\chi^2} := \sum_{j=1,\dots,J} (n_j - np_j)^2 / np_j$. Dieser Term wird $\widehat{\chi^2}$ -Statistik genannt. Zumindest asymptotisch für $n \rightarrow \infty$ führen große Werte von $\widehat{\chi^2}$ zu kleinen Wahrscheinlichkeiten. Dies motiviert, die χ^2 -Statistik approximativ als Anpassungsmaß zu benutzen. Der Name kommt aus der Verteilungsasymptotik:

THEOREM 1.9. (*Pearson*): Für $(p_j)_{j=1,\dots,J}$, $p_j > 0$ ist im Limes $n \rightarrow \infty$ die Statistik

$$\widehat{\chi^2} := \sum_{j=1,\dots,J} \frac{(n_j - np_j)^2}{np_j}$$

χ^2 -verteilt mit $J - 1$ Freiheitsgraden.

Für eine formale Entscheidungsregel können wir wieder einen kritischen Wert χ_{krit}^2 festlegen, und die Hypothese verwerfen, dass die Beobachtungen (X_1, \dots, X_n) identisch uniform verteilte Zufallszahlen sind, wenn die χ^2 -Statistik über diesem Wert liegt. Wählen wir als kritischen Wert das obere α -Quantil der χ^2 -Verteilung, so wissen wir, dass bei Zutreffen der Hypothese der Wert χ_{krit}^2 oder ein höherer Wert höchstens mit Wahrscheinlichkeit α erreicht wird. Damit können wir zumindest asymptotisch auch hier unsere Irrtumswahrscheinlichkeit für eine ungerechtfertigte Ablehnung der Hypothese kontrollieren.

Die χ^2 -Tests gehören zum Basisumfang von R als Funktion `chisq.test()`. Sie sind so ausgelegt, dass sie für allgemeinere "Kontingenztafeln" genutzt werden können. Wir benötigen sie hier nur für einen Spezialfall: die Tafel ist in unserem Fall

der (eindimensionale) Vektor der Besetzungszahlen für vorgewählte Zellen. (Hinweis: in der R-Implementierung sind allgemeinere Varianten in `library(loglin)` zu finden.)

Aufgabe 1.13	
	<p>Orientieren Sie sich mit <code>help(chisq.test)</code> über die Aufrufstruktur der χ^2-Tests. Wenden Sie ihn für die Hypothese ($p_j = 1/J$), $J = 5$ auf folgende Vektoren von Besetzungszahlen an:</p> <p style="text-align: center;">(3 3 3 3 3) (1 2 5 3 3) (0 0 9 0 6).</p>

Aufgabe 1.14	
	<p>Welche Resultate erwarten Sie, wenn Sie die folgenden Vektoren mit dem χ^2-Test auf uniforme Verteilung testen:</p> <p style="text-align: center;">1:100 runif(100) sin(1:100) rnorm(100)?</p> <p>Führen Sie diese Tests durch und diskutieren Sie die Resultate. <i>Hinweis:</i> Die Funktion <code>chisq.test()</code> erwartet als Eingabe eine Häufigkeitstabelle. Die Prozedur <code>table()</code> gibt die Möglichkeit, Besetzungstabellen direkt zu erstellen (siehe <code>help(chisq.test)</code>). Sie können aber auch die Funktion <code>hist()</code> benutzen, und den Eintrag <code>counts</code> aus dem Resultat benutzen.</p>

Die Approximationen für die χ^2 -Statistik gelten zunächst nur, wenn die Zellen fest gewählt sind, unabhängig von der Information aus der Stichprobe. Praktische Histogramm-Algorithmen bestimmen jedoch Zellenanzahl und Zellgrenzen aufgrund der Stichprobe. Dazu werden (implizit) Parameter der Verteilung geschätzt. Unter bestimmten Voraussetzungen gilt noch immer eine χ^2 -Asymptotik, wie z.B. nach dem folgenden Theorem aus [Rao73, Section 6b.2]:

THEOREM 1.10. (i) Let the cell probabilities be the specified functions $\pi_1(\boldsymbol{\theta}), \dots, \pi_k(\boldsymbol{\theta})$ involving q unknown parameters $(\theta_1, \dots, \theta_q) = \boldsymbol{\theta}'$. Further let

- (a) $\hat{\boldsymbol{\theta}}$ be an efficient estimator of $\boldsymbol{\theta}$ in the sense of (5c.2.6),
- (b) each $\pi_i(\boldsymbol{\theta})$ admit continuous partial derivatives of the first order (only) with respect to θ_j , $j = 1, \dots, q$ or each $\pi_i(\boldsymbol{\theta})$ be a totally differentiable function of $\theta_1, \dots, \theta_q$, and

(c) the matrix $M = (\pi_r^{-1/2} \partial \pi_r / \partial \theta_s)$ of order $(k \times q)$ computed at the true values of θ is of rank q . Then the asymptotic distribution of

$$(1.1) \quad \chi^2 = \sum \frac{(n_i - n\hat{\pi}_i)^2}{n\hat{\pi}_i} = \sum \frac{(O - E)^2}{E}$$

is $\chi^2(k - 1 - q)$, where $\hat{\pi}_i = \pi_i(\hat{\theta})$.

BEWEIS. Siehe [Rao73] Abschnitt 6b.2. □

Aufgabe 1.15	
*	<p>Entwerfen Sie vergleichbare Testumgebungen für feste und für adaptive Zellwahlen.</p> <p>Ziehen Sie für feste und für adaptive Zellwahlen jeweils $s = 1000$ Stichproben aus <code>runif()</code> vom Umfang 50; berechnen Sie formal die χ^2-Statistik und plotten Sie deren Verteilungsfunktion.</p> <p>Vergleichen Sie die Verteilungsfunktionen.</p>

Wiederholte Stichproben

Wir haben uns bis jetzt darauf konzentriert, die Verteilung einer Zufallsvariablen zu untersuchen. Wir können das Verfahren fortsetzen. Wenn (X_1, \dots, X_n) identisch uniform verteilte Zufallszahlen sind, dann ist bei vorgewählten Zellen die χ^2 -Statistik approximativ χ^2 -verteilt, und $\kappa := \sqrt{n} \sup |F_n - F|$ hat asymptotisch die Kolmogorov-Smirnov-Verteilung.

Wir können wiederholt Stichproben $(X_{1j}, \dots, X_{nj})_{j=1..m}$ ziehen und daraus Statistiken $\widehat{\chi^2}_j$ und κ_j berechnen. Bei unabhängig, identisch verteilten Ausgangsdaten müssen diese nach χ^2 bzw. Kolmogorov-Smirnov verteilt sein. Bei diesen wiederholten Stichproben wird nicht nur die Verteilung der einzelnen Beobachtungen, sondern die gemeinsame Verteilung der jeweils n Stichprobenelemente untersucht.

Aufgabe 1.16	
	<p>Ziehen Sie für $n = 10, 50, 100$ wiederholt jeweils 300 Stichproben nach <code>runif()</code>. Berechnen Sie dafür jeweils die χ^2- und Kolmogorov-Smirnov-Statistik.</p> <p>Welchen χ^2-Test benutzen Sie?</p> <p>Plotten Sie die Verteilungsfunktionen dieser Statistiken und vergleichen Sie sie mit den theoretischen asymptotischen Verteilungen.</p> <p>Sprechen irgendwelche Befunde gegen die Annahme unabhängig uniform verteilter Zufallszahlen?</p> <p style="text-align: right;">(Fortsetzung) →</p>

Aufgabe 1.16	(Fortsetzung)
	<p><i>Hinweis:</i> die Funktionen für den χ^2- und Kolmogorov-Smirnov-Test speichern ihre Information intern als Liste. Um die Namen der Listenelemente zu bekommen, kann man sich ein Testobjekt generieren. Benutzen Sie z.B.</p> <p style="text-align: center;"><code>names(chisq.test(runif(100)))</code>.</p>

Güte

Die uniforme Verteilung war in unserer Diskussion bislang die angezielte Modellverteilung, unsere "Hypothese". Wir haben diskutiert, wie die unterschiedlichen Verfahren sich verhalten müssten, wenn diese Hypothese gilt. Das daraus abgeleitete Verteilungsverhalten kann dazu dienen, kritische Grenzen für formale Tests festzulegen. Wir verwerfen die Hypothese, wenn die beobachteten Test-Statistiken zu extrem sind. Was "zu extrem" ist, wird anhand der abgeleiteten Verteilungen bestimmt. Dies führt zu Entscheidungsregeln wie:

verwerfe die Hypothese, wenn $F_{\chi^2}(\widehat{\chi^2}) \geq 1 - \alpha$

oder

verwerfe die Hypothese, wenn $F_{Kolmogorov-Smirnov}(\kappa) \geq 1 - \alpha$

für geeignet festzulegende (kleine) Werte von α .

Wenn wir ein Entscheidungsverfahren formal festgelegt haben, können wir im nächsten Schritt fragen, wie scharf das Verfahren ist, wenn die Hypothese tatsächlich abzulehnen ist. Eine genauere Analyse bleibt der Statistik-Vorlesung vorbehalten. Mit den bis jetzt diskutierten Möglichkeiten können wir jedoch schon das Verhalten mit einer Monte-Carlo-Strategie untersuchen.

Als Simulations-Szenario wählen wir eine Familie von Alternativen. Die uniforme Verteilung fügt sich in die Beta-Verteilungen mit den Dichten

$$p_{a,b}(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1} \quad \text{für } a > 0, b > 0 \text{ und } 0 < x < 1$$

ein. Wir wählen als Alternativen Verteilungen aus dieser Familie. Daraus ziehen wir wiederholt Stichproben, und wenden jeweils formal unsere Entscheidungsverfahren an. Wir registrieren, ob das Verfahren zu einer Ablehnung der Hypothese führt oder nicht. Zu gegebener Wahl eines Stichprobenumfangs n und einer Wiederholungszahl m und bei Wahl einer Grenzwahrscheinlichkeit α erhalten wir eine Tabelle

$(a, b) \mapsto \#$ Simulationen, bei denen die Hypothese verworfen wird.

Speziell für die uniforme Verteilung $(a, b) = (1, 1)$ erwarten wir annähernd $m \cdot \alpha$ Verwerfungen. Für andere Verteilungen ist ein Verfahren um so entscheidungsschärfer, je größer der Anteil der Verwerfungen ist.

Aufgabe 1.17	
**	<p>Untersuchen Sie die Trennschärfe des Kolmogorov-Smirnov-Test und des χ^2-Tests. Wählen Sie jeweils einen Wert für n, m und α, und wählen Sie 9 Paare für (a, b). Notieren Sie die Überlegungen hinter Ihrer Wahl.</p> <p>Ziehen Sie zu diesen Parametern mit <code>rbeta()</code> Zufallsstichproben.</p> <p>Führen Sie jeweils den Kolmogorov-Smirnov-Test und einen χ^2-Test mit 10 gleichgroßen Zellen auf $(0, 1)$ durch.</p>
	<p>Wählen Sie Alternativparameter (a, b) so, dass Sie entlang der folgenden Geraden die Entscheidungsverfahren vergleichen können:</p> <ul style="list-style-type: none"> i) $a = b$ ii) $b = 1$ iii) $a = 1$ <p>und führen Sie eine entsprechende Simulation durch.</p>
	<p>Wählen Sie Alternativparameter (a, b) so, dass Sie für den Bereich $0 < a, b < 5$ die Entscheidungsverfahren vergleichen können.</p> <p>Ihre Schlüsse?</p> <p><i>Hinweis:</i> Mit <code>outer(x, y, fun)</code> wird eine Funktion <code>fun()</code> auf alle Paare aus den Werten von x, y angewandt und das Ergebnis als Resultat zurückgeliefert.</p> <p>Mit</p> <p style="text-align: center;"><code>contour()</code></p> <p>können Sie einen Contour-Plot erzeugen. Siehe <code>demo("graphic")</code>.</p>

Aufgabe 1.18	
**	<p>Entwerfen Sie eine Prüfstrategie, um "Pseudozufallszahlen" zu entlarven.</p> <p>Testen Sie diese Strategie an einfachen Beispielen</p> <ul style="list-style-type: none"> i) $x = 1..100 \pmod m$ für geeignete m ii) $\sin(x)$ $x = 1..100$ iii) ... <p>Werden diese als "nicht zufällig" erkannt?</p> <p>Versuchen Sie dann, die bereitgestellten Zufallszahlengeneratoren zu entlarven.</p>

1.4. Momente und Quantile

Verteilungsfunktionen oder Dichten sind mathematisch nicht einfach zu handhaben: der Raum der Funktionen ist im allgemeinen unendlich-dimensional und endliche geometrische Argumente oder endliche Optimierungsargumente sind nicht direkt anwendbar. Um die Analyse zu vereinfachen, greift man bisweilen auf endliche Beschreibungen zurück.

Historisch haben die Momente eine wichtige Rolle gespielt: Wahrscheinlichkeiten werden als Masse-Verteilungen interpretiert, und die Momente analog zu den Momenten der Mechanik eingeführt. Das erste Moment, entsprechend dem Schwerpunkt, heißt in der Statistik **Erwartungswert**.

DEFINITION 1.11. Ist X eine reellwertige Zufallsvariable mit Verteilung P , so ist der Erwartungswert von X definiert als

$$E_P(X) := E(X) := \int X dP,$$

falls das Integral existiert.

Das zweite Moment und höhere Momente werden konventionell zentriert. Für das zweite (zentrale) Moment, die **Varianz**, haben wir die folgende Definition:

DEFINITION 1.12. Ist X eine reellwertige Zufallsvariable mit Verteilung P , so ist die Varianz von X definiert als

$$Var_P(X) := Var(X) := \int (X - E(X))^2 dP.$$

Die Integralausdrücke müssen nicht immer definiert sein, d.h. die Momente müssen nicht immer existieren. Existieren sie jedoch, so geben sie eine erste Information über die Verteilung. Der Erwartungswert wird oft als das “statistische Mittel” interpretiert; die Wurzel aus der Varianz, die **Standardabweichung**, als “Streuung”.

Die Definitionen können auch auf empirische Verteilungen angewandt werden. Dies gibt einen ersten Weg, die Momente einer unbekanntem theoretischen Verteilung aus den Daten zu schätzen. Für den Mittelwert gilt Konsistenz:

$$E_P(E_{P_n}(X)) = E_P(X),$$

d.h. im statistischen Mittel stimmen empirischer Erwartungswert und Erwartungswert der zu Grunde liegenden Verteilung überein (falls definiert).

Für die Varianz gilt diese Konsistenz nicht, sondern es gilt

$$\frac{n}{n-1} E_P(Var_{P_n}(X)) = Var_P(X),$$

falls $n > 1$. Der mathematische Hintergrund ist, dass der Erwartungswert ein linearer Operator ist. Er kommutiert mit linearen Operatoren. Aber die Varianz ist ein quadratischer Operator, und das macht eine Korrektur nötig, wenn man Konsistenz will. Die entsprechend korrigierte Varianz wird oft als **Stichprobenvarianz** bezeichnet.

Für die Schätzung der ersten beiden Momente eines Vektor von Zufallszahlen stehen in R Funktionen bereit: **mean()** schätzt den Mittelwert und **var()** die

(Stichproben-)Varianz. Die Funktion `sd()` schätzt die Standardabweichung eines Vektors.

Aufgabe 1.19	
	<p>Generieren Sie jeweils eine Stichprobe von 100 Zufallsvariablen aus den Verteilungen mit den folgenden Dichten:</p> $p(x) = \begin{cases} 0 & x < 0 \\ 1 & 0 \leq x \leq 1 \\ 0 & x > 1 \end{cases}$ <p>sowie</p> $p(x) = \begin{cases} 0 & x \leq 0 \\ 2 & 0 < x \leq 1/4 \\ 0 & 1/4 < x \leq 3/4 \\ 2 & 3/4 < x \leq 1 \\ 0 & x > 1 \end{cases}$ <p>Schätzen Sie dazu Mittelwert, Varianz und Standardabweichung.</p> <p>Wiederholen Sie die Schätzung für 1000 Stichproben. Analysieren Sie die Verteilung von geschätztem Mittelwert, Varianz und Standardabweichung bei wiederholten Stichproben.</p>

Momente sind durch einfache arithmetische Operationen zu berechnen und ihre Kombination folgt (exakt oder approximierbar) einfachen Gesetzen. Sie sind jedoch sehr sensitiv. Die Verschiebung einer beliebig kleinen Wahrscheinlichkeitsmasse kann sie zum Zusammenbruch bringen. Für die empirische Verteilung bedeutet dies: stammen die beobachteten Daten zu einem Anteil $1 - \varepsilon$ aus einer Modellverteilung und zu einem Anteil ε aus einer anderen Verteilung, so können die Momente jeden beliebigen Wert annehmen, für jeden beliebig kleinen Wert von ε . Quantile sind gegenüber einem Zusammenbruch robuster als Momente. So müssen 50% der Daten "Ausreißer" sein, bis der Median beeinflusst wird, während das erste Moment, der Erwartungswert, schon bei Veränderung nur eines Datenpunkte beliebige Werte annehmen kann.

Mit der Verfügbarkeit von programmierbaren Rechnern haben Quantile als beschreibende Größe an Bedeutung gewonnen. Ihre Berechnung setzt implizit eine Sortier-Operation voraus, ist also komplexer als die Berechnung von Momenten. Auch die Regeln zur Kombination sind nicht so einfach wie bei Momenten und setzt oft eine explizite Rechnung voraus. Aber mit den verfügbaren technischen Mitteln ist dies keine wesentliche Einschränkung.

R bietet eine Reihe von Funktionen, um mit Quantilen zu arbeiten. `quantile()` ist eine elementare Funktion, um Quantile zu bestimmen. Die Funktion `summary()`

gibt eine Zusammenfassung der Verteilungsinformation, die auch auf Quantilen basiert ist.

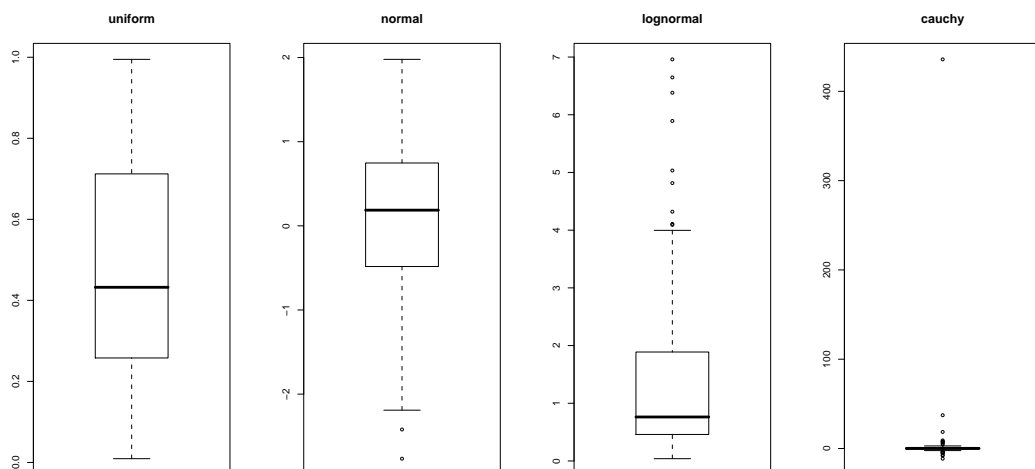
Aufgabe 1.20	
	<p>Generieren Sie jeweils eine Stichprobe von 100 Zufallsvariablen aus den Verteilungen von Aufgabe 1.19. Schätzen Sie dazu Median, oberes und unteres Quartil.</p> <p>Wiederholen Sie die Schätzung für 1000 Stichproben. Analysieren Sie die Verteilung von geschätztem Median, oberem und unterem Quartil bei wiederholten Stichproben.</p>

Mit `boxplot()` erhält man eine grafische Repräsentation dieser Zusammenfassung. Der hier benutzte “Box&Whisker-Plot” hat eine Reihe von Variationen. Deshalb ist es bei der Interpretation notwendig, sich jeweils über die benutzten Details zu informieren. Üblich ist eine Kennzeichnung durch eine “Box”, die den zentralen Teil der Verteilung beschreibt. In der Standardversion kennzeichnet eine Linie den Median, und eine “Box” darum reicht vom Median der oberen Hälfte bis zum Median der unteren Hälfte. Grob entspricht dies dem oberen und dem unteren Quartil. Die feinere Definition sorgt dafür, dass die Information auch noch verlässlich wieder gegeben wird, wenn Bindungen, d.h. vielfache Beobachtungen des selben Wertes auftreten. Die “Whisker” beschreiben die angrenzenden Bereiche. Ausreißer sind besonders gekennzeichnet.

Beispiel 1.10:

Eingabe

```
oldpar <- par(mfrow = c(1, 4))
boxplot(runif(100), main = "uniform")
boxplot(rnorm(100), main = "normal")
boxplot(exp(rnorm(100)), main = "lognormal")
boxplot(rcauchy(100), main = "cauchy")
par(oldpar)
```



Theorem 1.4 gibt eine Möglichkeit, Konfidenzintervalle für Quantile zu bestimmen, die allgemein gültig sind, unabhängig von der Form der zu Grunde liegenden Verteilung.

Um das p -Quantil x_p einer stetigen Verteilungsfunktion durch eine Ordnungsstatistik $X_{(k:n)}$ zum Konfidenzniveau $1 - \alpha$ nach oben abzuschätzen, suchen wir

$$\min_k : P(X_{(k:n)} \geq x_p) \geq 1 - \alpha.$$

Aber $X_{(k:n)} \geq x_p \iff F(X_{(k:n)}) \geq p$ und wegen Theorem 1.4 ist damit

$$P(X_{(k:n)} \geq x_p) = 1 - F_{beta}(p; k, n - k + 1).$$

Wir können also \min_k direkt aus der Beta-Verteilung ermitteln, oder wir benutzen die Beziehung zur Binomialverteilung und bestimmen k als

$$\min_k : P_{bin}(X \leq k - 1; n, p) \geq 1 - \alpha.$$

Aufgabe 1.21	
	<p>Für stetige Verteilungen und den Verteilungsmedian X_{med} ist $P(X_i \geq X_{med}) = 0.5$. Deshalb kann ein k so bestimmt werden, dass</p> $k = \min\{k : P(X_{(k)} \leq X_{med}) < \alpha\}$ <p>und $X_{(k)}$ als obere Abschätzung für den Median zum Konfidenzniveau $1 - \alpha$ gewählt werden.</p> <p>Konstruieren Sie mit dieser Idee ein Konfidenzintervall für den Median zum Konfidenzniveau $1 - \alpha = 0.9$.</p>
	<p>Modifizieren Sie den Box & Whiskerplot so, dass er dieses Intervall einzeichnet.</p> <p><i>Hinweis:</i> Sie benötigen dazu die Verteilungsfunktion F_X, ausgewertet an der durch die Ordnungsstatistik $X_{(k)}$ definierten Stelle. Die Verteilungen von $F_X(X_{(k)})$ wird in Theorem 1.4 diskutiert.</p>
	<p>Der Boxplot bietet eine Option <code>notch = TRUE</code>, um Konfidenzintervalle zu generieren. Versuchen Sie, mithilfe der Dokumentation herauszufinden, wie ein <code>notch</code> bestimmt wird. Vergleichen Sie Ihre Konfidenzintervalle mit den durch <code>notch</code> gekennzeichneten Intervallen.</p>
*	<p>Bestimmen Sie analog ein verteilungsunabhängiges Konfidenzintervall für den Interquartilsabstand.</p>
***	<p>Ergänzen Sie den Box & Whiskerplot so, dass er die Skaleninformation statistisch verlässlich darstellt.</p> <p><i>Hinweis:</i> Wieso reicht es nicht, Konfidenzintervalle für die Quartile einzuzeichnen?</p>

1.5. Ergänzungen

1.5.1. Ergänzung: Zufallszahlen. Wenn wir unabhängige identisch uniform verteilte Zufallszahlen hätten, könnten wir auch Zufallszahlen mit vielen anderen Verteilungen generieren. Z.B.

LEMMA 1.13. (*Inversionsmethode*): Ist (U_i) eine Folge unabhängiger $U[0, 1]$ verteilter Zufallsvariablen und F eine Verteilungsfunktion, so ist $(X_i) := (F^{-1}U_i)$ eine Folge unabhängiger Zufallsvariablen mit Verteilung F .

Analytisch ist dieses Lemma nur brauchbar, wenn F^{-1} bekannt ist. Numerisch hilft es jedoch viel weiter: anstelle von F^{-1} werden Approximationen benutzt, oft sogar nur eine Inversionstabelle.

Die Inversionsmethode ist eine Methode, aus gleichverteilten Zufallszahlen andere Zielverteilungen abzuleiten. Weitere (evtl. effektivere) Methoden, aus gleichverteilten Zufallszahlen andere Zielverteilungen abzuleiten, werden in der Literatur zur statistischen Simulation diskutiert.

Für eine Reihe von Verteilungen werden transformierte Zufallsgeneratoren bereitgestellt. Eine Liste ist im Anhang (Seite A-43) angegeben. Zu jeder Verteilungsfamilie gibt es dabei eine Reihe von Funktionen, deren Namen aus einem Kurznamen für die Verteilung abgeleitet sind. Für die Familie xyz ist $rxyz$ eine Funktion, die Zufallszahlen erzeugt. $dxyz$ berechnet die Dichte bzw. das Zählmaß für diese Familie, $pxyz$ die Verteilungsfunktion, und $qxyz$ die Quantile³.

Übersicht: einige ausgewählte Verteilungen. Weitere Verteilungen siehe A.23 (Seite A-43).

<i>Verteilung</i>	<i>Zufallszahlen</i>	<i>Dichte</i>	<i>Verteilungsfunktion</i>	<i>Quantile</i>
Binomial	<i>rbinom</i>	<i>dbinom</i>	<i>pbinom</i>	<i>qbinom</i>
Hypergeometrisch	<i>rhyper</i>	<i>dhyper</i>	<i>phyper</i>	<i>qhyper</i>
Poisson	<i>rpois</i>	<i>dpois</i>	<i>ppois</i>	<i>qpois</i>
Gauß	<i>rnorm</i>	<i>dnorm</i>	<i>pnorm</i>	<i>qnorm</i>
Exponential	<i>rexp</i>	<i>dexp</i>	<i>pexp</i>	<i>qexp</i>

1.5.2. Ergänzung: Grafische Vergleiche. Abweichungen von einfachen geometrischen Formen werden besser wahrgenommen als Abweichungen zwischen allgemeinen Grafen ähnlicher Form. Deshalb kann es hilfreich sein, Darstellungen zu wählen, die auf einfache Formen wie z.B. Geraden führen. So wählt man um zwei Verteilungsfunktionen F, G zu vergleichen anstelle der Funktionsgraphen den Graphen von

$$x \mapsto (F(x), G(x)).$$

³d.h. mit den in der Statistik üblichen Bezeichnungen ist verwirrender Weise $p_{xyz} \equiv dxyz$ und $F_{xyz} \equiv pxyz$.

Dieser Graph heißt *PP-Plot* oder *probability plot*. Stimmen die Verteilungen überein, so ist der Plot eine diagonale Gerade. Abweichungen von der Diagonalgestalt sind leicht zu erkennen.

Alternativ kann die Merkmalskala als Bezug genommen werden und der Graph von

$$p \mapsto (F^{-1}(p), G^{-1}(p))$$

betrachtet werden. Dieser Graph heißt *QQ-Plot* oder *Quantilplot*. Stimmen die Verteilungen überein, so zeigt auch dieser Plot eine diagonale Gerade.

Im Spezialfall der uniformen Verteilung auf $[0, 1]$ ist auf diesem Intervall $x = F(x) = F^{-1}(x)$, d.h. *QQ-Plot* und *PP-Plot* stimmen überein und sind der Graph der Verteilungsfunktion. Bei nicht-uniformen Verteilungen werden die Graphen im *PP-Plot* auf die Wahrscheinlichkeitsskala $[0, 1]$ standardisiert, und im *QQ-Plot* auf die Merkmalskala umskaliert.

Aufgabe 1.22	
	Erstellen Sie einen <i>PP-Plot</i> der $t(\nu)$ -Verteilung gegen die Standardnormalverteilung im Bereich $0.01 \leq p \leq 0.99$ für $\nu = 1, 2, 3, \dots$
	Erstellen Sie einen <i>QQ-Plot</i> der $t(\nu)$ -Verteilung gegen die Standardnormalverteilung im Bereich $-3 \leq x \leq 3$ für $\nu = 1, 2, 3, \dots$
	Wie groß muss ν jeweils sein, damit jeweils die t -Verteilung in diesen Plots kaum von der Normalverteilung zu unterscheiden ist?
	Wie groß muss ν sein, damit die t -Verteilung bei einem Vergleich der Verteilungsfunktionen kaum von der Normalverteilung zu unterscheiden ist?

Können die Verteilungen durch eine affine Transformation im Merkmalsraum ineinander überführt werden, so zeigt der *QQ-Plot* immer noch eine Gerade; Steigung und Achsenabschnitt repräsentieren die affine Transformation. Dies ist zum Beispiel so bei der Familie der Normalverteilungen: ist F die Standard-Normalverteilung $N(0, 1)$ und $G = N(\mu, \sigma^2)$, so ist der *QQ-Plot* eine Gerade mit Achsenabschnitt μ und Steigung σ .

Für empirische Verteilungen findet Korollar 1.5 Anwendung: anstelle von i/n wird ein für die Schiefe korrigierter Bezugspunkt gewählt, damit im Mittel eine Gerade erzeugt wird. Der Quantilplot mit dieser Korrektur für empirische Verteilungen ist als Funktion `qqplot()` bereitgestellt. Für den Spezialfall der Normalverteilung ist eine Variante von `qqplot()` als `qqnorm()` verfügbar, um eine empirische Verteilung mit der theoretischen Normalverteilung zu vergleichen.

Durch die Transformationen auf die Wahrscheinlichkeits- bzw. Merkmalskala gewinnen die graphischen Verfahren an Schärfe. So ist zum Beispiel selbst bei einem Stichprobenumfang von $n = 50$ die Verteilungsfunktion der Normalverteilung oft nur für den geübten Betrachter von der uniformen zu unterscheiden. Im Normal-*QQ-Plot*

hingegen zeigen sich uniforme Stichproben als deutlich nicht-linear, normalverteilte Daten hingegen geben weitgehend lineare Bilder.

Zur Illustration erzeugen wir uns zunächst zufällige Datensätze:

Eingabe

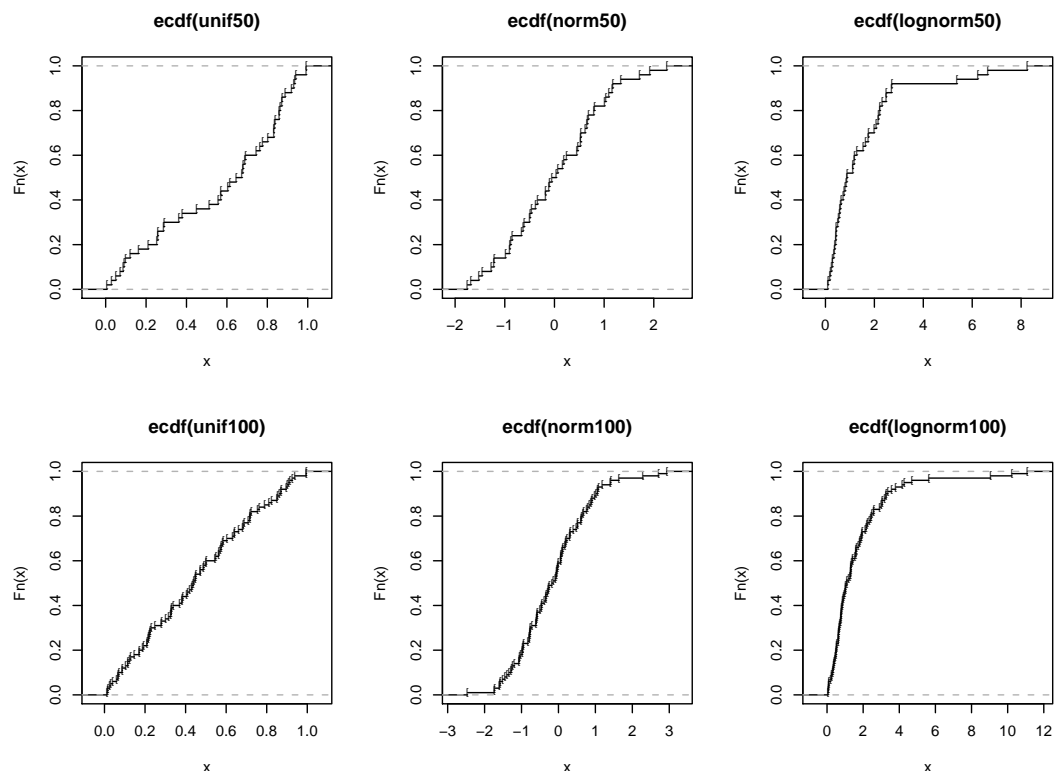
```
unif50 <- runif(50)
unif100 <- runif(100)
norm50 <- rnorm(50)
norm100 <- rnorm(100)
lognorm50 <- exp(rnorm(50))
lognorm100 <- exp(rnorm(100))
```

Mit diesen Datensätzen generieren wir Plots der Verteilungsfunktionen.

Beispiel 1.11:

Eingabe

```
oldpar <- par(mfrow = c(2, 3))
plot(ecdf(unif50), pch = "[")
plot(ecdf(norm50), pch = "[")
plot(ecdf(lognorm50), pch = "[")
plot(ecdf(unif100), pch = "[")
plot(ecdf(norm100), pch = "[")
plot(ecdf(lognorm100), pch = "[")
par(oldpar)
```

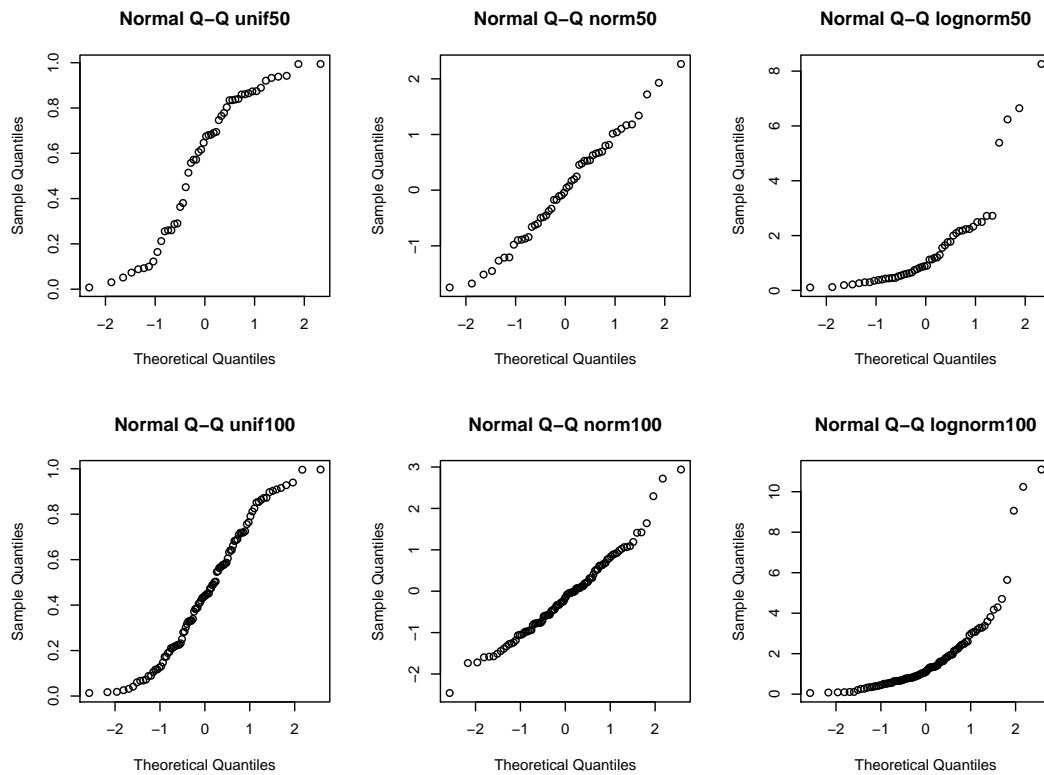


Zum Vergleich dazu die entsprechenden QQ -Plots für die selben Daten:

Beispiel 1.12:

Eingabe

```
oldpar <- par(mfrow = c(2, 3))
qqnorm(unif50, main = "Normal Q-Q unif50")
qqnorm(norm50, main = "Normal Q-Q norm50")
qqnorm(lognorm50, main = "Normal Q-Q lognorm50")
qqnorm(unif100, main = "Normal Q-Q unif100")
qqnorm(norm100, main = "Normal Q-Q norm100")
qqnorm(lognorm100, main = "Normal Q-Q lognorm100")
par(oldpar)
```



Aufgabe 1.23	
	Benutzen Sie <i>PP</i> -Plots anstelle von Verteilungsfunktionen, um die χ^2 - und Kolmogorov-Smirnov-Approximationen darzustellen.

Aufgabe 1.24	
	Benutzen Sie <i>QQ</i> -Plots anstelle von Verteilungsfunktionen. Können Sie in diesem Plot mit Hilfe der χ^2 - bzw. Kolmogorov-Smirnov-Statistik Konfidenzbereiche darstellen?

Um einen Eindruck über die Fluktuation zu bekommen, müssen wir empirische Plots mit typischen Plots einer Modellverteilung vergleichen. Eine Plot-Matrix ist ein einfacher Weg dazu. Wir geben hier ein Beispiel für den Normal- QQ -Plot, das wir gleich als Funktion implementieren:

```

Eingabe
qqnormx <- function(x, nrow = 5, ncol = 5, main = deparse(substitute(x))) {
  oldpar <- par(mfrow = c(nrow, ncol))
  qqnorm(x, main = main)
  for (i in 1:(nrow*ncol-1)) qqnorm(rnorm(length(x)), main = "N(0, 1)")
  par(oldpar)
}

```

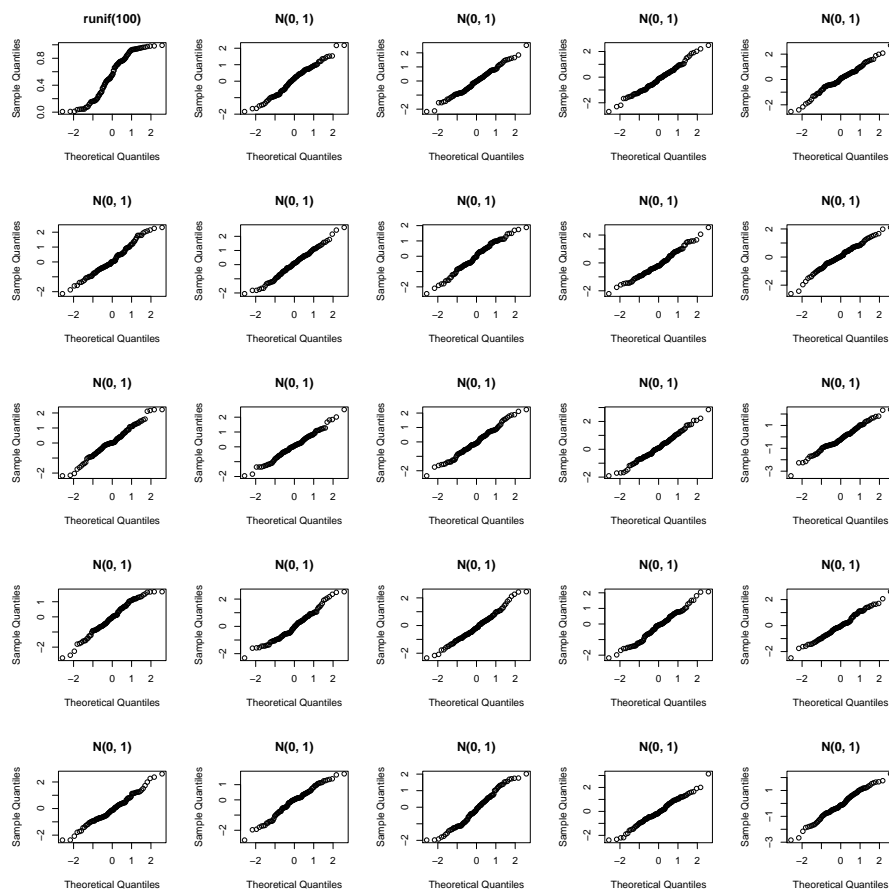
Wir haben in diesem Beispiel eine *for*-Schleife benutzt. Wie alle Programmiersprachen hat R Kontrollstrukturen, wie bedingte Anweisungen und Schleifen. Eine Übersicht über die Kontrollstrukturen in R ist im Anhang A.14 zu finden.

Abweichung von einer linearen Struktur ist als Fluktuation zu betrachten, wenn sie im Rahmen der simulierten Beispiele bleibt. Ist der zu untersuchende Datensatz extrem im Vergleich zu den simulierten Beispiele, so widerspricht das der Modellverteilung.

Beispiel 1.13:

```
qqnormx(runif(100))
```

Eingabe



Auf lange Sicht lohnt es sich, die Plot-Funktionen so zu modifizieren, dass sie auch Informationen über die zu erwartende Fluktuation wieder geben. In Beispiel 1.9 haben wir für die Verteilungsfunktion Monte-Carlo-Bänder konstruiert. Wir können diese Idee auf den *PP*-Plot und den *QQ*-Plot übertragen. Dazu ist es nur notwendig, die Bänder jeweils in der für den Plot geeigneten Skala darzustellen.

Aufgabe 1.25	
	<p>Erzeugen Sie sich mit <code>rnorm()</code> Pseudozufallszahlen für die Gaußverteilung zum Stichprobenumfang $n = 10, 20, 50, 100$.</p> <p>Erzeugen Sie jeweils einen <i>PP</i>-Plot und einen <i>QQ</i>-Plot, wobei die theoretische Gaußverteilung als Bezug dient.</p>
*	<p>Fügen Sie Monte-Carlo-Bänder aus der Einhüllenden von 19 Simulationen hinzu.</p> <p>Sie müssen zunächst anstelle der uniformen Verteilung die Normalverteilung zur Erzeugung der Monte-Carlo-Bänder benutzen. Sie müssen außerdem die Resultate im Koordinatensystem des <i>QQ</i>-Plots darstellen, d.h. die x-Achse repräsentiert die Quantile der Normalverteilung. Hinweis: inspizieren Sie dazu die Quelle von <code>qqnorm()</code>.</p> <p>Die Bänder sind zunächst Bänder für die Standard-Normalverteilung. Finden Sie Bänder für die vorliegenden Daten.</p>

1.5.3. Ergänzung: Grafik-Aufbereitung. Bislang wurde die R-Grafik in rudimentärer Form benutzt. Für ernsthafte Arbeit muss die Grafik so aufbereitet werden, dass ihre Bestandteile identifiziert und wiedererkennbar sind. Dazu gehören Überschriftungen, Achsenkennzeichnungen etc. R unterscheidet zwischen “high level”-Grafik und “low level”. “High level”-Funktionen erzeugen eine neue Grafik. Sie bieten darüber hinaus Möglichkeiten, allgemeine Grafikparameter zu steuern.

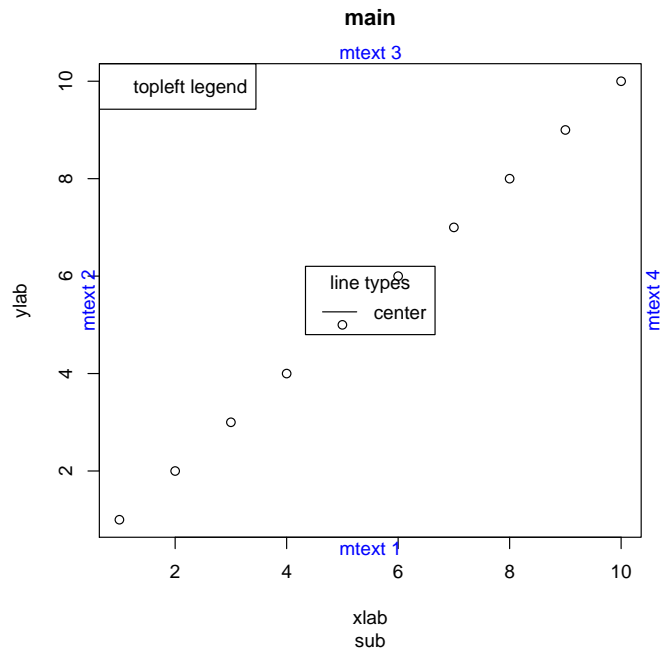
Die “low level”-Funktionen fügen Elemente zu vorhandenen Grafiken hinzu oder modifizieren die Grafik im Detail. Die Funktion `legend()` zum Beispiel kann Legenden innerhalb der Graphik hinzu fügen.

Beispiel 1.14:

```

Eingabe
plot(1:10, xlab = "xlab", ylab = "ylab", main = "main", sub = "sub")
mtext("mtext 1", side = 1, col = "blue")
mtext("mtext 2", side = 2, col = "blue")
mtext("mtext 3", side = 3, col = "blue")
mtext("mtext 4", side = 4, col = "blue")
legend("topleft", legend = "topleft legend")
legend("center", legend = "center" , lty = 1:4, title = "line types")

```

**Aufgabe 1.26**

Inspizieren Sie mit `help(plot)` die Steuerungsmöglichkeiten der plot-Funktion. Einige Detail-Information zu den Parametern erhalten Sie erst in `help(plot.default)`. Korrigieren Sie Ihren letzten Plot so, dass er eine korrekte Überschrift trägt.

Weitere Hinweise: [R D07a] Ch. 12.

1.5.4. Ergänzung: Funktionen. R-Kommandos können zu Funktionen zusammengefasst werden. Funktionen können parametrisiert sein. Funktionen erlauben eine flexible Wiederverwendbarkeit.

Beispiel für eine Funktion

Beispiel 1.15:

Eingabe

```

ppdemo <- function (x, samps = 19) {
  # samps: nr of simulations

  y <- (1:length(x))/length(x)
  plot(sort(x), y, xlab = substitute(x), ylab = expression(F[n]),
       main = "Verteilungsfunktion mit Monte-Carlo-Band (unif.)",
       type = "s")
  mtext(paste(samps, "Monte-Carlo-Stichproben"), side = 3)
  samples <- matrix(runif(length(x)* samps), nrow = length(x), ncol = samps)
  samples <- apply(samples, 2, sort)
  envelope <- t(apply(samples, 1, range))
  lines(envelope[, 1], y, type = "s", col = "red");
  lines(envelope[, 2], y, type = "s", col = "red")
}

```

Wir haben bei `ppdemo()` die Funktion `mtext()` benutzt, die Randbeschriftungen erlaubt.

Funktionen werden in der Form `<Name>(<Aktuelle Parameterliste>)` aufgerufen.

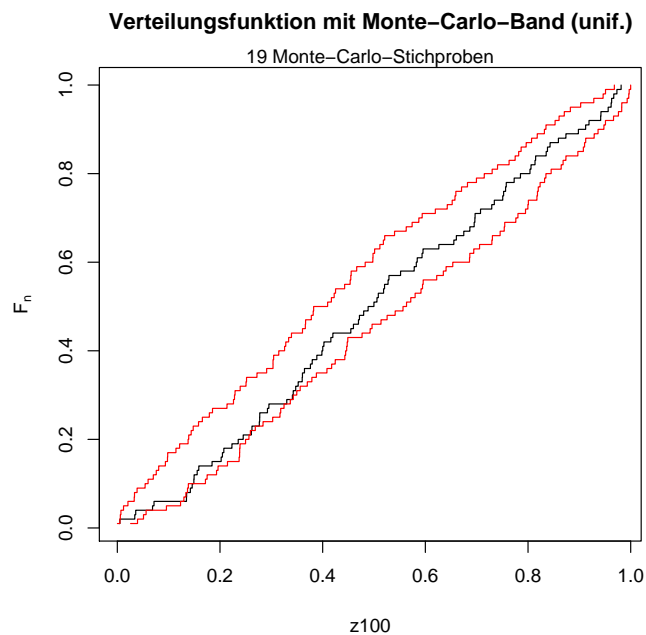
Beispiel 1.16:

Eingabe

```

z100 <- runif(100)
ppdemo(z100)

```



Wird nur der Name der Funktion eingegeben, so wird die Definition der Funktion zurück gegeben, d.h. die Funktion wird aufgelistet. Beispiel:

Beispiel 1.17:

<i>ppdemo</i>	Eingabe
<pre>function (x, samps = 19) { # samps: nr of simulations y <- (1:length(x))/length(x) plot(sort(x), y, xlab = substitute(x), ylab = expression(F[n]), main = "Verteilungsfunktion mit Monte-Carlo-Band (unif.)", type = "s") mtext(paste(samps, "Monte-Carlo-Stichproben"), side = 3) samples <- matrix(runif(length(x)* samps), nrow = length(x), ncol = samps) samples <- apply(samples, 2, sort) envelope <- t(apply(samples, 1, range)) lines(envelope[, 1], y, type = "s", col = "red"); lines(envelope[, 2], y, type = "s", col = "red") }</pre>	Ausgabe

Aufgabe 1.27	
	Inspizieren Sie <code>runif()</code> mit <code>qqplot()</code> und <code>plot()</code> . Überarbeiten Sie Ihre bisherigen Programmieraufgaben und schreiben Sie die wiederverwendbaren Teile als Funktionen.

Parameter bei Funktionen werden dem Wert nach übergeben. Jede Funktion erhält eine Kopie der aktuellen Parameterwerte. Dies sorgt für eine sichere Programmierumgebung. Auf der anderen Seite führt dies zu einer Speicherbelastung und bringt einen Zeitverlust mit sich. In Situationen, wo der Parameterumfang groß ist oder die Zeit eine kritische Größe ist, kann dieser Aufwand vermieden werden, indem direkt auf Variable zugegriffen wird, die in der Umgebung der Funktion definiert sind. Entsprechende Techniken sind in [GI00] beschrieben.

Funktionen in R können auch geschachtelt sein, d. h. innerhalb einer Funktion können auch wieder Funktionen definiert werden. Diese sind nur innerhalb der umgebenden Funktion sichtbar.

Funktionen können Objekte als Resultate haben. Ein Objekt wird explizit als Resultat übergeben mit `return(obj)`. Das Resultat kann auch implizit übergeben werden: wird das Ende einer Funktion erreicht, ohne dass `return()` aufgerufen wurden, so wird der Wert des letzten ausgewerteten Ausdruck übergeben.

<code>circlearea <- function(r) r^2 * pi</code> <code>circlearea(1:4)</code>	Eingabe
[1] 3.141593 12.566371 28.274334 50.265482	Ausgabe

Resultate können auch bereit gestellt werden, so dass sie nur auf Anfrage übergeben werden. Wir haben diese Technik beim Histogramm kennen gelernt. Der Aufruf `hist(x)` übergibt kein Resultat, sondern hat nur den (gewünschten) Seiteneffekt, ein Histogramm zu zeichnen. Benutzen wir `hist()` jedoch in einem Ausdruck, zum Beispiel in einer Zuweisung `xhist <- hist(x)`, so erhalten wir als Wert die Beschreibung des Histogramms. Um Resultate nur bei Bedarf zu übergeben, wird anstelle `return(obj)` der Aufruf `invisible(obj)` benutzt.

Aufgabe 1.28	
	<p>Schreiben Sie als Funktionen:</p> <ul style="list-style-type: none"> • Eine Funktion <code>ehist</code>, die ein Histogramm mit Ergänzungen zeigt. • Eine Funktion <code>ecdf</code>, die die empirische Verteilung zeigt. • Eine Funktion <code>eqnorm</code>, die einen <i>QQ</i>-Plot mit der Standard-Normalverteilung vergleicht. • Eine Funktion <code>eboxplot</code>, die einen Box&Whisker-Plot zeigt. • Eine zusammenfassende Funktion <code>eplot</code>, die eine Plot-Matrix mit diesen vier Plots zeigt. <p>Ihre Funktionen sollten die Standardfunktionen so aufrufen (oder modifizieren, falls notwendig), dass die Plots eine angemessene Beschriftung erhalten.</p>

Während Anweisungen in R schrittweise ausgeführt werden und so die Resultate bei jedem Schritt inspiziert werden können, werden beim Aufruf einer Funktion alle Anweisungen in der Funktion als Einheit ausgeführt. Dies kann eine Fehlersuche schwierig machen. R bietet Möglichkeiten, die Inspektion gezielt auf Funktionen zu ermöglichen. Details dazu finden sich im Anhang A.13 “Debugging und Profiling” auf Seite Seite A-21.

1.5.5. Ergänzung: Das Innere von R. Ein typischer Arbeitsabschnitt von R verarbeitet Kommandos in drei Teilschritten:

- `parse()` analysiert einen Eingabetext und wandelt ihn in ein interne Darstellung als R-Ausdruck. R-Ausdrücke sind spezielle R-Objekte.
- `eval()` interpretiert diesen Ausdruck und wertet ihn aus. Das Resultat ist wieder ein R-Objekt.
- `print()` zeigt das resultierende Objekt.

Details sind zu ergänzen:

1.5.5.1. *Parse*. Der erste Schritt besteht aus zwei Teilen: einem Leseprozess, der die Eingabe einscann und in Bausteine (Tokens) zerlegt, und dem eigentlichen Parsing, das die Bausteine falls möglich zu einem syntaktisch korrekten Ausdruck zusammenfasst. Die Funktion `parse()` fasst beide Schritte zusammen. Dabei kann `parse()` sowohl auf lokalen Dateien arbeiten, als auch auf externen, durch eine URL-Referenz bezeichneten Dateien.

Als inverse Funktion steht `substitute()` zur Verfügung. Eine typische Anwendung ist es, aktuelle Parameter eines Funktionsaufrufs zu entschlüsseln und informative Beschriftungen zu erzeugen.

1.5.5.2. *Eval.* Die Funktion `eval()` wertet einen R-Ausdruck aus. Dazu müssen die Referenzen im Ausdruck je nach den aktuellen Umgebungsbedingungen in entsprechende Werte übersetzt werden. Da R ein interpretiertes System ist, können die Umgebungsbedingungen variieren; je nach Umgebung kann derselbe Ausdruck zu unterschiedlichen Resultaten führen.

Jede Funktion definiert eine eigene lokale Umgebung. Funktionen können geschachtelt sein und somit auch die Umgebungen. Die Umgebung kann auch dadurch verändert werden, dass Zusatzpakete für R geladen werden. Die aktuelle unmittelbare Umgebung kann mit `environment()` erfragt werden. Mit `search()` erhält man eine Liste der Umgebungen, die sukzessive durchsucht werden, um Referenzen aufzulösen. Mit `ls()` erhält man eine Liste der Objekte in einer Umgebung.

Die Erweiterbarkeit von R bringt die Möglichkeit mit sich, dass Bezeichnungen kollidieren und damit die Übersetzung von Referenzen in aktuelle Werte fraglich wird. Als Schutz dagegen bietet R 2.x die Möglichkeit, Bezeichnungen in (geschützten) Namensräumen zusammen zu fassen. In den meisten Fällen ist dies transparent für den Benutzer; die Auflösung von Namen folgt der Suchreihenfolge, die durch die Kette der Umgebungen bestimmt ist. Um explizit auf Objekte eines bestimmten Namensraums zuzugreifen, kann dieser mit angegeben werden (z. B. `base::pi` als expliziter Name für die Konstante `pi` im Namensraum `base`).

1.5.5.3. *Print.* Die Funktion `print()` ist als *polymorphe* Funktion implementiert. Um `print()` auszuführen bestimmt R anhand der Klasse des zu druckenden Objekts eine geeignete Methode. Details folgen später in Abschnitt 2.6.5 Seite 2-39.

1.5.5.4. *Ausführung von Dateien.* Die Funktion `source()` steht bereit, um eine Datei als Eingabe für R zu benutzen. Dabei kann die Datei lokal sein, oder über eine URL-Referenz bezeichnet sein. Konventionell wird für die Namen von R-Kommandodateien die Endung `.R` benutzt.

Die Funktion `Sweave()` erlaubt es, Dokumentation und Kommandos miteinander zu verweben. Konventionell wird für die Namen von `Sweave()`-Eingabedateien die Endung `.Rnw` benutzt. Details zum Format finden sich in der `Sweave()`-Dokumentation <http://www.ci.tuwien.ac.at/~leisch/Sweave/Sweave-manual-20060104.pdf>.

1.5.6. Ergänzung: Pakete. Funktionen, Beispiele, Datensätze etc. können in R zu Paketen zusammengefasst werden, die bestimmten Konventionen entsprechen. Die Konventionen unterscheiden sich bei verschiedenen Implementierungen. Als aktuelle Referenz sollten die Konventionen von R [R D08] benutzt werden, denen wir auch hier folgen. Eine Reihe von Paketen sind Standardbestandteil von R. Pakete für spezielle Zwecke findet man im Internet z.B über <http://www.cran.r-project.org/src/contrib/PACKAGES.html>.

Nicht-Standard-Pakete müssen zunächst im R-System installiert werden. In der Regel gibt es dazu betriebssystem-spezifische Kommandos. Komfortabler ist jedoch die Installation aus R mit der Funktionen `install.packages()`. Ist keine spezielle Quelle angegeben, so greift `install.packages` dabei auf eine vorbereitete Adresse (in der Regel die oben angegebene) zurück. Sie können Pakete jedoch von jedem beliebigen Speicher laden. Insbesondere kann mit `install.packages(<package>, repos = NULL)` unter `<package>` ein direkter Zugriffspfad auf Ihrem Rechner angegeben werden.

Die Funktion `update.packages()` vergleicht installierte Versionen mit dem aktuellen Stand im Netz und frischt gegebenenfalls die installierte Version auf.

Installierte Pakete werden mit

library(pkgname)

geladen. Danach sind die im Paket definierten Objekte (Funktionen, Datensätze, ...) über den aktuellen Suchpfad auffindbar und direkt verwendbar.

Pakete werden mit

detach(pkgname)

wieder frei gegeben, d.h. ihre Objekte erscheinen nicht mehr im Suchpfad.

Technisch sind Pakete Verzeichnisse, die den R-Konventionen folgen. Üblich liegen sie in gepackter Form als .tar.gz-Files vor. In der Regel wird man zunächst als Benutzer vorbereitete Binärpakete installieren. Nur selten muss man auf die Quellpakete anderer Entwickler zurückgreifen.

Bei der Organisation der eigenen Arbeit lohnt es sich, den R-Konventionen zu folgen und zusammen gehörende Teile als R-Pakete zu organisieren. Dann stellt R eine ganze Reihe von Werkzeugen zur Unterstützung bereit. Die Konventionen und die bereitgestellten Werkzeuge sind in [R D08] dokumentiert. Für Unix/Linux/Mac OS X-Benutzer sind die wichtigsten Werkzeuge als Kommandos verfügbar:

R CMD check <directory> # überprüft ein Verzeichnis

R CMD build <directory> # generiert ein R-Paket

Als Einstieg: Die Funktion *package.skeleton()* hilft bei der Konstruktion neuer Pakete. *package.skeleton()* erzeugt dabei außer einem vorbereiteten Paket eine Hilfsdatei, die die weiteren Schritte zur Erzeugung eines ladbaren Pakets beschreibt.

Pakete müssen eine Datei DESCRIPTION mit bestimmter Information enthalten. Die Details sind in [R D08] beschrieben, und ein Prototyp wird von *package.skeleton()* erzeugt. Weiteres ist optional.

<i>Name</i>	<i>Art</i>	<i>Inhalt</i>
DESCRIPTION	Datei	eine Herkunftsbeschreibung nach Formatkonventionen.
R	Verzeichnis	R code. Dateien in diesem Verzeichnis sollten mit <i>source()</i> gelesen werden können. Empfohlene Namensendung: <i>.R</i> .
data	Verzeichnis	Zusätzliche Daten. Dateien in diesem Verzeichnis sollten mit <i>data()</i> gelesen werden können. Empfohlene Namensendungen und Formate: <i>.R</i> für R-Code. Alternativ: <i>.r</i> <i>.tab</i> für Tabellen. Alternativ: <i>.txt</i> , <i>.csv</i> <i>.RData</i> für Ausgaben von <i>save()</i> . Alternativ: <i>.rda</i> . Das Verzeichnis sollte eine Datei <i>00Index</i> mit einer Übersicht über die Datensätze enthalten.

(Fortsetzung)→

<i>Name</i>	<i>Art</i>	<i>Inhalt</i>
exec	Verzeichnis	Zusätzliche ausführbare Dateien, z.B. Perl- oder Shell-Skripte.
inst	Verzeichnis	Wird (rekursiv) in das Zielverzeichnis kopiert. Dieses Verzeichnis kann insbesondere eine Datei CITATION enthalten, die in R mit einer Funktion <code>citation()</code> ausgewertet wird.
man	Verzeichnis	Dokumentation im R-Dokumentationsformat (siehe: [R D08] "Writing R extensions", zugänglich über http://www.cran.r-project.org/). Empfohlene Namensendung: .Rd
src	Verzeichnis	Fortran, C und andere Quellen.
demo	Verzeichnis	ausführbare Beispiele. Dieses Verzeichnis sollte in einer Datei 00Index eine Beschreibung enthalten.

Aufgabe 1.29	
	<p>Installieren Sie die Funktionen der letzten Aufgaben als Paket. Das Paket sollte enthalten:</p> <ul style="list-style-type: none"> • Eine Funktion <code>ehist</code>, die ein Histogramm mit Ergänzungen zeigt. • Eine Funktion <code>ecdf</code>, die die empirische Verteilung zeigt. • Eine Funktion <code>eqnorm</code>, die einen <i>QQ</i>-Plot mit der Standard-Normalverteilung vergleicht • Eine Funktion <code>eboxplot</code>, die einen Box&Whisker-Plot zeigt. • Eine zusammenfassende Funktion <code>epplot</code>, die eine Plot-Matrix mit diesen vier Plots zeigt. <p>Sie können das Paket mit <code>package.skeleton()</code> vorbereiten, wenn Sie die einzelne Funktionen definiert haben.</p> <p>Laden Sie dieses Paket. Überprüfen Sie, ob Sie das Paket auch nach Neustart wieder mit <code>library()</code> laden können.</p> <p><i>Hinweis:</i> ist x ein Objekt, so erzeugt die Funktion <code>prompt(x)</code> ein Gerüst, aus dem eine Dokumentation für x entwickelt werden kann.</p>

1.6. Statistische Zusammenfassung

Als Leitbeispiel diente in diesem Kapitel die statistische Analyse einer (univariaten) Stichprobe. Dabei haben wir eine in der Statistik zentrale Modellvorstellung benutzt: die Werte der Stichprobe werden als Zufallsvariable aufgefasst, die aus einer zugrundeliegenden theoretischen Verteilung entstammen. Ziel der statistischen Analyse ist der Schluss aus der empirischen Verteilung der Stichprobe auf die unbekannte zu Grunde liegende theoretische Verteilung. Dieser Schluss kann zwei Formen annehmen: wir können die empirische

Verteilung mit einer hypothetischen Verteilung vergleichen. Dies ist das Vorgehen der klassischen Statistik. Oder wir können versuchen, aus der empirischen Verteilung Merkmale der zu Grunde liegenden Verteilung zu extrahieren. Dies ist das Vorgehen der Datenanalyse.

Beide Wege sind eng miteinander verwandt. Das wesentliche Werkzeug für beide war hier die Untersuchung der empirischen Verteilungsfunktion.

1.7. Literatur und weitere Hinweise:

[**R D08**] R Development Core Team (2000-2005): Writing R extensions.
Siehe: <http://www.r-project.org/manuals.html>.

[**GS77**] Gänßler, P; Stute, W.: Wahrscheinlichkeitstheorie. Heidelberg: Springer 1977.

[**GI00**] Gentleman, R.; Ihaka, R.: Lexical Scope and Statistical Computing. Journal of Computational and Graphical Statistics 9 (2000) 491–508.

KAPITEL 2

Regression

2.1. Allgemeines Regressionsmodell

Aus der Tradition experimenteller Wissenschaften stammt das Paradigma des (kontrollierten) Versuchs. Unter Versuchsbedingungen x wird ein Resultat y gemessen, zusammengesetzt aus einem systematischen Effekt $m(x)$ und einem Messfehler ε .

$$y = m(x) + \varepsilon.$$

Dies ist eine ganz spezielle Betrachtungsweise; es wird nicht unvoreingenommen das gemeinsame Verhalten von x und y untersucht, sondern eine Unsymmetrie hineingesteckt: x ist die "Ursache", y (oder eine Veränderung von y) der Effekt. Die "Ursache" x ist in dieser Vorstellung vorgegeben oder direkt kontrollierbar; y ist mittelbar (über die Versuchsbedingungen) beeinflusst. In dieser Vorstellung ist ε ein Messfehler, der durch geeignete Rahmenbedingungen des Versuchs möglichst klein gehalten wird und im Mittel verschwinden sollte: es sollte keinen systematischen Fehler geben.

Vom statistischen Standpunkt ist der wesentliche Unterschied der Rollen von x und y , dass für y das stochastische Verhalten mithilfe von ε modelliert wird, während x als "gegeben" angenommen wird und dafür keine Stochastik im Modell vorgesehen ist.

Um einen formal überschaubaren Rahmen zu bekommen, betrachten wir den Fall, dass x als Vektor von reellen Variablen repräsentiert werden kann, $x \in \mathbb{R}^p$, und dass die Messwerte eindimensionale reelle Werte sind, $y \in \mathbb{R}$. In einem stochastischen Modell kann die oben skizzierte Idee formal gefasst werden. Eine mögliche Formalisierung ist es, den Messfehler ε als Zufallsvariable zu modellieren. Nehmen wir ferner an, dass der Erwartungswert von ε existiert, so können wir die Annahme, dass der Messfehler im Mittel verschwindet, formalisieren als $E(\varepsilon) = 0$.

Um den systematischen Effekt m zu untersuchen, betrachten wir Messreihen. Der Index $i, i = 1, \dots, n$, kennzeichnet den Messpunkt, und das Modell ist damit

$$y_i = m(x_i) + \varepsilon_i \quad i = 1, \dots, n$$

mit $x_i \in \mathbb{R}^p$
 $E(\varepsilon_i) = 0$.

Das statistische Problem ist:

schätze die Funktion m aus den Messwerten y_i bei Messbedingung x_i .

Zu diesem Problem der Kurvenschätzung oder "Regression" gibt es eine umfangreiche Literatur in der Statistik. Wir wollen uns hier auf das "computing" konzentrieren. Dazu betrachten wir zunächst eine sehr vereinfachte Version des Regressionsproblems, die lineare Regression. Wesentliche Aspekte lassen sich bereits an diesem Problem illustrieren.

Einer einheitlichen Sprechweise zuliebe nennen wir y_i die **Respons** und die Komponenten von x_{ij} mit $j = 1, \dots, p$ die **Regressoren**. Die Funktion m heißt die **Modellfunktion**.

2.2. Lineares Modell

Wir beginnen mit dem Regressionsmodell - jetzt in Vektorschreibweise¹ -

$$(2.1) \quad \begin{aligned} Y &= m(X) + \varepsilon \\ Y &\text{ mit Werten in } \mathbb{R}^n \\ X &\in \mathbb{R}^{n \times p} \\ E(\varepsilon) &= 0 \end{aligned}$$

und setzen zusätzlich voraus, dass m linear ist. Dann gibt es (mindestens) einen Vektor $\beta \in \mathbb{R}^p$, so dass

$$m(X) = X\beta$$

und das Regressionsproblem ist jetzt reduziert auf die Aufgabe, β aus der Information (Y, X) zu schätzen.

Das so modifizierte Regressionsmodell

$$(2.2) \quad \begin{aligned} Y &= X\beta + \varepsilon \\ Y &\text{ mit Werten in } \mathbb{R}^n \\ X &\in \mathbb{R}^{n \times p} \\ \beta &\in \mathbb{R}^p \\ E(\varepsilon) &= 0 \end{aligned}$$

heißt **lineares Modell** oder auch **lineare Regression**. Die Matrix X , in der die Werte der Regressoren zusammengefasst ist, also die Information über die Versuchsbedingungen, heißt **Design-Matrix** des Modells.

BEISPIEL 2.1. (*Einfache lineare Regression*) Wird die Versuchsbedingung durch den Wert einer reellen Variablen beschrieben, von der der Versuchsausgang über eine lineare Modellfunktion

$$m(x) = a + b \cdot x,$$

abhängt, so können wir eine Versuchsserie mit Versuchsausgang y_i bei Versuchsbedingung x_i koordinatenweise schreiben als

$$(2.3) \quad y_i = a + b \cdot x_i + \varepsilon_i.$$

In Matrixschreibweise können wir die Versuchsserie zusammenfassen als lineares Modell

$$(2.4) \quad Y = \underbrace{\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}}_X \cdot \underbrace{\begin{pmatrix} a \\ b \end{pmatrix}}_\beta + \varepsilon.$$

Die einfache lineare Regression ist ein Beispiel für lineare Modelle. Die Einweg-Klassifikation ist das andere Basis-Beispiel.

BEISPIEL 2.2. (*Einweg-Klassifikation*) Zum Vergleich von k Behandlungen (insbesondere für den Spezialfall $k = 2$) benutzen wir Indikatorvariablen, die in einer Matrix zusammengefasst werden. Die Indikatorvariable für Behandlung i steht in Spalte i . In der Regel

¹Wir wechseln Konventionen und Schreibweisen, wenn es hilfreich ist. Die Verwirrung gehört zu den Konventionen: in einigen Konventionen kennzeichnen Großbuchstaben Zufallsvariable, in anderen Funktionen, in wieder anderen Vektoren. Die Auflösung bleibt jeweils dem Leser überlassen.

haben wir wiederholte Beobachtungen $j = 1, \dots, n_i$ unter Behandlung i , insgesamt also $n = \sum_{i=1}^k n_i$ Beobachtungen. Dem Modell

$$(2.5) \quad Y = \underbrace{\begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & \dots & 0 \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}}_X \cdot \underbrace{\begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{pmatrix}}_{\beta} + \varepsilon.$$

entspricht in Koordinaten

$$(2.6) \quad y_{ij} = \mu_i + \varepsilon_{ij}.$$

Dies ist das typische Modell, um die Hypothese “kein Unterschied” $\mu_1 = \dots = \mu_k$ gegen die Alternative zu testen, dass sich die Behandlungen im Mittel unterscheiden.

Der selbe Zusammenhang kann auch dargestellt werden, wenn wir die Messwerte als Summe eines Grundwertes μ_0 und dazu eines Behandlungseffekts $\mu'_i = \mu_i - \mu_0$ interpretieren. Dies entspricht in Koordinaten

$$(2.7) \quad y_{ij} = \mu_0 + \mu'_i + \varepsilon_{ij}.$$

In Matrixschreibweise ist dies

$$(2.8) \quad Y = \underbrace{\begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & \dots & 0 \\ 1 & 0 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{pmatrix}}_{X'} \cdot \underbrace{\begin{pmatrix} \mu_0 \\ \mu'_1 \\ \mu'_2 \\ \vdots \\ \mu'_k \end{pmatrix}}_{\beta'} + \varepsilon.$$

Beispiel 2.2 illustriert, dass die Darstellung eines Problems als lineares Modell nicht eindeutig ist. (2.5) und (2.8) sind gleichwertige Darstellungen und nur aus der Anwendung kann entschieden werden, welche den Vorrang hat.

Für die mathematische Analyse ist die Design-Matrix X ein wesentliches Hilfsmittel. Für die Datenanalyse können wir R zu Hilfe nehmen, um diese Matrix (implizit) für uns zu erstellen. R versteht eine spezielle Notation, die **Wilkinson-Rogers-Notation**, mit der Modelle beschrieben werden können. In dieser Notation schreiben wir

$$y \sim x.$$

Der Fehlerterm wird in diesem Modell nicht notiert.

Der konstante Term wird implizit angenommen. Für die Einweg-Klassifikation erhalten wir also das Modell (2.8). Wenn wir keinen konstanten Term wollen (also bei der Regression die Regressionsgerade durch den Ursprung geht, bzw. bei der Einweg-Klassifikation das Modell (2.5) benutzt und kein Gesamtmittel vorgegeben sein soll), so haben wir in Koordinatenschreibweise

$$y_i = b \cdot x_i + \varepsilon_i.$$

In der Wilkinson-Rogers-Notation muss der konstante Term explizit auf Null gesetzt werden:

$$y \sim 0 + x.$$

Weitere Regressoren können mit dem Operator `+` gekennzeichnet werden. So entspricht $y \sim u + v$ in Koordinaten dem Modell

$$y_i = a + b \cdot u_i + c \cdot v_i + \varepsilon_i.$$

Wir kommen in den Abschnitten 2.2.4 und 2.3 noch auf diese Notation zurück.

Es gibt eine umfangreiche Literatur über lineare Modelle. Das Buch “The Theory of Linear Models” von Bent Jørgensen [Jør93] ist besonders zu empfehlen. Es deckt den mathematischen Hintergrund dieses Kapitels weitgehend ab und enthält zahlreiche illustrierende Beispiele.

2.2.1. Faktoren. Mit Hilfe der Notation zur Design- und Modellbeschreibung kann die Übersetzung zwischen einer fallorientierten Beschreibung eines Designs in eine Design-Matrix für ein lineares Modell in kanonischer Form automatisch geschehen. Bisweilen braucht die Übersetzung etwas Nachhilfe. Betrachten Sie z.B. einen Datensatz

```
y <- c( 1.1, 1.2, 2.4, 2.3, 1.8, 1.9)
x <- c( 1, 1, 2, 2, 3, 3).
```

Der Vektor x kann als quantitativer Vektor für das Regressions-Modell

$$y_i = a + b x_i + \varepsilon_i$$

als Regressor gemeint sein, oder es kann in der Einweg-Klassifikation, dem Modell der Einweg-Varianzanalyse,

$$y_{ix} = \mu + \alpha_x + \varepsilon_{ix}$$

die Kennzeichnung einer Behandlungsgruppe sein. Um beide Möglichkeiten zu unterscheiden, können Vektoren in R als **Faktoren** definiert werden. Vektoren, die keine Faktoren sind, werden als quantitative Variable behandelt wie im ersten Beispiel. Faktoren werden als Kennzeichner behandelt und in der Design-Matrix in entsprechende Indikatorvariable übersetzt. So ergibt

$$y \sim x$$

das Regressionsmodell, jedoch

$$y \sim \text{factor}(x)$$

das Varianz-Modell für das Einweg-Layout.

Durch einen Parameter `ordered = TRUE` kann beim Aufruf der Funktion `factor()` die erzeugte Variable als geordnet gekennzeichnet werden. Die erzeugte Variable wird dann bei den Auswertungen als ordinal skaliert behandelt.

$$y \sim \text{factor}(x, \text{ordered} = \text{TRUE})$$

Ohne diese Kennzeichnung werden Faktoren als kategorial skaliert betrachtet.

Die Zahlenwerte von Faktoren brauchen keine aufsteigende Folge zu sein. Sie werden (auch für ordinale Faktoren) als bloße Namen benutzt und durch eine laufende Nummer ersetzt. So ergibt

```
factor( c(2, 2, 5, 5, 4, 4) )
```

einen Vektor mit drei Faktorwerten 1, 2, 3, die die Namen "2", "5" und "4" haben. Faktoren können auch durch Namen bezeichnet werden, z.B.

```
y ~ factor ( c("Beh1", "Beh1", "Beh2", "Beh2", "Beh3", "Beh3") )
```

Die unterschiedlichen Werte eines Faktors nennt man **Stufen** des Faktors. Sie können mit `levels()` erfragt werden, z.B.

```
levels(factor( c(2, 2, 5, 5, 4, 4) ))
```

```
levels(factor( c("Beh1", "Beh1", "Beh2", "Beh2", "Beh3", "Beh3") ))
```

2.2.2. Kleinste-Quadrate-Schätzung. Eine erste Idee zur Schätzung im linearen Regressionsmodell kann so gewonnen werden: Bei gegebenem X ist $E(Y) = X\beta$, also $X^T E(Y) = X^T X\beta$ und damit $(X^T X)^{-1} X^T E(Y) = \beta$. Dabei bedeutet X^T die transponierte Matrix zu X und $(X^T X)^{-1}$ die (generalisierte) Inverse von $(X^T X)$. Die Gleichung motiviert das folgende Schätzverfahren:

$$(2.9) \quad \hat{\beta} = (X^T X)^{-1} X^T Y.$$

Setzt man aus 2.2 die Modellbeziehung $Y = X\beta + \varepsilon$ ein und benutzt, dass $E(\varepsilon) = 0$, so erhält man

$$(2.10) \quad E(\hat{\beta}) = E\left((X^T X)^{-1} X^T (X\beta + \varepsilon) \right) = \beta,$$

d.h. $\hat{\beta}$ ist ein erwartungstreuer Schätzer für β . Ob und wie weit dieser Schätzer neben dieser Konsistenz auch noch statistische Qualitäten hat, wird in Statistik-Vorlesungen diskutiert. Ein Satz zur Charakterisierung dieses Schätzers ist dort als Gauß-Markov-Theorem bekannt. Wir werden auf diesen Schätzer häufig zurückkommen und geben ihm deshalb einen Namen: **Gauß-Markov-Schätzer**. Im Fall eines linearen Modells, wie dem Regressionsmodell, hat dieser Schätzer eine Reihe von Optimalitätseigenschaften. So minimiert dieser Schätzer die mittlere quadratische Abweichung, ist also in diesem Modell ein **Kleinste-Quadrate-Schätzer**.

Der Kleinste-Quadrate-Schätzer für lineare Modelle wird durch die Funktion `lm()` berechnet.

Zu Illustration erzeugen wir uns einen Beispiel-Datensatz.

```
x <- 1:100
err <- rnorm(100, mean = 0, sd = 10)
y <- 2.5*x + err
```

Den Kleinste-Quadrate-Schätzer erhalten wir nun durch

Beispiel 2.1:

```
lm(y ~ x)
```

Eingabe

Call:

```
lm(formula = y ~ x)
```

Ausgabe

Coefficients:

```
(Intercept)          x
      1.460         2.492
```

Aufgabe 2.1

Aufgabe 2.1	
	Wir haben die Daten ohne konstanten Term erzeugt, dies aber bei der Schätzung nicht vorausgesetzt. Wiederholen Sie die Schätzung im Modell ohne konstanten Term. Vergleichen Sie die Resultate.

Der Schätzer $\hat{\beta}$ führt unmittelbar zu einer Schätzung \hat{m} für die Funktion m in unserem ursprünglichen Modell:

$$\hat{m}(x) = x^T \cdot \hat{\beta}.$$

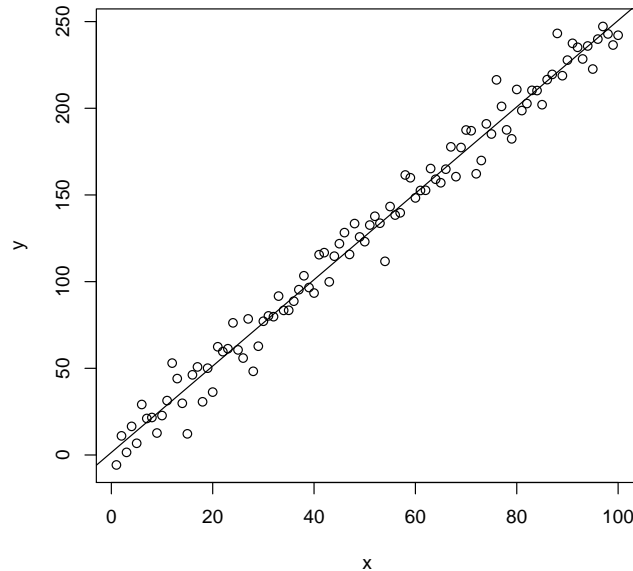
Die Auswertung an einem Punkt x ergibt Werte $\hat{y} := \hat{m}(x)$, den **Fit** an der Stelle x . Die Auswertung an den Messpunkten ergibt den Vektor der gefitteten Werte $\hat{Y} = X\hat{\beta}$.

In unserem Beispiel ist dies eine Regressionsgerade. Mit `plot()` können wir die Datenpunkte zeichnen. Wenn wir das Resultat der Regression speichern, können wir mit `abline()` die Regressionsgerade hinzufügen.

Beispiel 2.2:

Eingabe

```
lmres <- lm(y ~ x)
plot(x, y)
abline(lmres)
```



`abline()` ist eine Funktion, die Geraden anhand von unterschiedlichen Parametrisierungen zeichnen kann. Weiter Information erhalten sie mit `help(abline)`.

Die Schätzgleichung (2.9) gibt uns an, wie der Fit an den Messpunkten berechnet wird.

$$(2.11) \quad \hat{Y} = X(X^T X)^{-1} X^T \cdot Y.$$

Die Matrix

$$(2.12) \quad H := X(X^T X)^{-1} X^T$$

nennt man **Hut-Matrix**². Sie ist das wesentliche Werkzeug, um den Gauß-Markov-Schätzer für eine bestimmte Design-Matrix X zu untersuchen. Die Design-Matrix, und damit die Hutmatrix, hängt nur von den Versuchsbedingungen ab, nicht aber von dem Ausgang des Versuchs. Der Fit hingegen bezieht sich auf eine bestimmte Stichprobe, die in den in den beobachteten Stichprobenwerten Y repräsentiert ist.

Im linearen Modell ist ein Term ε enthalten, der den Messfehler oder die Versuchsvariabilität repräsentiert. Diesen stochastischen Fehler können wir nicht direkt beobachten - sonst könnten wir ihn subtrahieren und damit die Modellfunktion exakt bestimmen. Wir können nur mittelbar darauf schließen.

Die Werte der Zufallsbeobachtung Y unterscheidet sich in der Regel vom Fit \hat{Y} . Die Differenz

$$R_X(Y) := Y - \hat{Y}$$

²sie setzt dem Y den Hut auf.

heißt **Residuum**. Das Residuum kann als Schätzer für den nicht-beobachtbaren Fehlerterm ε angesehen werden. Die Residuen sind nicht wirklich der Fehlerterm. Dies wäre nur der Fall, wenn die Schätzung exakt wäre. Für den allgemeinen Fall zeigt uns die Beziehung

$$\begin{aligned}
 R_X(Y) &= Y - \hat{Y} \\
 &= (I - H)Y \\
 &= (I - H)(X\beta + \varepsilon) \\
 &= (I - H)\varepsilon,
 \end{aligned}
 \tag{2.13}$$

dass die Residuen Linearkombinationen der Fehler sind. Wir müssen aus diesen Linearkombinationen auf die Fehler zurück schließen.

Existiert die Varianz der Fehlerterme, so ist durch die Varianzmatrix Σ der Fehlerterme $\text{Var}(\varepsilon) = \Sigma$ die Varianz der Residuen bestimmt:

$$\begin{aligned}
 \text{Var}(R_X(Y)) &= \text{Var}((I - H)\varepsilon) \\
 &= (I - H)\Sigma(I - H)^\top.
 \end{aligned}
 \tag{2.14}$$

Bislang haben wir nur vorausgesetzt, dass kein systematischer Fehler vorliegt, modelliert als die Annahme

$$E(\varepsilon) = 0.$$

Wir sprechen von einem **einfachen linearen Modell**, wenn darüber hinaus gilt:

$$\begin{aligned}
 (\varepsilon_i)_{i=1, \dots, n} &\quad \text{sind unabhängig} \\
 \text{Var}(\varepsilon_i) = \sigma^2 &\quad \text{für ein } \sigma \text{ das nicht von } i \text{ abhängt.}
 \end{aligned}$$

Im linearen Modell versuchen wir, den Parametervektor β zu schätzen. Die Varianzstruktur des Fehlervektors bringt dabei Störparameter mit sich, die die Schätzung verkomplizieren können. Im einfachen linearen Modell reduziert sich die Situation auf nur einen unbekannt Störparameter σ . Formeln wie eq:02-varerr vereinfachen sich, denn in diesem Fall ist $\Sigma = \sigma^2 I$ und der Parameter σ kann aus der Formel herausgezogen werden. Wir können diesen Parameter aus den Residuen schätzen, denn die **residuelle Varianz**

$$s^2 := \frac{1}{n - \text{Rk}(X)} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2
 \tag{2.15}$$

ist ein erwartungstreuer Schätzer für σ^2 . Wir schreiben deshalb auch $\widehat{\sigma^2} := s^2$. (Das Wurzelziehen ist keine lineare Operation und erhält deshalb nicht den Erwartungswert. Die residuelle Standardabweichung $\sqrt{s^2}$ ist kein erwartungstreuer Schätzer für σ .) Wieder in die Schätzformel (2.9) eingesetzt liefert uns dies auch eine Schätzung für die Varianz/Covarianzmatrix des Schätzers für β , denn im einfachen Modell ist

$$\text{Var}(\hat{\beta}) = \sigma^2 X^\top X
 \tag{2.16}$$

und kann durch $s^2 X^\top X$ geschätzt werden.

Die Standard-Ausgabe in Beispiel 2.1 listet nur minimale Information über den Schätzer. Mehr Information über Schätzer, Residuen und daraus abgeleitete Kenngrößen erhalten wir, wenn wir eine zusammengefasste Darstellung anfordern.

Beispiel 2.3:

```
summary(lm( y ~ x))
```

Eingabe

Call:

```
lm(formula = y ~ x)
```

Ausgabe

Residuals:

```
      Min       1Q   Median       3Q      Max
-26.6301  -4.8625   0.2448   6.7120  25.5667
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   1.4602     1.9487   0.749   0.455
x              2.4918     0.0335  74.380 <2e-16 ***
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 9.671 on 98 degrees of freedom
```

```
Multiple R-squared:  0.9826,    Adjusted R-squared:  0.9824
```

```
F-statistic:  5532 on 1 and 98 DF,  p-value: < 2.2e-16
```

Aufgabe 2.2	
	<p>Analysieren Sie die in Beispiel 2.3 (Seite 2-8) gezeigten Ausgaben von <code>lm()</code>. Welche Terme können Sie interpretieren? Stellen Sie diese Interpretationen schriftlich zusammen. Für welche Terme fehlt Ihnen noch Information?</p> <p>Erstellen Sie eine kommentierte Version der Ausgabe.</p>

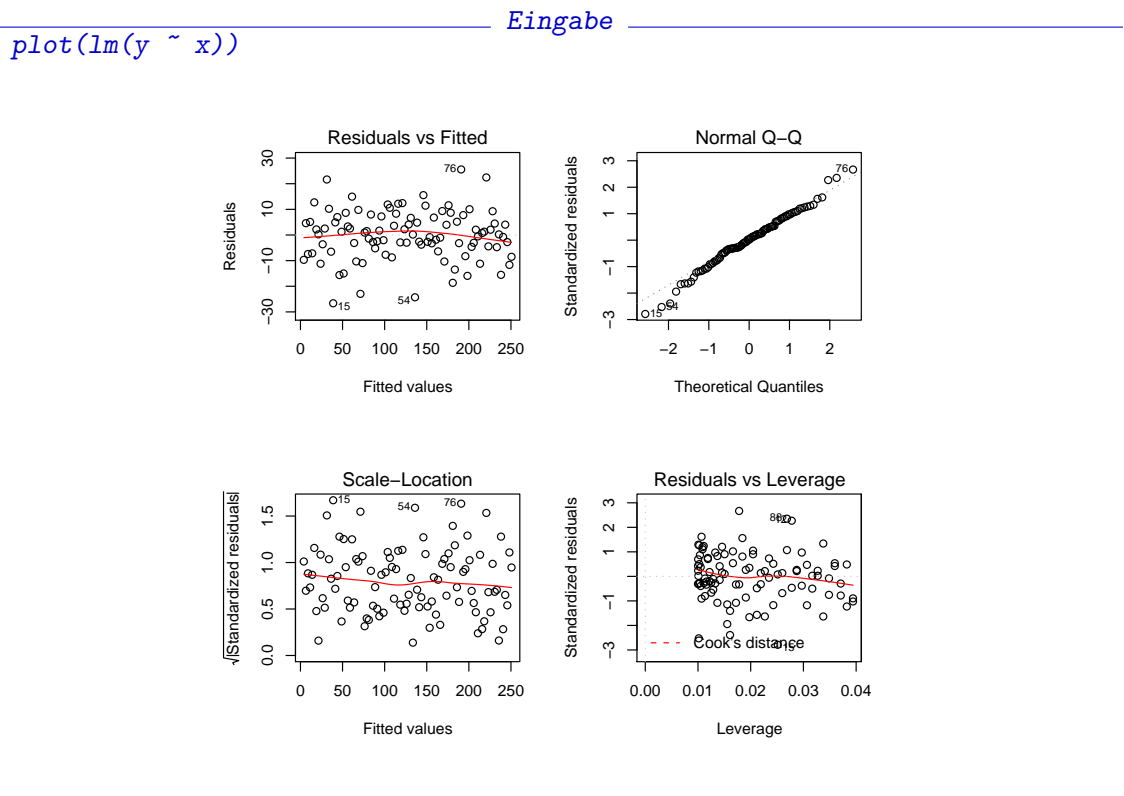
In Abschnitt 2.3 werden wir den theoretischen Hintergrund bereitstellen, der uns hilft, die noch offenen Terme zu interpretieren.

Der Aufruf der Funktion `lm()` liefert immer ein Resultat - wenn es den Daten angemessen ist, aber es gibt auch ein Resultat, wenn das lineare Modell gar nicht angemessen ist. Wir brauchen deshalb eine Diagnostik, die uns hilft, zu erkennen, ob das Modell verlässlich und brauchbar ist.

Aufgabe 2.3	
	<p>Sei $yy < -2.5 * x + 0.01x^2 + err$. Welches Resultat erhalten Sie, wenn Sie eine Regression mit dem (falschen) Modell <code>yy ~ x</code> rechnen? Gibt es Hinweise darauf, dass dieses Modell nicht angemessen ist?</p>

Die Funktion `lm()` führt nicht nur die Schätzung im linearen Modell durch, sondern liefert eine ganze Reihe von Diagnostiken, die helfen können zu beurteilen, ob die Modellvoraussetzungen vertretbar erscheinen. Eine Darstellung mit `plot()` zeigt vier Aspekte davon.

Beispiel 2.4:



Der obere linke Plot zeigt die Residuen gegen den Fit. Die Verteilung der gefitteten Werte hängt vom Design ab.

Die Residuen sollten annähernd wie ein Scatterplot von unabhängigen Variablen aussehen. Die Verteilung der Residuen sollte nicht vom Fit abhängen. Sind hier systematische Strukturen zu erkennen, so ist das ein Warnzeichen dass das Modell oder die Modellvoraussetzungen nicht erfüllt sind.

Nach der vorausgegangen Diskussion können wir noch genauer sein: die Residuen sollten nach (2.13) Linearkombinationen von unabhängig identisch verteilten Variablen sein. Falls die Modellvoraussetzungen erfüllt sind, ist die Varianz durch (2.14) beschrieben.

In der eindimensionalen Situation würde ein Plot der Residuen gegen den Regressor ausreichen. Für p Regressoren wird die graphische Darstellung problematisch. Der Plot der Residuen gegen den Fit verallgemeinert sich auch auf höhere Dimensionen.

Verteilungsaussagen über die Schätzer und Residuen können wir machen, wenn wir mit Verteilungsaussagen über die Fehlerterme beginnen. Die kräftigsten Aussagen sind möglich, wenn die Fehlerterme unabhängig identisch normalverteilt sind. Der obere rechte Plot sollte annähernd wie der “normal probability plot” von normalverteilten Variablen aussehen, wobei das “annähernd” wiederum bedeutet: bis auf Transformation mit der Matrix $I - H$.

Die beiden übrigen Plots sind spezielle Diagnostiken für lineare Modelle (siehe `help(plot.lm)`).

Aufgabe 2.4	
	Inspizieren Sie das Resultat von Aufgabe 2.3 grafisch. Welche Hinweise gibt es jetzt, dass das lineare Modell nicht angemessen ist?

`plot()` stellt für lineare Modelle noch weitere diagnostische Plots bereit. Diese müssen explizit mit dem Parameter `which` angefordert werden.

help(lm)

lm

Fitting Linear Models

Description.

`lm` is used to fit linear models. It can be used to carry out regression, single stratum analysis of variance and analysis of covariance (although `aov` may provide a more convenient interface for these).

Usage.

```
lm(formula, data, subset, weights, na.action,
   method = "qr", model = TRUE, x = FALSE, y = FALSE, qr = TRUE,
   singular.ok = TRUE, contrasts = NULL, offset, ...)
```

Arguments.

<code>formula</code>	a symbolic description of the model to be fit. The details of model specification are given below.
<code>data</code>	an optional data frame, list or environment (or object coercible by <code>as.data.frame</code> to a data frame) containing the variables in the model. If not found in <code>data</code> , the variables are taken from <code>environment(formula)</code> , typically the environment from which <code>lm</code> is called.
<code>subset</code>	an optional vector specifying a subset of observations to be used in the fitting process.
<code>weights</code>	an optional vector of weights to be used in the fitting process. Should be <code>NULL</code> or a numeric vector. If non- <code>NULL</code> , weighted least squares is used with weights <code>weights</code> (that is, minimizing $\sum(w \cdot e^2)$); otherwise ordinary least squares is used.
<code>na.action</code>	a function which indicates what should happen when the data contain NAs. The default is set by the <code>na.action</code> setting of <code>options</code> , and is <code>na.fail</code> if that is unset. The “factory-fresh” default is <code>na.omit</code> . Another possible value is <code>NULL</code> , no action. Value <code>na.exclude</code> can be useful.
<code>method</code>	the method to be used; for fitting, currently only <code>method = "qr"</code> is supported; <code>method = "model.frame"</code> returns the model frame (the same as with <code>model = TRUE</code> , see below).
<code>model, x, y, qr</code>	logicals. If <code>TRUE</code> the corresponding components of the fit (the model frame, the model matrix, the response, the QR decomposition) are returned.

`singular.ok` logical. If `FALSE` (the default in S but not in R) a singular fit is an error.

`contrasts` an optional list. See the `contrasts.arg` of `model.matrix.default`.

`offset` this can be used to specify an *a priori* known component to be included in the linear predictor during fitting. This should be `NULL` or a numeric vector of length either one or equal to the number of cases. One or more `offset` terms can be included in the formula instead or as well, and if both are specified their sum is used. See `model.offset`.

... additional arguments to be passed to the low level regression fitting functions (see below).

Details.

Models for `lm` are specified symbolically. A typical model has the form `response ~ terms` where `response` is the (numeric) response vector and `terms` is a series of terms which specifies a linear predictor for `response`. A terms specification of the form `first + second` indicates all the terms in `first` together with all the terms in `second` with duplicates removed. A specification of the form `first:second` indicates the set of terms obtained by taking the interactions of all terms in `first` with all terms in `second`. The specification `first*second` indicates the *cross* of `first` and `second`. This is the same as `first + second + first:second`.

If the formula includes an `offset`, this is evaluated and subtracted from the response.

If `response` is a matrix a linear model is fitted separately by least-squares to each column of the matrix.

See `model.matrix` for some further details. The terms in the formula will be re-ordered so that main effects come first, followed by the interactions, all second-order, all third-order and so on: to avoid this pass a `terms` object as the formula (see `aov` and `demo(glm.vr)` for an example).

A formula has an implied intercept term. To remove this use either `y ~ x - 1` or `y ~ 0 + x`. See `formula` for more details of allowed formulae.

`lm` calls the lower level functions `lm.fit`, etc, see below, for the actual numerical computations. For programming only, you may consider doing likewise.

All of `weights`, `subset` and `offset` are evaluated in the same way as variables in `formula`, that is first in `data` and then in the environment of `formula`.

Value.

`lm` returns an object of class `"lm"` or for multiple responses of class `c("mlm", "lm")`.

The functions `summary` and `anova` are used to obtain and print a summary and analysis of variance table of the results. The generic accessor functions `coefficients`, `effects`, `fitted.values` and `residuals` extract various useful features of the value returned by `lm`.

An object of class `"lm"` is a list containing at least the following components:

`coefficients` a named vector of coefficients

`residuals` the residuals, that is response minus fitted values.

`fitted.values` the fitted mean values.

`rank` the numeric rank of the fitted linear model.

`weights` (only for weighted fits) the specified weights.

`df.residual` the residual degrees of freedom.

`call` the matched call.

`terms` the `terms` object used.

`contrasts` (only where relevant) the contrasts used.
`xlevels` (only where relevant) a record of the levels of the factors used in fitting.
`offset` the offset used (missing if none were used).
`y` if requested, the response used.
`x` if requested, the model matrix used.
`model` if requested (the default), the model frame used.

In addition, non-null fits will have components `assign`, `effects` and (unless not requested) `qr` relating to the linear fit, for use by extractor functions such as `summary` and `effects`.

Using time series.

Considerable care is needed when using `lm` with time series.

Unless `na.action = NULL`, the time series attributes are stripped from the variables before the regression is done. (This is necessary as omitting NAs would invalidate the time series attributes, and if NAs are omitted in the middle of the series the result would no longer be a regular time series.)

Even if the time series attributes are retained, they are not used to line up series, so that the time shift of a lagged or differenced regressor would be ignored. It is good practice to prepare a `data` argument by `ts.intersect(..., dframe = TRUE)`, then apply a suitable `na.action` to that data frame and call `lm` with `na.action = NULL` so that residuals and fitted values are time series.

Note.

Offsets specified by `offset` will not be included in predictions by `predict.lm`, whereas those specified by an offset term in the formula will be.

Author(s).

The design was inspired by the S function of the same name described in Chambers (1992). The implementation of model formula by Ross Ihaka was based on Wilkinson & Rogers (1973).

References.

Chambers, J. M. (1992) *Linear models*. Chapter 4 of *Statistical Models in S* eds J. M. Chambers and T. J. Hastie, Wadsworth & Brooks/Cole.

Wilkinson, G. N. and Rogers, C. E. (1973) Symbolic descriptions of factorial models for analysis of variance. *Applied Statistics*, **22**, 392–9.

See Also.

`summary.lm` for summaries and `anova.lm` for the ANOVA table; `aov` for a different interface.

The generic functions `coef`, `effects`, `residuals`, `fitted`, `vcov`.

`predict.lm` (via `predict`) for prediction, including confidence and prediction intervals; `confint` for confidence intervals of *parameters*.

`lm.influence` for regression diagnostics, and `glm` for **generalized linear models**.

The underlying low level functions, `lm.fit` for plain, and `lm.wfit` for weighted regression fitting.

Examples.

```
## Annette Dobson (1990) "An Introduction to Generalized Linear Models".
## Page 9: Plant Weight Data.
ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
group <- gl(2,10,20, labels=c("Ctl","Trt"))
weight <- c(ctl, trt)
anova(lm.D9 <- lm(weight ~ group))
summary(lm.D90 <- lm(weight ~ group - 1))# omitting intercept
summary(resid(lm.D9) - resid(lm.D90)) #- residuals almost identical

opar <- par(mfrow = c(2,2), oma = c(0, 0, 1.1, 0))
plot(lm.D9, las = 1)      # Residuals, Fitted, ...
par(opar)

## model frame :
stopifnot(identical(lm(weight ~ group, method = "model.frame"),
                    model.frame(lm.D9)))
```

Nicht erwähnt in der Help-Information: mit *first-second* werden Terme der ersten Gruppe ins Modell aufgenommen, die der zweiten Gruppe aber ausgeschlossen. Ausführlichere Information zur Formel-Darstellung erhält man mit *help(formula)*. Eine Zusammenfassung ist im Anhang A.53 (Seite A-35) zu finden.

Die Hut-Matrix ist eine Besonderheit linearer Modelle. Fit und Residuum jedoch sind allgemeine Konzepte, die bei allen Arten der Schätzung angewandt werden können. Die Anwender sind oft mit dem Fit (oder der Schätzung) zufrieden. Für den ernsthaften Anwender und für den Statistiker sind die Residuen oft wichtiger: sie weisen darauf hin, was vom Modell oder der Schätzung noch nicht erfasst ist.

2.2.3. Weitere Beispiele für lineare Modelle.

Die Matrix X heißt die *Design-Matrix* des Modells. Sie kann die Matrix sein, die mit den ursprünglichen Messbedingungen x_i als Zeilenvektoren gebildet wird. Aber sie ist nicht auf diesen Spezialfall beschränkt. Unter der scheinbar so einfachen Modellklasse der linearen Modelle lassen sich viele wichtige Spezialfälle einordnen. Ein paar davon sind im folgenden zusammengestellt.

Einfache lineare Regression:

$$y_i = a + b x_i + \varepsilon_i \quad \text{mit } x_i \in \mathbb{R}, a, b \in \mathbb{R}$$

kann als lineares Modell mit

$$X = (1 \ x)$$

geschrieben werden, wobei $1 = (1, \dots, 1)^\top \in \mathbb{R}^n$.

Polynomiale Regression:

$$y_i = a + b_1 x_i + b_2 x_i^2 + \dots + b_k x_i^k + \varepsilon_i \quad \text{mit } x_i \in \mathbb{R}, a, b_j \in \mathbb{R}$$

kann als lineares Modell mit

$$X = (1 \ x \ x^2 \ \dots \ x^k)$$

geschrieben werden, wobei $x^j = (x_1^j \ \dots \ x_n^j)^\top$.

Analog für eine Vielzahl von Modellen, die durch andere Transformationen erreicht werden können.

Varianzanalyse: Einweg-Layout

Gemessen wird unter m Versuchsbedingungen, dabei n_j Messungen unter Versuchsbedingung $j, j = 1, \dots, m$. Die Messung setze sich additiv zusammen aus einem Grundeffekt μ , einem für die Bedingung j spezifischen Beitrag α_j , und einem Messfehler nach

$$y_{ij} = \mu + \alpha_j + \varepsilon_{ij} \quad \text{mit } \mu, \alpha_j \in \mathbb{R}, j = 1, \dots, m.$$

Mit $n = \sum n_j$ und

$$X = (1 \ I_1 \ \dots \ I_m),$$

wobei I_j die Vektoren der Indikatorvariablen für die Zugehörigkeit zur Versuchsgruppe j sind, lässt sich dies als lineares Modell schreiben.³

Covarianzanalyse

Analog zur Varianzanalyse werden Unterschiede zwischen Gruppen untersucht, aber zusätzliche (linear eingehende) Einflussfaktoren werden korrigierend berücksichtigt. Unter Versuchsbedingung j bei Beobachtung i hängt die Messung zusätzlich von Einflussfaktoren x_{ij} der Versuchseinheit ij ab.

$$y_{ij} = \mu + \alpha_j + b x_{ij} + \varepsilon_{ij} \quad \text{mit } \mu, \alpha_j \in \mathbb{R}.$$

2.2.4. Modellformeln. R erlaubt es, Modelle auch dadurch zu spezifizieren, dass die Regeln angegeben werden, nach denen die Design-Matrix gebildet wird. Die Syntax, nach denen die Regeln notiert werden, ist sehr kurz in der Beschreibung von `lm()` angegeben. Wir diskutieren sie jetzt etwas ausführlicher. Diese Modell-Spezifikation ist auch für allgemeinere, nicht lineare Modelle möglich. Die Modell-Spezifikationen werden als Attribut mit dem Namen "formula" gespeichert. Sie können mit `formula()` manipuliert werden.

Beispiele

<code>y ~ 1 + x</code>	entspricht $y_i = (1 \ x_i)(\beta_1 \ \beta_2)^\top + \varepsilon$
<code>y ~ x</code>	Kurzschreibweise für <code>y ~ 1 + x</code> (ein konstanter Term wird implizit angenommen)
<code>y ~ 0 + x</code>	entspricht $y_i = x_i \beta + \varepsilon$
<code>log(y) ~ x1 + x2</code>	entspricht $\log(y_i) = (1 \ x_{i1} \ x_{i2})(\beta_1 \ \beta_2 \ \beta_3)^\top + \varepsilon$ (ein konstanter Term wird implizit angenommen)
<code>lm(y ~ poly(x, 4), data = Experiment)</code>	analysiert den Datensatz "Experiment" mit einem linearen Modell für polynomiale Regression vom Grade 4 in x .

Wichtige Spezialfälle für faktorielle Designs sind:

³Es ist Konvention, dass bei Varianzanalysen der letzte Index die Beobachtung zählt, und Indizes in alphabetischer Folge vergeben werden. Konventionell werden also im Vergleich zu unserer Notation die Rollen von i und j vertauscht.

$y \sim A$	Einweg-Varianzanalyse mit Faktor A ,
$y \sim A + x$	Covarianzanalyse mit Faktor A und Regressions-Covariable x ,
$y \sim A + B$	Zwei-Faktor-Kreuz-Layout mit Faktoren A und B ohne Interaktion,
$y \sim A * B$	Zwei-Faktor-Kreuz-Layout mit Faktoren A und B und allen Interaktionen (Kombinationen der Stufen von A und B),
$y \sim A/B$	Zwei-Faktor hierarchisches Layout mit Faktor A und Subfaktor B .

Eine Übersicht über alle Operatoren zur Modellspezifikation ist im Anhang A.53 (Seite A-35) zu finden.

Aufgabe 2.5	
	Schreiben Sie die vier oben in Abschnitt 2.2.3 genannten Modelle als R-Modellformeln.
	Erzeugen Sie sich für jedes dieser Modelle ein Beispiel durch Simulation und wenden Sie <code>lm()</code> auf diese Beispiele an. Vergleichen Sie die durch <code>lm()</code> geschätzten Parameter mit den Parametern, die Sie in der Simulation benutzt haben.

Die Modellformel wird in einem Eintrag im Resultat von `lm()` gespeichert. Sie kann also aus dem Resultat zurück gewonnen werden. Anhand der Formel-Notation generiert R implizit eine Design-Matrix. Mit `model.matrix()` kann diese Design-Matrix inspiziert werden.

Aufgabe 2.6	
	Generieren Sie drei Vektoren mit je 10 $N(\mu_j, 1)$ -verteilten Zufallsvariablen $\mu_j = j$, $j = 1, 3, 9$. Verketteten Sie diese zu einem Vektor y .
	Generieren Sie sich einen Vektor x aus je 10 wiederholten Werten j , $j = 1, 3, 9$.
	Berechnen Sie die Gauß-Markov-Schätzer in den linearen Modellen $y \sim x$ und $y \sim \text{factor}(x)$.
	Lassen Sie sich das Resultat jeweils als Tabelle mit <code>summary()</code> und als Grafik mit <code>plot()</code> anzeigen und vergleichen Sie die Resultate.

2.2.5. Gauß-Markov-Schätzer und Residuen. Wir werfen nun einen genaueren Blick auf den Gauß-Markov-Schätzer. Kenntnisse aus der linearen Algebra, langes Nachdenken oder andere Quellen sagen uns:

BEMERKUNG 2.3.

- (1) Die Design-Matrix X definiert eine Abbildung $\mathbb{R}^p \rightarrow \mathbb{R}^n$ mit $\beta \mapsto X\beta$.
Der Bild-Raum dieser Abbildung sei \mathcal{M}_X , $\mathcal{M}_X \subset \mathbb{R}^n$. \mathcal{M}_X ist der von den Spaltenvektoren von X aufgespannte Vektorraum.
- (2) Sind die Modell-Annahmen erfüllt, so ist $E(Y) \in \mathcal{M}_X$.
- (3) $\hat{Y} = \pi_{\mathcal{M}_X}(Y)$, wobei $\pi_{\mathcal{M}_X} : \mathbb{R}^n \rightarrow \mathcal{M}_X$ die (euklidische) Orthogonalprojektion ist.
- (4) $\hat{\beta} = \arg \min_{\beta} |Y - \hat{Y}_{\beta}|^2$ wobei $\hat{Y}_{\beta} = X\beta$.

Die Charakterisierung (3) des Gauß-Markov-Schätzers als Orthogonalprojektion hilft für das Verständnis oft weiter: der Fit ist die Orthogonalprojektion des Beobachtungsvektors auf den Erwartungswertraum des Modells (und minimiert damit den quadratischen Abstand). Das Residuum ist das orthogonale Komplement.

In der Statistik ist die Charakterisierung als Orthogonalprojektion auch ein Ausgangspunkt, um den Schätzer systematisch zu analysieren. In einfachen Fällen helfen Kenntnisse aus der Wahrscheinlichkeitstheorie schon weiter, etwa zusammengefasst im folgenden Satz:

THEOREM 2.4. *Sei Z eine Zufallsvariable mit Werten in \mathbb{R}^n , die nach $N(0, \sigma^2 I_{n \times n})$ verteilt ist und sei $\mathbb{R}^n = L_0 \oplus \dots \oplus L_r$ eine Orthogonalzerlegung. Sei $\pi_i = \pi_{L_i}$ die Orthogonalprojektion auf L_i , $i = 0, \dots, r$.*

Dann gilt

- (i) $\pi_0(Z), \dots, \pi_r(Z)$ sind unabhängige Zufallsvariablen.
- (ii) $\frac{|\pi_i(Z)|^2}{\sigma^2} \sim \chi^2(\dim L_i)$ für $i = 0, \dots, r$.

BEWEIS. \rightarrow Wahrscheinlichkeitstheorie. Siehe z.B. [Jørgensen 1993, 2.5 Theorem 3]. \square

Mit $\varepsilon = Y - X\beta$ können daraus theoretische Verteilungsaussagen für Schätzer $\hat{\beta}$ und Residuen $Y - \hat{Y}$ abgeleitet werden.

Insbesondere erhalten wir für einfache lineare Modelle aus der residuellen Varianz auch einen Schätzer für die Varianz (bzw. Standardabweichung) jeder einzelnen Komponente $\hat{\beta}_k$. Die entsprechende t -Statistik und der p -Wert für den Test der Hypothese $\hat{\beta}_k = 0$ sind in der Ausgabe von `summary()` angegeben.

Aufgabe 2.7	
	Welche Verteilung hat $ R_X(Y) ^2 = Y - \hat{Y} ^2$, wenn ε nach $N(0, \sigma^2 I)$ verteilt ist?

Auf den ersten Blick ist $|R_X(Y)|^2 = |Y - \hat{Y}|^2$ ein geeignetes Maß, um die Qualität eines Modells zu beurteilen: kleine Werte sprechen für den Fit, große Werte zeigen, dass der Fit schlecht ist. Dies ist jedoch mit Sorgfalt zu betrachten. Zum einen hängt diese Größe von linearen Skalenfaktoren ab. Zum anderen muss die Dimensionen der jeweiligen Räume mit in Betracht gezogen werden. Was passiert, wenn weitere Regressoren ins Modell aufgenommen werden? Wir haben z.B. gesehen, dass "linear" auch die Möglichkeit gibt, nichtlineare

Beziehungen zu modellieren, zum Beispiel dadurch, dass geeignet transformierte Variable in die Design-Matrix mit aufgenommen werden. Die Charakterisierung (3) aus Bemerkung 2.3 sagt uns, dass effektiv nur der von der Design-Matrix aufgespannte Raum relevant ist. Hier sind die Grenzen des Gauß-Markov-Schätzers im linearen Modell erkennbar: wenn viele transformierte Variablen aufgenommen werden, oder generell wenn der durch die Design-Matrix bestimmte Bildraum zu groß wird, gibt es eine Überanpassung. Im Extrem ist $\hat{Y} = Y$. Damit werden alle Residuen zu null, aber die Schätzung ist nicht brauchbar.

Wir benutzen $|R_X(Y)|^2 / \dim(L_X)$, wobei L_X das orthogonale Komplement von \mathcal{M}_X in \mathbb{R}^n ist (also $\dim(L_X) = n - \dim(\mathcal{M}_X)$), um die Dimensionsabhängigkeit zu kompensieren.

Aufgabe 2.8	
	Modifizieren Sie die Plot-Ausgabe <code>plot.lm()</code> für die linearen Modelle so, dass anstelle des Tukey-Anscombe-Plots die studentisierten Residuen gegen den Fit aufgetragen werden.
*	Ergänzen Sie den <i>QQ</i> -Plot durch Monte-Carlo-Bänder für unabhängige Gauß'sche Fehler. <i>Hinweis:</i> Sie können die Bänder nicht direkt aus der Gaußverteilung generieren - Sie brauchen die Residuenverteilung, nicht die Fehlerverteilung.

Aufgabe 2.9	
	Schreiben Sie eine Prozedur, die für die einfache lineare Regression $y_i = a + bx_i + \varepsilon_i$ mit $x_i \in \mathbb{R}, a, b \in \mathbb{R}$ den Gauß-Markov-Schätzer berechnet und vier Plots darstellt: <ul style="list-style-type: none"> • Respons gegen Regressor, mit geschätzter Geraden • studentisierte Residuen gegen Fit • Verteilungsfunktion der studentisierten Residuen im <i>QQ</i>-Plot mit Bändern • Histogramm der studentisierten Residuen

2.3. Streuungszerlegung und Varianzanalyse

Wenn ein einfaches lineares Modell mit gaußverteilten Fehlern vorliegt, sind die *t*-Tests geeignet, eindimensionale Probleme (Tests oder Konfidenzintervalle für einzelne Parameter, punktweise Konfidenzintervalle) zu lösen. Um simultane oder mehrdimensionale Probleme zu lösen brauchen wir andere Werkzeuge. Anstelle der Differenzen oder Mittelwerte, die den *t*-Tests zu Grunde liegen, benutzen wir Norm-Abstände (bzw. quadratische Abstände), die auch auf höhere Dimensionen generalisieren.

Die Interpretation des Gauß-Markov-Schätzers als Orthogonalprojektion (Bem. 2.3 3) zeigt eine Möglichkeit, Modelle zu vergleichen: Für X, X' Design-Matrizen mit $\mathcal{M}_{X'} \subset \mathcal{M}_X$, betrachten wir die Zerlegung $\mathbb{R}^n = L_0 \oplus \dots \oplus L_r$ mit $L_0 := \mathcal{M}_{X'}$, und die orthogonale

Komplemente $L_1 := \mathcal{M}_X \ominus \mathcal{M}_{X'}$, $L_2 := \mathbb{R}^n \ominus \mathcal{M}_X$. Wieder bezeichnet π jeweils die entsprechende Projektion.

$$F := \frac{\frac{1}{\dim(L_1)} |\pi_{\mathcal{M}_X} Y - \pi_{\mathcal{M}_{X'}} Y|^2}{\frac{1}{\dim(L_2)} |Y - \pi_{\mathcal{M}_X} Y|^2}.$$

Diese Statistik, die F -Statistik (nach R.A. Fisher) ist die Basis für die **Varianzanalyse**, einer klassischen Strategie, Modelle zu vergleichen. **Streuungszerlegung** ist ein anderer Name für diesen Ansatz.

Die Idee wird auf Ketten von Modellen verallgemeinert. Ist $\mathcal{M}_0 \subset \dots \subset \mathcal{M}_r = \mathbb{R}^n$, so liefert $L_0 := \mathcal{M}_0$, $L_i := \mathcal{M}_i \ominus \mathcal{M}_{i-1}$ für $i = 1, \dots, r$ eine Orthogonalzerlegung. Mit den Bezeichnungen von oben ist dann

$$\frac{\frac{1}{\dim L_{i-1}} |\pi_{\mathcal{M}_i} Y - \pi_{\mathcal{M}_{i-1}} Y|^2}{\frac{1}{\dim L_i} |Y - \pi_{\mathcal{M}_i} Y|^2}$$

eine Teststatistik, die zum Test für das Modell \mathcal{M}_{i-1} im Vergleich zum Obermodell \mathcal{M}_i herangezogen wird.

Aufgabe 2.10	
	Welche Verteilung hat F , wenn $E(Y) \in \mathcal{M}_{X'}$ gilt und ε nach $N(0, \sigma^2 I)$ verteilt ist?

Aufgabe 2.11	
	Geben Sie eine explizite Formel für die F -Statistik zur Varianzanalyse im Einweg-Layout $y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$ im Vergleich zum homogenen Modell $y_{ij} = \mu + \varepsilon_{ij}.$

Die Varianzanalyse gibt eine andere Darstellung und Interpretation der linearen Modelle. Hier im Vergleich zu Beispiel 2.3 die Varianzanalyse-Darstellung:

Beispiel 2.5:	
Eingabe	
<code>summary(aov(lmres))</code>	
Ausgabe	
	Df Sum Sq Mean Sq F value Pr(>F)
x	1 517386 517386 5532.4 < 2.2e-16 ***
Residuals	98 9165 94

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1	

Aufgabe 2.12	
	Analysieren Sie die in Beispiel 2.3 (Seite 2-8) gezeigten Ausgaben von <code>lm()</code> . Welche Terme können Sie jetzt interpretieren? Stellen Sie diese Interpretationen schriftlich zusammen. Für welche Terme fehlt Ihnen noch Information?

In der Ausgabe finden Sie noch einen Hinweis auf “R-squared”. Der Term, der hier angegeben wird, ist ein Schätzer für den Anteil von $Var(Y)$, der durch das Modell erklärt wird:

$$R^2 = \frac{mss}{mss + rss}$$

mit $mss := \frac{1}{n} \sum (\hat{Y}_i - \bar{Y})^2$ und $rss := \frac{1}{n} \sum (R_X(Y)_i - \overline{R_X(Y)})^2$. Die Bezeichnung R^2 kommt von der einfachen linearen Regression. Dort ist konventionell die Korrelation $Cor(X, Y)$ mit R bezeichnet, und $R^2 = Cor(X, Y)^2$. R^2 berücksichtigt nicht die Anzahl der geschätzten Parameter und kann deshalb zu optimistisch sein. Der Term “adjusted R-squared” hat eine Gewichtung, die die Freiheitsgrade berücksichtigt.

[help\(anova\)](#)

anova

Anova Tables

Description.

Compute analysis of variance (or deviance) tables for one or more fitted model objects.

Usage.

`anova(object, ...)`

Arguments.

object an object containing the results returned by a model fitting function (e.g., `lm` or `glm`).
... additional objects of the same type.

Value.

This (generic) function returns an object of class `anova`. These objects represent analysis-of-variance and analysis-of-deviance tables. When given a single argument it produces a table which tests whether the model terms are significant.

When given a sequence of objects, `anova` tests the models against one another in the order specified.

The print method for `anova` objects prints tables in a “pretty” form.

Warning.

The comparison between two or more models will only be valid if they are fitted to the same dataset. This may be a problem if there are missing values and R’s default of `na.action = na.omit` is used.

References.

Chambers, J. M. and Hastie, T. J. (1992) *Statistical Models in S*, Wadsworth & Brooks/Cole.

See Also.

`coefficients, effects, fitted.values, residuals, summary, drop1, add1.`

Modelle für die Varianzanalyse können als Regeln angegeben werden. Dieselbe Syntax zur Modellbeschreibung wird benutzt wie schon bei der Regression. Wenn Terme auf der rechten Seite der Modellbeschreibung Faktoren sind, wird automatisch ein Varianzanalyse-Modell anstelle eines Regressionsmodells generiert.

Die Modellbeschreibung bestimmt die linearen Räume, in denen die Erwartungswerte liegen. Die Streuungszersetzungen sind dadurch jedoch nicht eindeutig bestimmt: die Angabe der Räume lässt evtl. noch verschiedene Orthogonalzerlegungen zu (z.B. abhängig von der Reihenfolge). Mehr noch: die Angabe der Faktoren bestimmt ein Erzeugendensystem der Räume. Die Faktoren brauchen nicht orthogonal zu sein, noch nicht einmal unabhängig.

Dies gilt für alle linearen Modelle. In der Regression ist Abhängigkeit eher die Ausnahme. Bei faktoriellen Designs taucht dieser Fall häufig auf. Die Einweg-Varianzanalyse in Koordinatendarstellung illustriert dieses Problem: mit

$$y_{ij} = \mu + \alpha_j + \varepsilon_{ij} \quad \text{mit } \mu, \alpha_j \in \mathbb{R}$$

ist für $n_j > 0$ die Zerlegung in μ und α_j nicht eindeutig. Der tieferliegende Grund ist: der globale Faktor μ definiert den vom Einheitsvektor 1 aufgespannten Raum, und dieser liegt in dem von den Gruppenindikatoren aufgespannten Raum.

Die Modellformel definiert eine Designmatrix X und damit einen Modellraum. Eine zusätzliche Matrix C wird benutzt, um die Matrix zu reduzieren und damit eine eindeutige Streuungszerslegung zu spezifizieren. Die effektive Designmatrix ist dann $[1 \ X \ C]$; C heißt **Kontrastmatrix**. Die Funktionen zur Varianzanalyse wie z.B. `lm()` oder `aov()` erlauben es, die Kontraste zu spezifizieren.

Die Funktion `anova()` operiert wie eine spezielle Formatierung der Ausgabe und wird analog `summary()` benutzt, also z.B. in der Form `anova(lm())`.

Aufgabe 2.13	
	<p>Die Datei "micronuclei" enthält einen Datensatz aus einem Mutagenitätstest. Zellkulturen (je 50 Einheiten) wurden in einer Kontrollgruppe und unter 5 chemischen Behandlungen beobachtet. Der Effekt der Substanzen ist, die Chromosomen aufzubrechen und Mikronuklei zu bilden. Registriert wurde die Größe der Mikronuklei (relativ zum Eltern-Nukleus).</p> <p>Lesen Sie die Datei "micronuclei" und berechnen Sie für jede Gruppe Mittelwert und Varianz.</p> <p>Hinweis: Sie können die Datei mit <code>data()</code> einlesen. Für Dateien mit Tabellenformat gibt es die spezielle Anweisung <code>read.table()</code>. Informieren Sie sich mit <code>help()</code> über beide Funktionen.</p> <p>Einige ausgewählte statistische Funktionen (z.B. Mittelwert) finden Sie in Tabelle A.22 im Anhang.</p>
	<p>Vergleichen Sie die Resultate. Sind Behandlungseffekte nachweisbar? Hinweise: Versuchen Sie zunächst, die Aufgabe als Einweg-Varianzanalyse zu formulieren. Den Datensatz müssen Sie zunächst z.B. mit Hilfe von <code>c()</code> auf eine geeignete Form bringen.</p>

Aufgabe 2.14	
*	Schreiben Sie eine Funktion <code>oneway()</code> , die als Argument eine Datentabelle nimmt und eine Einweg-Varianzanalyse als Test auf die Differenz zwischen den Spalten durchführt.
*	Ergänzen Sie <code>oneway()</code> durch die notwendigen diagnostischen Plots. Welche Diagnostiken sind notwendig?

Aufgabe 2.15	
	Das Industrieunternehmen Kiwi-Hopp ⁴ möchte einen neuen Hubschrauber auf den Markt bringen. Die Hubschrauber müssen also danach beurteilt werden, wie lange sie sich in der Luft halten, bis sie aus einer gegebenen Höhe (ca. 2m) den Boden erreichen ⁵ . Eine Konstruktionszeichnung ist unten (Abbildung 2.1, Seite 2-24) angegeben. Welche Faktoren könnten die Variabilität der Flug(Sink)zeiten beeinflussen? Welche Faktoren könnten die mittlere Flugzeit beeinflussen?
	Führen Sie 30 Versuchsflüge mit einem Prototyp durch und messen Sie die Zeit in 1/100s. (Sie müssen vielleicht zusammenarbeiten, um die Messungen durchzuführen.) Würden Sie die gemessene Zeit als normalverteilt ansehen? Die Anforderung ist, dass die mittlere Flugdauer mindestens 2.4s erreicht. Erfüllt der Prototyp diese Anforderung?
	Sie haben die Aufgabe, einen Entwurf für die Produktion auszusuchen. Folgende Varianten stehen zur Diskussion: Rotorbreite 45mm Rotorbreite 35mm Rotorbreite 45mm mit Zusatzfalte als Stabilisierung Rotorbreite 35mm mit Zusatzfalte als Stabilisierung. Ihr Haushalt erlaubt ca. 40 Testflüge. (Wenn Sie mehr Testflüge benötigen, müssen Sie dies gut begründen.) Bauen Sie 4 Prototypen und führen sie Testflüge durch, bei denen Sie die Zeit messen. Finden Sie diejenige Konstruktion, die die längste Flugdauer ergibt. Erstellen Sie einen Bericht. Der Bericht sollte folgende Details enthalten: <ul style="list-style-type: none"> • eine Liste der erhobenen Daten und eine Beschreibung des experimentellen Vorgehens. • geeignete Plots für jede Konstruktion • eine Varianzanalyse • eine klare Zusammenfassung Ihrer Schlüsse. <p style="text-align: right;">(Fortsetzung)→</p>

⁴Nach einer Idee von Alan Lee, Univ. Auckland, Neuseeland

⁵Kiwis können nicht fliegen.

Aufgabe 2.15	(Fortsetzung)
	<i>Weitere Hinweise:</i> Randomisieren Sie die Reihenfolge Ihrer Experimente. Reduzieren Sie die Variation, indem Sie gleichmässige Bedingungen für das Experiment schaffen (gleiche Höhe, gleiche Abwurftechnik etc.).
	Die Zusatzfaltung verursacht zusätzliche Arbeitskosten. Schätzen sie den Effekt ab, den diese Zusatzinvestition bringt.

Aufgabe 2.16	
	Benutzen Sie den Quantil-Quantil-Plot, um paarweise die Resultate des Helikopter-Experiments aus dem letzten Kapitel zu vergleichen. Formulieren Sie die Resultate.

Aufgabe 2.17	
	Inspizieren Sie die Implementierung von <code>qqnorm()</code> . Programmieren Sie eine analoge Funktion für den <i>PP</i> -Plot und wenden Sie diese auf die Helikopter-Daten an.

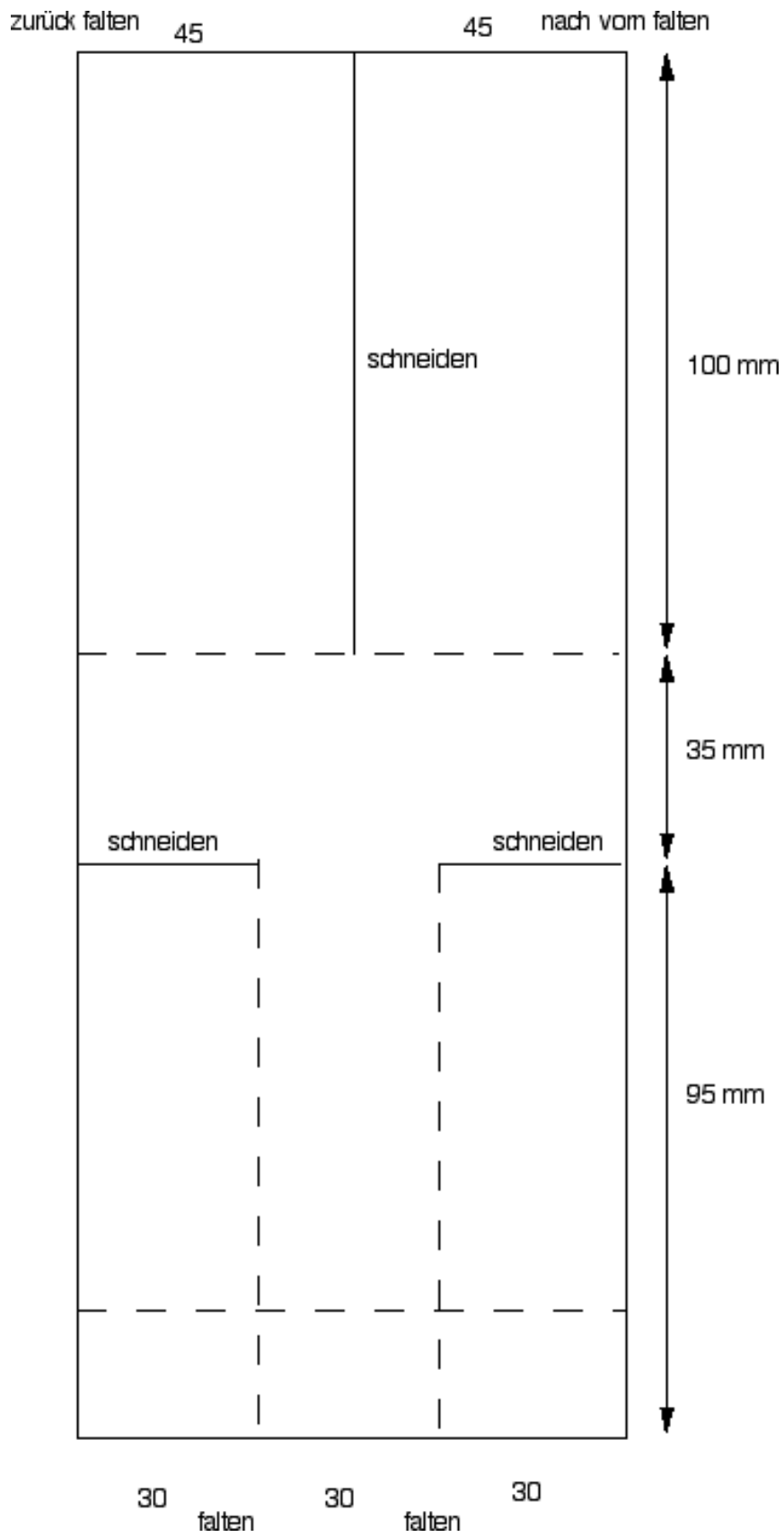


ABBILDUNG 2.1. KiwiHopp

2.4. Simultane Schätzung

2.4.1. Scheffé's Konfidenz-Bänder. Der Kleinste-Quadrate-Schätzer schätzt im Prinzip alle Komponenten des Parameter-Vektors simultan. Die Optimalitäts-Aussagen des Gauß-Markov-Theorems beziehen sich nur auf eindimensionale lineare Statistiken. Die Konfidenzaussagen gelten jedoch multivariat. Es gilt: Der mithilfe der F -Verteilung gewonnene Konfidenzbereich zum Konfidenzniveau $1 - \alpha$ hat die Form

$$\{\widehat{\beta} \in \mathbb{R}^k : (\sum_{j=1}^k (\widehat{\beta}_j - \beta_j)^2 \|x_j\|^2 / k) / \widehat{\sigma}^2 \leq F_{1-\alpha}(k, n - k)\},$$

d.h. der Konfidenzbereich ist eine Ellipse. Wir können die Ellipse auch als den Bereich definieren, der durch alle Tangenten der Ellipse begrenzt wird. Dies übersetzt die (eine) quadratischen Bedingung an die Punkte im Konfidenzbereich durch (unendlich viele) lineare Bedingungen. Diese geometrische Beziehung ist der Kern für den folgenden Satz:

THEOREM 2.5. *Sei $\mathcal{L} \subset \mathbb{R}^k$ ein linearer Unterraum der Dimension d ; $EY = Xb$ mit $rk(X) = p < n$. Dann ist*

$$P\{\ell^t \beta \in \ell^t \widehat{\beta} \pm (dF_{d, n-d}^\alpha)^{1/2} s(\ell^t (X^t X)^{-1} \ell)^{1/2} \forall \ell \in \mathcal{L}\} = (1 - \alpha).$$

BEWEIS. [Mil81, 2.2, p. 48] □

Dies ist ein simultaner Konfidenzbereich für alle Linearkombinationen aus \mathcal{L} . Als Test übersetzt ergibt dies einen simultanen Test für alle linearen Hypothesen aus \mathcal{L} . Im Falle $d = 1$ reduziert sich dieser Scheffé-Test auf den üblichen F -Test. Üblicherweise ist es nicht möglich, am selben Datenmaterial mehrere Tests durchzuführen, ohne dadurch das Konfidenzniveau zu verschlechtern. Der F -Test ist eine Ausnahme. Nach einem globalen F -Test können diese Linearkombinationen oder Kontraste einzeln getestet werden, ohne das Niveau zu verletzen.

Im Falle der einfachen linearen Regression übersetzt sich das Konfidenz-Ellipsoid im Parameterraum so in ein Hyperboloid als Konfidenzbereich für die Regressionsgeraden im Regressor/Respons-Raum.

Geht man zur Interpretation im Regressor/Respons-Raum, also dem Raum der Versuchsbedingungen und Beobachtungen über, so ist man häufig nicht so sehr an einem Konfidenzbereich für die Regressionsgerade interessiert, sondern daran, einen Prognosebereich (Toleranzbereich) für weitere Beobachtungen anzugeben. Für diesen muss zur Streuung der Regression noch die Fehlerstreuung addiert werden. Der Toleranzbereich ist entsprechend größer. Konfidenzbereich für die Regressionsgerade und Toleranzbereich für Beobachtungen können mit der Funktion `predict()` berechnet werden. Die folgende Abbildung zeigt beide Bereiche. Die Funktion `predict()` ist eine generische Funktion. Für lineare Modelle ruft sie `predict.lm()` auf. `predict()` erlaubt es, neue Stützstellen als Parameter `newdata` vorzugeben, an denen anhand des geschätzten Modells ein Fit berechnet wird. Die Variablen werden hier dem Namen nach zugeordnet. Deshalb muss `newdata` ein `data.frame` sein, dessen Komponenten-Namen den ursprünglichen Variablen entsprechen.

Wir bereiten einen Beispieldatensatz vor.

```
n <- 100
sigma <- 1
x <- (1:n)/n-0.5
err <- rnorm(n)
```

Eingabe

```
y <- 2.5 * x + sigma*err
lmxy <- lm(y ~ x)
```

Um bessere Kontrolle über die Grafik zu bekommen, berechnen wir die Plot-Grenzen und Stützpunkte vorab.

Eingabe

```
plotlim <- function(x){
  xlim <- range(x)
  # check implementation of plot. is this needed?
  del <- xlim[2]-xlim[1]
  if (del>0)
    xlim <- xlim+c(-0.1*del, 0.1*del)
  else xlim <- xlim+c(-0.1, 0.1)
  return(xlim)
}
xlim <- plotlim(x)
ylim <- plotlim(y)
#newx <- data.frame(x = seq(1.5*min(x), 1.5*max(x), 1/(2*n)))
newx <- data.frame(x = seq(xlim[1], xlim[2], 1/(2*n)))
```

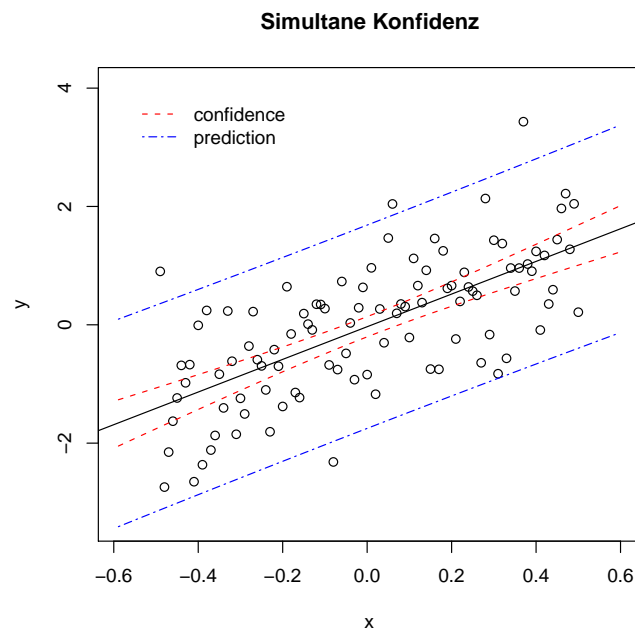
Für diese Daten berechnen wir nun Konfidenzbänder und zeichnen sie.

Beispiel 2.6:

```

plot(x, y, xlim = xlim, ylim = ylim)
abline(lmxy)
pred.w.plim <- predict(lmxy, newdata = newx, interval = "prediction")
pred.w.clim <- predict(lmxy, newdata = newx, interval = "confidence")
matplot(newx$x,
        cbind(pred.w.clim[, -1], pred.w.plim[, -1]),
        lty = c(2, 2, 6, 6),
        col = c(2, 2, 4, 4),
        type = "l", add = TRUE)
title(main = "Simultane Konfidenz")
legend("topleft",
      lty = c(2, 6),
      legend = c("confidence", "prediction"),
      col = c(2, 4),
      inset = 0.05, bty = "n")

```



2.4.2. Tukey's Konfidenz-Intervalle. Geometrisch ist das Konfidenz-Ellipsoid also durch seine (unendlich vielen) Tangentialebenen gekennzeichnet. Übersetzt als Test werden hier unendlich viele lineare Tests simultan durchgeführt. In vielen Anwendungen ist es jedoch möglich, gezieltere Fragestellungen anzugehen, etwa im Zwei-Stichprobenfall nur die Hypothese $\beta_1 - \beta_2 = 0$. Diese reduzierten Fragestellungen können in linearen Modellen formuliert werden und zu schärferen Tests führen. Dies geschieht durch die Spezifizierung von **Kontrasten** und wird in R auch für die Varianzanalyse unterstützt.

2.4.2.1. *Fallstudie: Titrierplatte.* Eine typisches Werkzeug in der Biologie und Medizin sind Tritrierplatten, die z.B. bei Versuchen mit Zellkulturen eingesetzt werden. Die Platte enthält in einem rechteckigen Raster kleine Vertiefungen. Auf die Platte insgesamt können

Substanzen aufgebracht. Mit einer Multipipette können auch spaltenweise oder zeilenweise Substanzen aufgebracht werden (Abbildung 2.2).

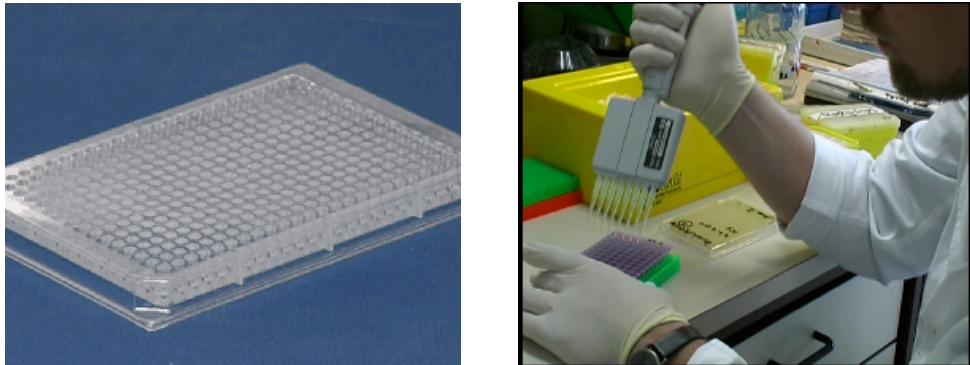


ABBILDUNG 2.2. Titrierplatten. Mit Mulipipetten können zeilenweise oder spaltenweise Substanzen aufgebracht werden.

Die Experimente werden oft in Serien durchgeführt. Aus einer Serie benutzen wir als Beispiel nur die Daten einer Platte.

```
p35 <- read.delim("../data/p35.tab")
```

Für die Analyse mit `lm()` müssen wir die Daten aus der Matrix-Form in eine lange Form überführen, die die Behandlung in einer Spaltenvariablen aufführt. Die Spalte `H` in diesem Versuch enthält keine Behandlung, sondern dient nur zur Qualitätskontrolle zwischen den Platten.

```
s35 <- stack(p35[,3:9]) # ignore column H
s35 <- data.frame(y=s35$values,
  Tmt=s35$ind,
  Lane=rep(1:12, length.out=dim(s35)[1])) # rename
lmres <- lm(y ~ 0 + Tmt, data= s35) # we do not want an overall mean
```

Die Zusammenfassung als lineares Modell enthält Tests für die einzelnen Koeffizienten.

```
summary(lmres)
```

```
Call:
```

```
lm(formula = y ~ 0 + Tmt, data = s35)
```

```
Residuals:
```

Min	1Q	Median	3Q	Max
-0.084833	-0.016354	0.009125	0.022729	0.073083

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
TmtA	0.19383	0.01035	18.73	<2e-16 ***
TmtB	0.24892	0.01035	24.06	<2e-16 ***
TmtC	0.23783	0.01035	22.99	<2e-16 ***
TmtD	0.24117	0.01035	23.31	<2e-16 ***

```
TmtE  0.24392    0.01035    23.57    <2e-16 ***
TmtF  0.23558    0.01035    22.77    <2e-16 ***
TmtG  0.22367    0.01035    21.62    <2e-16 ***
```

```
----
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.03584 on 77 degrees of freedom
Multiple R-squared:  0.9787,    Adjusted R-squared:  0.9768
F-statistic: 506.2 on 7 and 77 DF,  p-value: < 2.2e-16
```

Für dieses Beispiel sind die Tests für die einzelnen Koeffizienten nicht angebracht. `anova()` listet eine Zusammenfassung, die auf die Varianzanalyse zugeschnitten ist.

```
anova(lmres) Eingabe
```

```
Analysis of Variance Table Ausgabe
```

```
Response: y
      Df Sum Sq Mean Sq F value    Pr(>F)
Tmt      7  4.5513   0.6502  506.15 < 2.2e-16 ***
Residuals 77  0.0989   0.0013
```

```
----
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Wenn die Voraussetzungen des Gauß-linearen Modells gegeben sind, so ist der Behandlungseffekt signifikant. Es stellt sich sofort die Frage, zwischen welchen der Behandlungen ein signifikanter Unterschied besteht, d.h. uns interessieren die Kontraste, die die Behandlungsunterschiede beschreiben. Ohne das Niveau zu verletzen können diese post-hoc mit Tukey's Ansatz untersucht werden. Dazu brauchen wir die Funktion `glht()` für den Test für generalisierte lineare Hypothesen, die in der Bibliothek `multcomp` für multiples Testen bereit gestellt ist.

Beispiel 2.7:

```

----- Eingabe -----
library(multcomp)
lhtres<-glht(lmres,linfct=mcp(Tmt="Tukey"))
summary(lhtres)      # multiple tests
----- Ausgabe -----
Simultaneous Tests for General Linear Hypotheses

Multiple Comparisons of Means: Tukey Contrasts

Fit: lm(formula = y ~ 0 + Tmt, data = s35)

Linear Hypotheses:
      Estimate Std. Error t value p value
B - A == 0  0.055083   0.014632   3.765 0.00577 **
C - A == 0  0.044000   0.014632   3.007 0.05277 .
D - A == 0  0.047333   0.014632   3.235 0.02834 *
E - A == 0  0.050083   0.014632   3.423 0.01651 *
F - A == 0  0.041750   0.014632   2.853 0.07793 .
G - A == 0  0.029833   0.014632   2.039 0.39861
C - B == 0 -0.011083   0.014632  -0.757 0.98819
D - B == 0 -0.007750   0.014632  -0.530 0.99832
E - B == 0 -0.005000   0.014632  -0.342 0.99986
F - B == 0 -0.013333   0.014632  -0.911 0.96971
G - B == 0 -0.025250   0.014632  -1.726 0.60126
D - C == 0  0.003333   0.014632   0.228 0.99999
E - C == 0  0.006083   0.014632   0.416 0.99958
F - C == 0 -0.002250   0.014632  -0.154 1.00000
G - C == 0 -0.014167   0.014632  -0.968 0.95930
E - D == 0  0.002750   0.014632   0.188 1.00000
F - D == 0 -0.005583   0.014632  -0.382 0.99974
G - D == 0 -0.017500   0.014632  -1.196 0.89361
F - E == 0 -0.008333   0.014632  -0.570 0.99748
G - E == 0 -0.020250   0.014632  -1.384 0.80861
G - F == 0 -0.011917   0.014632  -0.814 0.98279
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Adjusted p values reported)

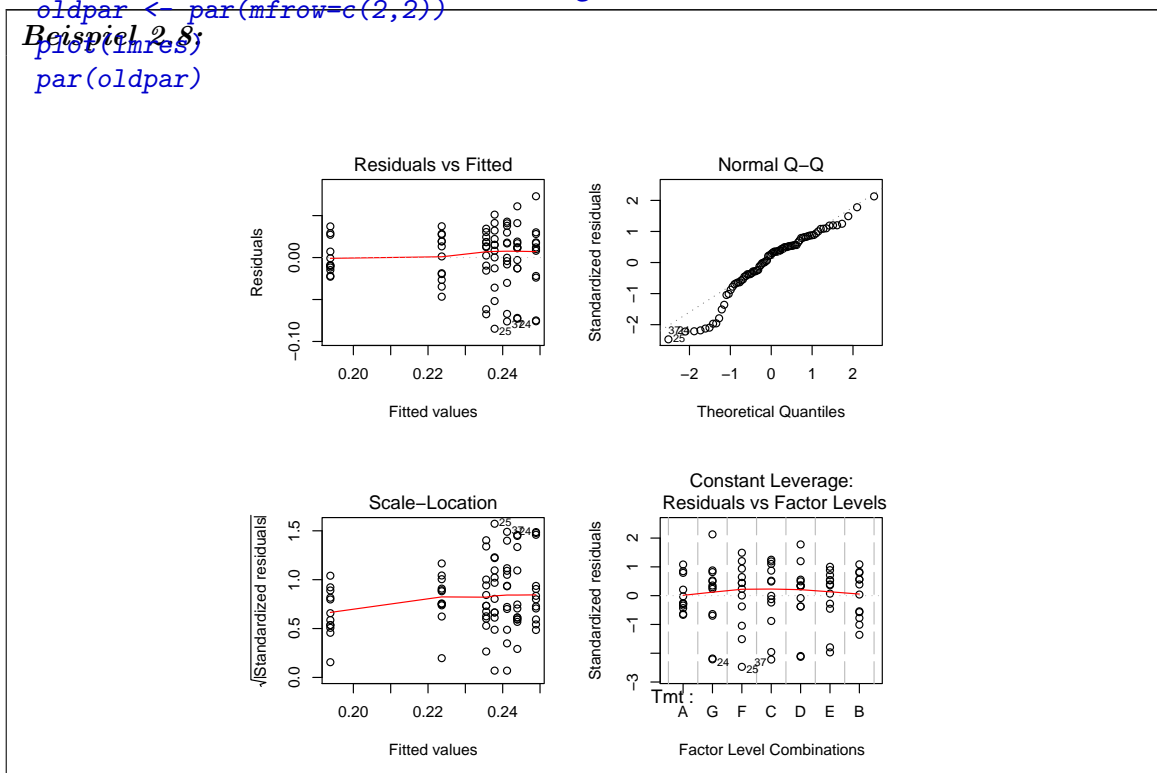
```

Unter den Voraussetzungen des Modells ist damit die Signifikanz der Unterschiede von *A* zu *B* und *D* gesichert.

Um die Voraussetzungen zu prüfen, stehen uns die Residuen zur Verfügung, die wir mit `plot()` inspizieren können.

Eingabe

```
oldpar <- par(mfrow=c(2,2))
lmres <- lm(y ~ x)
plot(lmres)
par(oldpar)
```



Die Verteilung zeigt eine deutliche Abweichung von der Normalverteilung, insbesondere bei kleinen Werten. Wir können diese selektiv inspizieren.

Beispiel 2.9:

Eingabe

```
#diagnostic
library(MASS)
s35$studres <- studres(lmres)
s35[s35$studres < -1,]
```

Ausgabe

	y	Tmt	Lane	studres
13	0.174	B	1	-2.239390
24	0.173	B	12	-2.271296
25	0.153	C	1	-2.559766
33	0.202	C	9	-1.044865
36	0.186	C	12	-1.523409
37	0.165	D	1	-2.279286
48	0.174	D	12	-1.994858
49	0.171	E	1	-2.175828
60	0.172	E	12	-2.144169
61	0.168	F	1	-2.007887
72	0.174	F	12	-1.821449
73	0.177	G	1	-1.367608
84	0.189	G	12	-1.010381

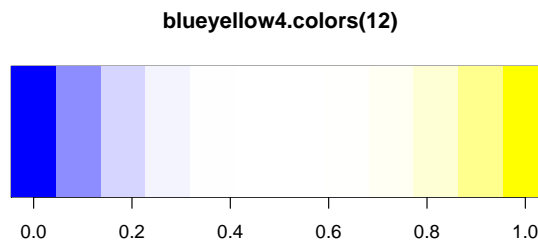
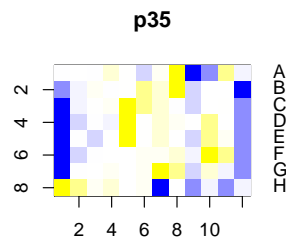
Das Muster ist auffällig. Fast alle besonders kleinen Werte sind am Rand der Platte.

Dieses Muster hätten wir auch mit einer visuellen Inspektion erkennen können:

Eingabe

```
# visualisation
# image, the easy way
a35 <- as.matrix(p35[3:10])
a35rk <- apply(a35,2, rank)
#image(a35rk)

# enhanced image, using bertin
library(bertin)
oldpar<- par(mfrow=c(2,1))
imagem(t(a35rk), col=blueyellow4.colors(12), main="p35")
colramp(blueyellow4.colors(12),12,horizontal=TRUE)
par(oldpar)
```



Bei unabhängigen Fehlern hätten wir eine zufällige Verteilung der Ränge in den Zeilen. Die Konzentrierung der extremen Werte in den extremen Spalten zeigt eine Inhomogenität im Produktionsprozeß.

In diesem Beispiel können wir also berichten, dass anscheinend ein Unterschied zwischen der Behandlung A und speziellen anderen Behandlungen besteht. Die Beurteilung ist aber mit Vorbehalt zu betrachten: die Modellvoraussetzungen sind nicht erfüllt. Es gibt eine erkennbare Inhomogenität zwischen den Zeilen. Wichtiger ist also der Hinweis im Produktionsprozess nach Ursachen dieser Inhomogenität zu suchen.

2.5. Nichtparametrische Regression

2.5.0.2. Transformationen.

2.5.0.3. *Box-Cox-Transformationen.* Lage- und Skalenparameter können auch als Versuch verstanden werden, die Verteilung auf eine Referenzgestalt zu transformieren. Lage- und Skalenparameter erfassen nur lineare Transformationen.

Die **Box -Cox -Transformationen**

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \text{für } \lambda \neq 0, \\ \log(y) & \text{für } \lambda = 0 \end{cases}$$

sind eine Familie, die so skaliert ist, dass die Logarithmus-Transformation stetig in Potenz-Transformationen eingebettet ist. Die Funktion `boxcox()` in `library(MASS)` kann benutzt werden, um λ zu wählen.

Generalisierte lineare Modelle sind so erweitert, dass sie bestimmte Transformationen schon im Modell berücksichtigen können. Dazu finden Sie weitere Information in [VR02].

2.5.1. Zwischenspiel: verallgemeinerte lineare Modelle. Wir wollen schnell zur praktischen Arbeit kommen. An dieser Stelle sollte jedoch eine Ausblick nicht fehlen, wie wir über die einschränkenden Annahmen des linearen Modells hinauskommen. Die linearen Modelle gehören zu den am besten untersuchten Modellen. Theorie und Algorithmen hierfür sind weit entwickelt. Von daher ist es naheliegend, zu probieren, wieweit sich die Modellklasse so erweitern lässt, dass theoretische und algorithmische Erfahrungen noch nutzbar sind.

Wir notierten das lineare Modell als

$$\begin{aligned} Y &= m(X) + \varepsilon \\ Y &\text{ mit Werten in } \mathbb{R}^n \\ X &\in \mathbb{R}^{n \times p} \\ E(\varepsilon) &= 0 \\ \text{mit } m(X) &= X\beta, \quad \beta \in \mathbb{R}^p. \end{aligned}$$

Eine wichtige Erweiterung ist, die Bedingung der Linearität aufzuheben. Sie wird abgemildert mit einer Zwischenstufe. Wir setzen also nicht mehr voraus, dass m linear ist, sondern nur, dass es sich über eine lineare Funktion faktorisieren lässt. Dies ergibt ein verallgemeinertes lineares Modell

$$\begin{aligned} Y &= m(X) + \varepsilon \\ Y &\text{ mit Werten in } \mathbb{R}^n \\ X &\in \mathbb{R}^{n \times p} \\ E(\varepsilon) &= 0 \\ m(X) &= \bar{m}(\eta) \text{ mit } \eta = X\beta, \quad \beta \in \mathbb{R}^p. \end{aligned}$$

Die nächste naheliegende Verallgemeinerung ist, eine Transformation für Y zu berücksichtigen. Zahlreiche weitere Abschwächungen sind diskutiert worden; eine kleine Anzahl hat sich als handhabbar erwiesen. Die verbliebenen Modelle werden als generalisierte lineare Modelle bezeichnet. Generalisierte lineare Modelle haben in R eine weitgehende Unterstützung. In der Regel findet sich zu den hier diskutierten R-Funktionen für lineare Modelle eine Entsprechung für generalisierte lineare Modelle. Weitere Information mit `help(glm)`.

2.5.2. Lokale Regression. Wir machen nun einen großen Sprung. Wir haben lineare Modelle diskutiert. Wir wissen, dass damit auch nichtlineare Funktionen modelliert werden können. Aber die Terme, die in die Funktion eingehen, müssen vorab spezifiziert werden. Zu viele Terme führen zu einer Überanpassung. Die statistische Behandlung von Regressionsproblemen mit geringen Einschränkungen an die Modellfunktion bleibt ein Problem.

Ein partieller Lösungsansatz kommt aus der Analysis. Dort ist es eine Standard-Technik, Funktionen lokal zu approximieren. Das analoge Vorgehen in der Statistik ist, anstelle eines globalen Schätzverfahrens eine lokalisierte Variante zu wählen. Wir nehmen immer noch an, dass

$$\begin{aligned} Y &= m(X) + \varepsilon & Y \in \mathbb{R}^n \\ X &\in \mathbb{R}^{n \times p} \\ E(\varepsilon) &= 0, \end{aligned}$$

aber wir nehmen Linearität nur lokal an:

$$m(x) \approx x' \beta_{x_0} \quad \beta_{x_0} \in \mathbb{R}^p \text{ und } x \approx x_0.$$

Wenn wir praktisch arbeiten wollen, reicht abstrakte Asymptotik nicht. Wir müssen das \approx spezifizieren. Dies kann skalenspezifisch geschehen (z.B. $x \approx x_0$ wenn $|x - x_0| < 3$) oder designabhängig (z.B. $x \approx x_0$ wenn $\#i : |x - x_i| \leq |x - x_0| < n/3$). Die heute üblichen Implementierungen haben feinere Varianten, die hier noch nicht diskutiert werden können. Der Illustration halber kann die folgende Vergrößerung reichen:

Lokalisierter Gauß-Markov-Schätzer:

Für $x \in \mathbb{R}^p$, bestimme

$$\delta = \min_d (\#i : |x - x_i| \leq d) \geq n \cdot f$$

wobei f ein gewählter Anteil (z.B. 0.5) ist.

Bestimme den Gauß-Markov-Schätzer $\hat{\beta}_x$, wobei nur diejenigen Beobachtungen einbezogen werden, für die $|x - x_i| \leq \delta$.

Schätze

$$\hat{m}(x) = x' \hat{\beta}_x.$$

Diese Vergrößerung ignoriert alle Messpunkte, die einen Abstand über δ haben. Feinere Methoden benutzen eine Gewichtung, um den Einfluss entfernter Messpunkte zunehmend zu reduzieren.

[help\(loess\)](#)

loess

Local Polynomial Regression Fitting

Description.

Fit a polynomial surface determined by one or more numerical predictors, using local fitting.

Usage.

```
loess(formula, data, weights, subset, na.action, model = FALSE,
      span = 0.75, enp.target, degree = 2,
      parametric = FALSE, drop.square = FALSE, normalize = TRUE,
      family = c("gaussian", "symmetric"),
      method = c("loess", "model.frame"),
      control = loess.control(...), ...)
```


Arguments.

<code>formula</code>	a formula specifying the numeric response and one to four numeric predictors (best specified via an interaction, but can also be specified additively).
<code>data</code>	an optional data frame, list or environment (or object coercible by <code>as.data.frame</code> to a data frame) containing the variables in the model. If not found in <code>data</code> , the variables are taken from <code>environment(formula)</code> , typically the environment from which <code>loess</code> is called.
<code>weights</code>	optional weights for each case.
<code>subset</code>	an optional specification of a subset of the data to be used.
<code>na.action</code>	the action to be taken with missing values in the response or predictors. The default is given by <code>getOption("na.action")</code> .
<code>model</code>	should the model frame be returned?
<code>span</code>	the parameter α which controls the degree of smoothing.
<code>enp.target</code>	an alternative way to specify <code>span</code> , as the approximate equivalent number of parameters to be used.
<code>degree</code>	the degree of the polynomials to be used, up to 2.
<code>parametric</code>	should any terms be fitted globally rather than locally? Terms can be specified by name, number or as a logical vector of the same length as the number of predictors.
<code>drop.square</code>	for fits with more than one predictor and <code>degree=2</code> , should the quadratic term (and cross-terms) be dropped for particular predictors? Terms are specified in the same way as for <code>parametric</code> .
<code>normalize</code>	should the predictors be normalized to a common scale if there is more than one? The normalization used is to set the 10% trimmed standard deviation to one. Set to false for spatial coordinate predictors and others known to be a common scale.
<code>family</code>	if "gaussian" fitting is by least-squares, and if "symmetric" a re-descending M estimator is used with Tukey's biweight function.
<code>method</code>	fit the model or just extract the model frame.
<code>control</code>	control parameters: see <code>loess.control</code> .
<code>...</code>	control parameters can also be supplied directly.

Details.

Fitting is done locally. That is, for the fit at point x , the fit is made using points in a neighbourhood of x , weighted by their distance from x (with differences in 'parametric' variables being ignored when computing the distance). The size of the neighbourhood is controlled by α (set by `span` or `enp.target`). For $\alpha < 1$, the neighbourhood includes proportion α of the points, and these have tricubic weighting (proportional to $(1 - (\text{dist}/\text{maxdist})^3)^3$). For $\alpha > 1$, all points are used, with the 'maximum distance' assumed to be $\alpha^{1/p}$ times the actual maximum distance for p explanatory variables.

For the default family, fitting is by (weighted) least squares. For `family="symmetric"` a few iterations of an M-estimation procedure with Tukey's biweight are used. Be aware that as the initial value is the least-squares fit, this need not be a very resistant fit.

It can be important to tune the control list to achieve acceptable speed. See `loess.control` for details.

Value.

An object of class "loess".

Note.

As this is based on the `cloess` package available at `netlib`, it is similar to but not identical to the `loess` function of S. In particular, conditioning is not implemented.

The memory usage of this implementation of `loess` is roughly quadratic in the number of points, with 1000 points taking about 10Mb.

Author(s).

B.D. Ripley, based on the `cloess` package of Cleveland, Grosse and Shyu available at <http://www.netlib.org/a/>.

References.

W.S. Cleveland, E. Grosse and W.M. Shyu (1992) Local regression models. Chapter 8 of *Statistical Models in S* eds J.M. Chambers and T.J. Hastie, Wadsworth & Brooks/Cole.

See Also.

`loess.control`, `predict.loess`.

`lowess`, the ancestor of `loess` (with different defaults!).

Examples.

```
cars.lo <- loess(dist ~ speed, cars)
predict(cars.lo, data.frame(speed = seq(5, 30, 1)), se = TRUE)
# to allow extrapolation
cars.lo2 <- loess(dist ~ speed, cars,
  control = loess.control(surface = "direct"))
predict(cars.lo2, data.frame(speed = seq(5, 30, 1)), se = TRUE)
```

Während die lineare Regression durch die Modellannahmen verpflichtet ist, immer ein lineares (oder linear parametrisiertes) Bild zu geben, können bei einer lokalisierten Variante auch Nichtlinearitäten dargestellt werden. Die Untersuchung dieser Familie von Verfahren bildet ein eigenes Teilgebiet der Statistik, die nichtparametrische Regression.

Wir bereiten wieder ein Beispiel vor:

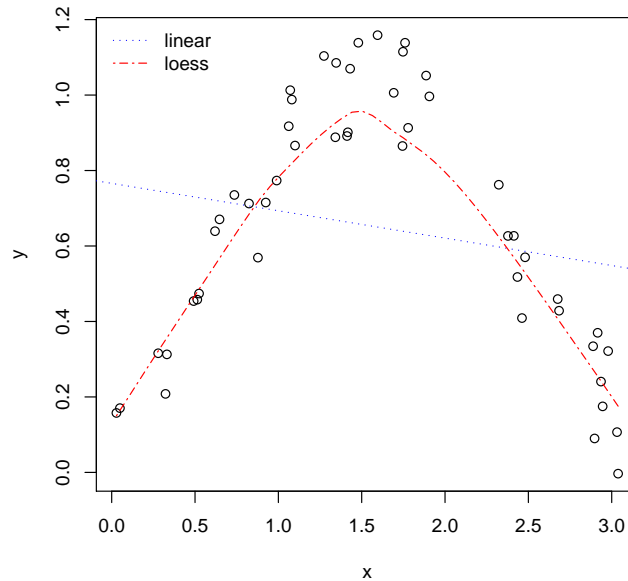
```
x <- runif(50) * pi
y <- sin(x)+rnorm(50)/10
```

Eingabe

Beispiel 2.11:

Eingabe

```
plot(x, y)
abline(lm(y ~ x), lty = 3, col = "blue")
lines(loess.smooth(x, y), lty = 6, col = "red")
legend("topleft",
      legend = c("linear", "loess"),
      lty = c(3, 6), col = c("blue", "red"), bty = "n")
```

**2.6. Ergänzungen**

2.6.1. Ergänzung: Diskretisierungen. Analog zum Vorgehen bei den Histogrammen können wir wieder diskretisieren. Im Hinblick auf die Regressoren haben wir dies beim Helikopter-Beispiel getan. Die Diskretisierung können wir auch bei der Respons vornehmen. Damit wird aus dem Regressionsproblem ein Kontingenztafel-Problem. Wir gehen hier nicht weiter auf diese Möglichkeit ein.

2.6.2. Ergänzung: Externe Daten. Daten, wie auch andere R-Objekte können mit `save()` in eine externe Datei geschrieben und mit `load()` wieder daraus gelesen werden. In diesen Dateien werden die Daten komprimiert; die Dateien sind zwischen verschiedenen R-Systemen austauschbar.

<i>Ein- Ausgabe von Daten für R</i>	
<code>save()</code>	Speichert Daten in eine externe Datei. Aufruf: <code>save(<Namen der zu speichernden Objekte>, file = <Dateiname>, ...)</code>

(Fortsetzung)→

Ein- Ausgabe von Daten für R (Fortsetzung)	
<code>load()</code>	Lädt Daten aus einer externe Datei. <i>Aufruf:</i> <code>load(file = <Dateiname>, ...)</code>

Häufig werden Daten in anderen Systemen vorbereitet. R stellt eine Reihe von Funktionen bereit, die Daten in unterschiedlichen Formaten einlesen können. Siehe dazu im Anhang Abschnitt A.17 auf Seite Seite A-29. Weitere Information findet sich im Manual “Data Import/Export” ([R D07b]).

Die Funktion `data()` bündelt verschieden Zugriffsroutinen, wenn die Zugriffspfade und Datei-Namen den R-Konventionen folgen.

In der Regel müssen eingelesene Daten noch nachbearbeitet werden, um das Format an die Aufrufkonventionen der gewünschten R-Funktionen anzupassen.

So erwartet z. B. `lm` die Regressoren als getrennte Variable. Für faktorielle Designs ist es hingegen üblich, die Resultate in einer Tafel zusammenzufassen, die die Faktor-Stufen als Zeilen- oder Spaltenlabels enthält. Die Funktion `stack()` überführt Tafeln in Spalten.

2.6.3. Ergänzung: Software-Test. Alle vorbereiteten Algorithmen, wie hier die Algorithmen zu den linearen Modellen und deren Varianten, sollten mit derselben Vorsicht behandelt werden wie mathematische Veröffentlichungen und Zitate. Selbst einfache Programme jedoch haben schnell eine semantische Komplexität, die weit über die mathematischer Beweise hinaus geht. Die übliche Strategie des “Nachrechnens” oder des schrittweisen Nachvollziehens verbietet sich dadurch. Anstelle einer vollständigen Überprüfung muss ein selektives Testen treten. Eine Teststrategie ist z.B. in [Sawitzki, 1994] beschrieben.

Die Überprüfung ist sowohl für die Implementierung als auch für den zu Grunde liegenden abstrakten Algorithmus nötig.

Aufgabe 2.18	
	Für diese Aufgabenserie sei $y_i = a + bx_i + \varepsilon_i$ mit $\varepsilon_i \text{ iid } \sim N(0, \sigma^2)$ und $x_i = i, i = 1, \dots, 10$.
	Wählen Sie eine Strategie, um <code>lm()</code> im Hinblick auf den Parameterraum (a, b, σ^2) zu überprüfen. Gibt es eine naheliegende Zellzerlegung für die einzelnen Parameter a, b, σ^2 ? Welche trivialen Fälle gibt es? Welche (uniforme) Asymptotik? Wählen Sie Testpunkte jeweils in der Mitte jeder Zelle und an den Rändern. Führen Sie diese Test durch und fassen Sie die Resultate zusammen.
	Welche Symmetrien/Antisymmetrien gibt es? Überprüfen Sie diese Symmetrien.
	Welche Invarianten/welches Covariate Verhalten gibt es? Überprüfen Sie diese Invarianten/Covariaten.

Aufgabe 2.19	
	Für diese Aufgabenserie sei $y_i = a + bx_i + \varepsilon_i$ mit ε_i iid $\sim N(0, \sigma^2)$.
	Welche extremen Designs (x_i) gibt es? Überprüfen Sie das Verhalten von $lm()$ bei vier extremalen Designs.
	Führen Sie die Aufgaben aus der letzten Gruppe aus, jetzt mit variablem Design. Fassen Sie Ihren Bericht zusammen.

Aufgabe 2.20	
	Für diese Aufgabenserie sei $y_i = a + bx_i + \varepsilon_i$ mit ε_i iid $\sim N(0, \sigma^2)$.
	Modifizieren Sie $lm()$ so, dass eine gesicherte Funktion für das einfache lineare Modell entsteht, die auch Abweichungen von den Modellannahmen untersucht.

2.6.4. R-Datentypen. R ist eine interpretierte Programmiersprache. Sie will es dem Anwender erlauben, Definitionen und Konkretisierungen flexibel zu handhaben. Aus Geschwindigkeitsgründen versucht R, Auswertungen so spät wie möglich durchzuführen. Dies erfordert einige Einschränkungen an die Sprache, die R von anderen Programmiersprachen unterscheidet.

R kennt keine abstrakten Datentypen. Ein Datentyp ist durch seine Instanzen, die Variablen, definiert.

Der Datentyp einer Variablen ist dynamisch: derselbe Name in denselben Kontext kann zu unterschiedlichen Zeiten unterschiedliche Variablenwerte und Variablentypen kennzeichnen.

Dennoch hat zu jeder Zeit eine Variable einen bestimmten Typ. Das R-Typensystem versteht man jedoch am besten in seiner historischen Entwicklung und die entsprechenden Funktionen. In der ersten Stufe war der Typ beschrieben durch `mode()` (z.B. "numeric") und `storage.mode()` (z.B. "integer" oder "real").

Beide Funktionen sind weitgehend durch `typeof()` abgelöst. Eine Zusammenfassung der Typen, die durch `typeof()` derzeit berichtet werden, ist in [R D07c] zu finden.

Komplexere Datentypen werden auf die in [R D07c] definierten zurückgeführt, indem die Variablen mit Attributen versehen werden. Dies geschieht mit der Funktion `attr()`, die auch benutzt werden kann, um Attribute zu inspizieren. So sind eine Matrix oder ein Array nur spezielle Vektoren, die sich dadurch auszeichnen, dass sie ein `dim`-Attribut haben. Das `class`-Attribut dient dazu, die Klasse explizit zu festzulegen.

Für die wesentlichen Typen sind Inspektionsprozeduren und Umwandlungsprozeduren: `is.<typ>()` prüft auf Typenzugehörigkeit, `as.<typ>()` wandelt den Typ.

2.6.5. Klassen und polymorphe Funktionen. Im Zuge der Weiterentwicklung wurde eine Anleihe an objekt-orientierte Programmierung gemacht. Dafür wurde ein spezielles Attribut mit dem Namen `class` genutzt: der Name des Typs (oder der "Klasse") wird hier gespeichert. Multiple Klassenzugehörigkeit in einer Hierarchie von Klassen ist auch möglich. In diesem Fall enthält `class` einen Vektor von Klassennamen. So hat zum Beispiel ein geordneter Faktor die Kennzeichnung `class = c("ordered", "factor")`. Zur Verwaltung der Klassen stehen Funktionen `class()`, `unclass()`, `inherits()` zur Verfügung.

Die Klassenzuordnung basiert dabei auf Vertrauen. R überprüft nicht, ob die Datenstruktur mit der angegebenen Klasse konsistent ist.

Funktionen wie `plot()`, `print()` und viele weitere überprüfen die Typen- und Klassenzugehörigkeit ihrer Argumente und verzweigen dann zu entsprechenden spezialisierten Funktionen. Dieses nennt man Polymorphismus. Wenn man eine polymorphe Funktion auflistet, erhält man zunächst nur den Hinweis, dass eine Dispatch-Funktion `UseMethod()` aufgerufen wird. Beispiel:

Beispiel 2.12:

	Eingabe	
<code>plot</code>		
<pre>function (x, y, ...) { if (is.function(x) && is.null(attr(x, "class"))) { if (missing(y)) y <- NULL hasylab <- function(...) !all(is.na(pmatch(names(list(...)), "ylab"))) if (hasylab(...)) plot.function(x, y, ...) else plot.function(x, y, ylab = paste(deparse(substitute(x)), "(x)", ...)) } else UseMethod("plot") } <environment: namespace:graphics></pre>		

`UseMethod()` bestimmt die Klasse des ersten Argumentes, mit dem die Funktion aufgerufen wurde, sucht dann eine Spezialisierung für diese Klasse und ruft schließlich die gefundene Funktion auf. Für *polymorphe* Funktionen findet man die entsprechenden Spezialisierungen mit Hilfe von `methods()`, z.B. `methods(plot)`.

2.6.6. Extraktor-Funktionen. Funktionen wie `lm()` liefern komplexe Datentypen mit umfangreicher Information. In einer rein objekt-orientierten Umgebung würden Zugriffsmethoden mit den Daten gemeinsam verkapselt. In R ist Objekt-Orientierung in Ansätzen und auf verschiedene Weisen realisiert. Dies spiegelt zum Teil die Entwicklung wieder. Bei genügend verallgemeinerbaren Strukturen werden Zugriffsmethoden wie in Abschnitt 2.6.5 bereitgestellt. Für die Objekte wie die von `lm()` gelieferten gibt es eine Reihe von Extraktor-Funktionen, die auf Komponenten zugreifen und diese geeignet aufbereiten.

<i>Extraktor-Funktionen für lm</i>	
<code>coef()</code>	extrahiert geschätzte Koeffizienten
<code>effects()</code>	extrahiert sukzessiv orthogonale Komponenten

(Fortsetzung) →

<i>Extraktor-Funktionen für lm</i> (Fortsetzung)	
<code>residuals()</code>	Roh-Residuen
<code>stdres()</code>	(in <code>library(MASS)</code>) standardisierte Residuen
<code>studres()</code>	(in <code>library(MASS)</code>) extern studentisierte Residuen
<code>fitted()</code>	
<code>vcov()</code>	Varianz/Kovarianzmatrix der geschätzten Parameter
<code>predict()</code>	Konfidenz- und Toleranzintervalle
<code>confint()</code>	Konfidenz-Intervalle für Parameter
<code>influence()</code>	extrahiert Einfluss-Diagnostiken
<code>model.matrix()</code>	bildet die Design-Matrix

2.7. Statistische Zusammenfassung

Als Leitbeispiel diente in diesem Kapitel die statistische Analyse eines funktionalen Zusammenhangs. Die betrachteten Modelle sind *finit* in dem Sinne, dass ein endlich-dimensionaler Funktionenraum den in Betracht gezogenen Zusammenhang zwischen Regressoren und Respons beschreibt. Die stochastische Komponente in diesen Modellen ist noch auf eine (eindimensionale) Zufallsverteilung beschränkt. Die Dimensionsbegriffe verdienen hier eine genauere Betrachtung. Wir haben zum einen die Regressor-Dimension. Dies ist die Dimension des Raumes der beobachteten oder abgeleiteten Parameter. Nicht alle Parameter sind identifizierbar oder schätzbar. Genauer gefasst ist die Dimension die Vektorraum-Dimension des gewählten Modell-Raums. Die Modelle werden durch Parameter in diesem Raum beschrieben. Diese Parameter können unbekannt oder hypothetisch sein. In jedem Fall aber haben wir sie als deterministisch betrachtet. Zum anderen haben wir die stochastische Komponente, repräsentiert durch den Fehler-Term. In diesem Kapitel sind wir von homogenen Fehlern ausgegangen. Damit bestimmt der Fehler-Term im Prinzip eine Dimension, die allerdings aus einem Raum von Verteilungen stammt. Für den Spezialfall der einfachen Gauß-linearen Modell sind die Verteilungen mit zwei Parametern präzisiert, dem Erwartungswert und der Varianz. Von dem Erwartungswert haben wir uns durch die Annahme befreit, dass das Modell im Mittel alle systematischen Effekte erfasst, also der Erwartungswert null ist. Die Varianz ist in unseren Problemen noch ein unbekannter Störparameter. Wir haben die dadurch entstehenden Problemen vermieden, indem wir uns auf Probleme beschränkt haben, in denen diese Störparameter durch einen geschätzten Wert ersetzt und so eliminiert wird.

2.8. Literatur und weitere Hinweise:

[CH92] Chambers, J.M.; Hastie, T.J. (eds.) (1992): *Statistical Models in S*. New York: Chapman & Hall.

[Jør93] Jørgensen, B. (1993): *The Theory of Linear Models*. New York: Chapman & Hall.

[R D07c] R Development Core Team (2004): *The R language definition*.

- [**Saw94a**] Sawitzki, G. (1994): Numerical Reliability of Data Analysis Systems. Computational Statistics & Data Analysis 18.2 (1994) 269-286. <<http://www.statlab.uni-heidelberg.de/reports/>>.
- [**Saw94b**] Sawitzki, G. (1994): Report on the Numerical Reliability of Data Analysis Systems. Computational Statistics & Data Analysis/SSN 18.2 (1994) 289-301. <<http://www.statlab.uni-heidelberg.de/reports/>>.

KAPITEL 3

Vergleich von Verteilungen

Wir beginnen mit der Konstruktion eines kleinen Werkzeugs, das uns Beispieldaten liefern wird. Basis ist ein kleiner Reaktionstester. Wir zeichnen einen “zufälligen” Punkt, warten auf einen Maus-Klick, und registrieren die Position des Mauszeigers. Damit bei wiederholten Aufrufen das Bild stabil bleibt, fixieren wir das Koordinatensystem.

Beispiel 3.1:

```
Eingabe -----  
plot(x = runif(1), y = runif(1),  
      xlim = c(0, 1), ylim = c(0, 1),  
      main = "Bitte auf den Punkt klicken",  
      xlab = '', ylab = '',  
      axes = FALSE, frame.plot = TRUE)  
locator(1)
```

```
Ausgabe -----  
$x  
[1] 0.6956522  
  
$y  
[1] 0.1260563
```

Bitte auf den Punkt klicken



Wir verpacken nun diesen Basistester. Wir merken uns die Koordinaten, versuchen, die Reaktionszeit des Benutzers zu messen, und liefern alle Resultate als Liste zurück.

Beispiel 3.2:

```

Eingabe
-----
click1 <- function(){
  x <- runif(1);y <- runif(1)
  plot(x = x, y = y, xlim = c(0, 1), ylim = c(0, 1),
       main = "Bitte auf den Punkt klicken",
       xlab = '', ylab = '',
       axes = FALSE, frame.plot = TRUE)
  clicktime <- system.time(xyclick <- locator(1))
  list(timestamp = Sys.time(),
       x = x, y = y,
       xclick = xyclick$x, yclick = xyclick$y,
       tclick = clicktime[3])
}

```

Zur weiteren Verarbeitung können wir die Liste in einen `data.frame` integrieren und diesen `data.frame` schrittweise mit Hilfe von `rbind` erweitern.

Beispiel 3.3:

```

Eingabe
-----
dx <- as.data.frame(click1())
dx <- rbind(dx, data.frame(click1()))
dx

```

```

Ausgabe
-----
      timestamp      x      y  xclick  yclick  tclick
elapsed 2008-03-17 21:40:49 0.29683 0.43955 0.69565 0.12606 0.261
elapsed1 2008-03-17 21:40:50 0.29617 0.58226 0.69565 0.12606 0.262

```

Aufgabe 3.1	
	<p>Definieren Sie eine Funktion <code>click(runs)</code>, die zu vorgegebener Anzahl <code>runs</code> die Aufgabe von <code>click1()</code> wiederholt und das Resultat als <code>data.frame</code> übergibt. Eine erste (zusätzliche) Messung sollte als Warmlaufen betrachtet werden und nicht in die Auswertung mit einbezogen werden.</p> <p>Wählen Sie eine Anzahl <code>runs</code>. Begründen Sie Ihre Wahl von <code>runs</code>. Führen Sie <code>click(runs)</code> durch und speichern Sie das Resultat mit Hilfe von <code>write.table()</code> in einer Datei.</p> <p>Stellen Sie die Verteilung der Komponente <code>tclick()</code> mit den Methoden aus Kapitel 1 (Verteilungsfunktion, Histogramm, Boxplot) dar.</p>

3.1. Shift/Skalenfamilien

Ein Vergleich von Verteilungen kann eine sehr anspruchsvolle Aufgabe sein. Der mathematische Raum, in dem Verteilungen angesiedelt sind, ist nicht mehr ein Zahlenraum oder ein (endlichdimensionaler) Vektorraum. Der eigentliche Raum, in dem Verteilungen beheimatet sind, ist ein Raum von Maßen. In einfachen Fällen, etwa bei Verteilungen auf \mathbb{R} , können wir alles auf Verteilungsfunktionen reduzieren und sind damit immerhin bei einem Funktionenraum. Selbst hier kann ein Vergleich noch große Schwierigkeiten machen. Wir haben keine einfache Ordnungsrelation.

Aufgabe 3.2	
	<p>Führen Sie Aufgabe 3.1 einmal mit der rechten und einmal mit der linken Hand durch. Vergleichen Sie die empirischen Verteilungen von <code>tclick()</code>.</p> <p>Die erhobenen Daten enthalten auch Information über die Positionen. Definieren Sie ein Maß <i>dist</i> für die Abweichung. Begründen Sie Ihre Definition. Führen Sie auch für <i>dist</i> einen rechts/links Vergleich durch.</p>

Wir konzentrieren uns hier auf den Vergleich von nur zwei Verteilungen, etwa der von Messungen in zwei Behandlungsgruppen. Wie nehmen wieder einen einfachen Fall: die Beobachtungen seien jeweils unabhängig identisch verteilt (jetzt mit der für den Vergleich von Behandlungen üblichen Index-Notation).

Y_{ij} unabhängig identisch verteilt mit Verteilungsfunktion F_i

$i = 1, 2$ Behandlungen

$j = 1, \dots, n_i$ Beobachtungen in Behandlungsgruppe i .

Wie vergleichen wir die Beobachtungen in den Gruppen $i = 1, 2$? Die (einfachen) linearen Modelle

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

betrachten als Unterschied häufig nur eine Verschiebung $\Delta = \alpha_1 - \alpha_2$.

Bezeichnungen: Zu einer Verteilung mit Verteilungsfunktion F heißt die Familie mit

$$F_a(x) = F(x - a)$$

die **Shift-Familie** zu F . Die Verschiebung a heißt Shift- oder Lage-Parameter.

Die Behandlung kann aber, in Wahrscheinlichkeiten gesprochen, die Wahrscheinlichkeitsmassen auch in anderer Weise verschieben, als es ein additiver Term im Modell bewirken kann. Wir brauchen allgemeinere Vergleichsmöglichkeiten als die durch einen Shift definierten.

Bezeichnung: Eine Verteilung mit Verteilungsfunktion F_1 heißt **stochastisch kleiner** als eine mit Verteilungsfunktion F_2 ($F_1 \prec F_2$), wenn F_1 eher bei kleineren Werten liegt als F_2 . Das bedeutet, dass F_1 eher ansteigt.

$$F_1(x) \geq F_2(x) \quad \forall x$$

und

$$F_1(x) > F_2(x) \quad \text{für mindestens ein } x.$$

Für Shift-Familien gilt: Ist $a < 0$, so ist $F_a \prec F$. Der Shift bewirkt eine Parallelverschiebung der Verteilungsfunktionen.

Ein typisches Resultat des Click-Experiments (Aufgabe 3.1 ist zum Beispiel in Abbildung 3.1 dargestellt. Die Zeiten für die rechte Seite sind stochastisch kleiner als die für die linke Seite. Die Verteilungen gehören jedoch nicht zu einer Shift-Familie, denn die Verteilungsfunktionen sind nicht parallel.

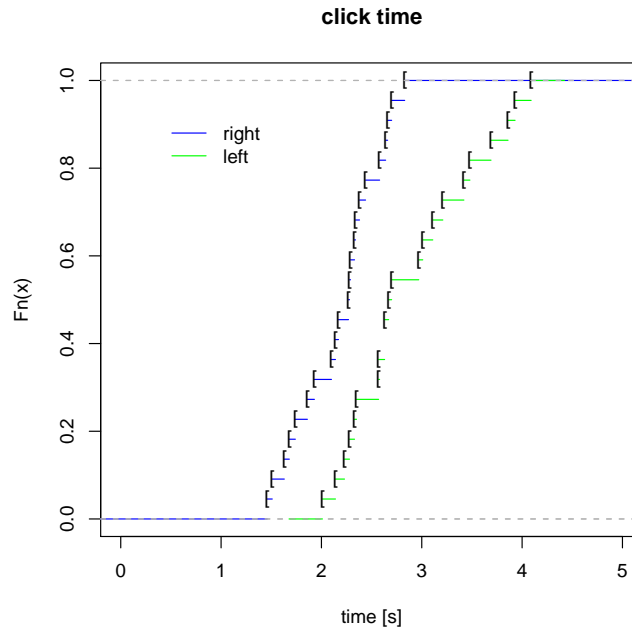


ABBILDUNG 3.1. Verteilungsfunktion für die rechts/links-Klickzeit

Aufgabe 3.3	
	Wie sieht ein <i>PP</i> -Plot für F_1 gegen F_2 aus, wenn $F_1 \prec F_2$?
	Wie sieht ein <i>QQ</i> -Plot für F_1 gegen F_2 aus, wenn $F_1 \prec F_2$?

Leider ist die dadurch definierte stochastische Ordnung nur von beschränktem Wert. Sie definiert keine vollständige Ordnung. Für Shift-Familien ist sie ausreichend. Aber Gegenbeispiele kann man sich konstruieren, wenn man die Shift-Familien nur geringfügig erweitert.

Bezeichnungen: Zu einer Verteilung mit Verteilungsfunktion F heißt die Familie mit

$$F_{a,b}(x) = F\left(\frac{x-a}{b}\right)$$

die *Skalen-Shiftfamilie* zu F .

Aufgabe 3.4	
	Die Skalen-Shiftfamilien zur $N(0, 1)$ -Verteilung sind die $N(\mu, \sigma^2)$ -Verteilungen. Welche $N(\mu, \sigma^2)$ -Verteilungen sind stochastisch kleiner als die $N(0, 1)$ -Verteilung? Welche sind stochastisch größer? Für welche ist die Ordnungsrelation undefiniert?

Die aus der linearen Theorie kommende Einordnung nach Lage/Skalen und die stochastische Ordnung klaffen auseinander, und beide Aspekte müssen oft getrennt betrachtet werden. Viele statistische Methoden konzentrieren sich auf Aspekte, die durch Skalen-Shiftfamilien motiviert sind. Unterschiede jenseits dessen, was durch Skala und Shift beschrieben werden kann, bedürfen oft besonderer Aufmerksamkeit.

In Kapitel 2 haben wir eine typische Situation für lineare Modelle betrachtet. Im Prinzip haben wir es mit Skalen-Shiftfamilien zu tun. Der (stochastische) Skalenparameter in diesen Modellen ist jedoch nur ein Störparameter, der eliminiert werden kann. Dazu benutzen wir einen Schätzer für diesen Skalenparameter, die residuelle Varianz, die wir dann heraus gekürzt haben. Als eine Besonderheit bei Gauß-linearen Modellen erhalten wir hier unabhängige Schätzer für Erwartungswert und Varianz. Dadurch können wir im Falle der einfachen Gauß-linearen Modelle Statistiken gewinnen, die nicht mehr vom Skalenparameter abhängen.

Im allgemeinen Fall haben wir jedoch eine aufsteigende Leiter von Problemen:

- Shift-Alternativen
- Shift/Skalen-Alternativen
- stochastische Ordnung
- allgemeinere Alternativen

Test- und Schätzprobleme konzentrieren sich oft nur auf einen Aspekt des Problems, die Lage. Der Skalenparameter ist hier nur eine Störgröße, ein “nuisance parameter”. Unterschiede im Shift-Parameter führen zu stochastisch monotonen Beziehungen. Unterschiede im Skalenparameter sind nicht so einfach einzuordnen und Test-Statistiken müssen erst von diesem, Störparameter bereinigt werden, wenn ausser dem Shift-Parameter auch der Skalenparameter variieren kann.

3.2. QQ-Plot, PP-Plot

Als Vergleichsdarstellung für Verteilungsfunktionen haben wir den *PP*-Plot und den *QQ*-Plot kennengelernt. So lange man innerhalb einer Skalen-Shiftfamilie bleibt, hat der *QQ*-Plot zumindest in einer Hinsicht einen Vorteil gegenüber dem *PP*-Plot:

BEMERKUNG 3.1. Sind F_1, F_2 Verteilungsfunktionen aus einer gemeinsamen Skalen-Shiftfamilie, so ist der *QQ*-Plot von F_1 gegen F_2 eine Gerade.

Insbesondere für die Gaußverteilungen ist der *QQ*-Plot gegen $N(0, 1)$ ein wichtiges Hilfsmittel. Jede Gaußverteilung gibt in diesem Plot eine Gerade. Der *QQ*-Plot ist für diese Situation bereits als Funktion `qqnorm()` vorbereitet.

Für den Vergleich von zwei Stichproben mit gleichem Stichprobenumfang kann die entsprechende Funktion `qqplot()` genutzt werden: bezeichnen wir die empirischen Quantile mit $Y_{1,(i:n)}$ bzw. $Y_{2,(i:n)}$, so ist dieser Plot der Graph $(Y_{1,(i:n)}, Y_{2,(i:n)})_{i=1..n}$. Sind die Stichprobenumfänge verschieden, so behilft sich R und generiert die Markierungspunkte durch lineare Interpolation, wobei der kleinere der beiden Stichprobenumfänge die Anzahl der Interpolationspunkte bestimmt.

Der *PP*-Plot hat keine dem *QQ*-Plot vergleichbare Äquivarianzeigenschaften. Wenn wir Skalen-Shiftparameter eliminieren wollen, müssen wir die Daten zunächst entsprechend transformieren. Die mathematische Theorie ist jedoch für den *PP*-Plot einfacher. Insbesondere gibt es auch hier einen entsprechenden Kolmogorov-Smirnov-Test (siehe Abschnitt 3.2.1).

Der Äquivarianz des QQ -Plots als Vorteil stehen auf der anderen Seite strukturelle Nachteile entgegen. In Bereichen niedriger Dichte bestimmen empirisch wenige Datenpunkte den Plot. Entsprechend hat er hier eine große Varianz. Gleichzeitig sind hier der Wahrscheinlichkeit nach benachbarte Quantile im Wertebereich weit entfernt: die große Varianz kombiniert sich ungünstig mit einer großen Variabilität, und der QQ -Plot zeigt entsprechend große Fluktuation. Für die meisten Lehrbuch-Verteilungen bedeutet dies, dass der QQ -Plot in den Randbereichen kaum zu interpretieren ist. Der PP -Plot hat keine entsprechenden Skalendefizite, aber auch nicht die Äquivarianzeigenschaft des QQ -Plots. Er wird deshalb in der Regel auf geeignet standardisierte Variable angewandt.

help(qqplot)

qqnorm

Quantile-Quantile Plots

Description.

`qqnorm` is a generic function the default method of which produces a normal QQ plot of the values in `y`. `qqline` adds a line to a normal quantile-quantile plot which passes through the first and third quartiles.

`qqplot` produces a QQ plot of two datasets.

Graphical parameters may be given as arguments to `qqnorm`, `qqplot` and `qqline`.

Usage.

```
qqnorm(y, ...)
## Default S3 method:
qqnorm(y, ylim, main = "Normal Q-Q Plot",
       xlab = "Theoretical Quantiles", ylab = "Sample Quantiles",
       plot.it = TRUE, datax = FALSE, ...)

qqline(y, datax = FALSE, ...)

qqplot(x, y, plot.it = TRUE, xlab = deparse(substitute(x)),
       ylab = deparse(substitute(y)), ...)
```

Arguments.

<code>x</code>	The first sample for <code>qqplot</code> .
<code>y</code>	The second or only data sample.
<code>xlab</code> , <code>ylab</code> , <code>main</code>	plot labels. The <code>xlab</code> and <code>ylab</code> refer to the y and x axes respectively if <code>datax = TRUE</code> .
<code>plot.it</code>	logical. Should the result be plotted?
<code>datax</code>	logical. Should data values be on the x-axis?
<code>ylim</code> , ...	graphical parameters.

Value.

For `qqnorm` and `qqplot`, a list with components

<code>x</code>	The x coordinates of the points that were/would be plotted
<code>y</code>	The original y vector, i.e., the corresponding y coordinates <i>including NAs</i> .

References.

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

See Also.

`ppoints`, used by `qqnorm` to generate approximations to expected order statistics for a normal distribution.

Examples.

```
y <- rt(200, df = 5)
qqnorm(y); qqline(y, col = 2)
qqplot(y, rt(300, df = 5))
```

```
qqnorm(precip, ylab = "Precipitation [in/yr] for 70 US cities")
```

Aufgabe 3.5	
	Benutzen Sie den Quantil-Quantil-Plot, um die Resultate des rechts/links <i>click</i> -Experiments zu vergleichen. Formulieren Sie die Resultate.
	Fassen Sie die rechts/links <i>tclick</i> -Daten zu einem Vektor zusammen. Vergleichen Sie den Quantil-Quantil-Plot mit dem von Monte-Carlo-Stichproben aus dem zusammengefassten Vektor. Erinnerung: Zufallsstichproben können Sie mit <code>sample()</code> ziehen. Mit <code>par(mfrow = c(2, 2))</code> teilen Sie den Zeichenbereich so ein, dass Sie vier Plots gleichzeitig sehen können.
**	Benutzen Sie bei <code>sample()</code> den Parameterwert <code>replace = FALSE</code> . Wie müssen Sie jetzt <code>sample()</code> anwenden, um den zusammengefassten Vektor in zwei Vektoren mit Monte-Carlo-Stichproben aufzuteilen? Welche Unterschiede zu <code>replace = TRUE</code> sind zu erwarten?

Aufgabe 3.6	
	Bestimmen Sie für die <i>tclick</i> -Daten des rechts/links <i>click</i> -Experiments Skalen- und Shiftparameter so, dass die Verteilungen in den Gruppen nach Skalen-Shift-Transformation möglichst gut übereinstimmen. Beschreiben Sie die Unterschiede anhand der Skalen-Shiftparameter. Verwenden Sie dazu eine Modellierung mit einem linearen Modell.
	Benutzen Sie die Funktion <code>boxplot()</code> , um Quartile und Flankenverhalten darzustellen. Vergleichen Sie die Information mit den Skalen-Shiftparametern. <i>Hinweis:</i> was entspricht dem Shift(Lage)parameter? Was entspricht dem Skalenparameter?

Wenn Darstellungen affin invariant sind, können Skalen-Shiftparameter ignoriert werden. Wenn Darstellungen nicht affin invariant sind, ist es häufig hilfreich, zunächst Skalen-Shiftparameter geeignet zu schätzen, die Verteilungen zu standardisieren, und dann die standardisierten Verteilungen zu untersuchen.

Das Problem, das wir uns damit potentiell einhandeln, ist, dass dann das stochastische Verhalten der Schätzung für die Skalen-Shiftparameter berücksichtigt werden muss. Der übliche Ausweg ist es, vorsichtigerweise “konservative” Tests und robuste Schätzer zu benutzen. Die folgende Transformation versucht, Skala und Lage an eine Standard-Normalverteilung anzupassen.

```

ScaleShiftStd <- function (x) { Eingabe
  xq <- quantile(x[!is.na(x)], c(0.25, 0.75))
  y <- qnorm(c(0.25, 0.75))
  slope <- diff(y)/diff(xq)
  (x-median(x, na.rm = FALSE)) * slope
}

```

Um Verteilungen direkt miteinander vergleichen zu können, greifen wir auf Techniken aus dem ersten Kapitel zurück. Was dort über den Vergleich zu einer theoretischen Verteilung gesagt worden ist, kann analog auf den Vergleich von zwei Verteilungen, z.B. aus zwei Behandlungsgruppen, übertragen werden. Die statistischen Aussagen müssen jedoch revidiert werden. Nun betrachten wir nicht mehr eine feste und eine zufällige Verteilung, sondern wir vergleichen zwei zufällige (empirische) Verteilungen.

Die für den Einstichproben-Fall (eine Stichprobe im Vergleich zu einer hypothetischen Verteilung) benutzte Idee von Monte-Carlo-Bändern kann nicht unmittelbar übertragen werden: wir wollen zwei Verteilungen miteinander vergleichen, aber wir haben keine ausgezeichnete Modellverteilung, aus der wir Referenzstichproben ziehen können.

Wir können jedoch die Idee modifizieren und bedingte Monte-Carlo-Bänder konstruieren. Bedingt bedeutet hier: die Konstruktion hängt von beobachteten Stichprobenwerten ab. Wir nehmen an, dass wir zwei Stichproben Y_{11}, \dots, Y_{1n_1} und Y_{21}, \dots, Y_{2n_2} von insgesamt unabhängigen und innerhalb der Gruppen identisch nach F_1 bzw. F_2 verteilten Beobachtungen haben. Falls kein Unterschied zwischen den Verteilungen besteht, so ist $(Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2})$ eine iid-Stichprobe aus einer gemeinsamen Verteilung $F = F_1 = F_2$ mit Stichprobenumfang $n = n_1 + n_2$. Bei einer iid-Stichprobe hätte jede Permutation der Indizes die gleiche Wahrscheinlichkeit.

Die motiviert das folgende Verfahren: wir permutieren das Tupel $(Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2})$ und ordnen die ersten n_1 Werte (nach Permutation) der ersten Gruppe zu, die anderen der zweiten.

Die Permutationsgruppe ist schnell so groß, dass sie nicht mehr vollständig ausgewertet werden kann. Anstelle dessen benutzen wir eine zufällige Auswahl von Permutationen. Wir benutzen die so generierten Werte, um Monte-Carlo-Bänder zu generieren.

Aufgabe 3.7	
	Modifizieren Sie die Funktionen für <i>PP</i> -Plot und <i>QQ</i> -Plot so, dass Monte-Carlo-Bänder für den Vergleich von zwei Stichproben hinzugefügt werden. <p style="text-align: right;">(Fortsetzung)→</p>

Aufgabe 3.7	(Fortsetzung)
	<i>Hinweis:</i> mit der Funktion <code>sample()</code> können Sie zufällige Permutationen generieren.

Bei größerem Stichprobenumfang kann der Aufwand Permutationen zu generieren zu zeitaufwendig sein. Um Verwaltungsaufwand zu sparen, können wir die Permutation durch ein Ziehen aus den n Werten $(Y_{11}, \dots, Y_{1n_1}, Y_{12}, \dots, Y_{1n_2})$ *mit Zurücklegen* ersetzen. Diese approximative Lösung wird als **Bootstrap-Approximation**¹ bezeichnet.

Da es nur endlich viele Permutationen gibt, können wir bei kleinem Stichprobenumfang auch alle Permutationen durchgehen. Wir wählen die Bänder dann so, dass ein hinreichend großer Anteil (etwa mehr als 95 %) aller Kurven innerhalb der Bänder liegt. Permutationen, die sich nur innerhalb der Gruppen unterscheiden, ergeben dieselben Kurven. Diese Zusatzüberlegung zeigt, dass wir nicht alle $n!$ Permutationen überprüfen müssen, sondern nur die $\binom{n}{n_1}$ Auswahlen für die Zuteilung zu den Gruppen.

Aufgabe 3.8	
**	Ergänzen Sie <i>PP</i> -Plot und <i>QQ</i> -Plot für die <i>click</i> -Experimente durch Permutations-Bänder, die 95 % der Permutationen abdecken.
*	Erzeugen Sie neue Plots, in denen Sie die <i>PP</i> -Plots und <i>QQ</i> -Plots durch Monte-Carlo-Bänder aus den Permutationen ergänzen. Benutzen Sie die Einhüllende von 19 Monte-Carlo-Stichproben. <i>Hinweis:</i> benutzen Sie die Funktion <code>sample()</code> um eine Stichprobe vom Umfang n_1 aus $x = (Y_{11}, \dots, Y_{1n_1}, Y_{12}, \dots, Y_{1n_2})$ zu ziehen.
	<i>Hinweis:</i> Siehe <code>help(sample)</code> .

Aufgabe 3.9	
*	Versuchen Sie, die Eigenschaften der Permutationsbänder, Monte-Carlo-Bänder und Bootstrap-Bänder zu vergleichen, wenn $F_1 = F_2$ gilt.

Wenn nicht die Verteilungen verglichen werden sollen, sondern nur einzelne festgelegte Kenngrößen, so können diese Strategien analog eingesetzt werden. Wenn wir uns z.B. auf die Shift-Alternative beschränken (d.h. F_1 und F_2 sind aus eine Shiftfamilie, d.h. $F_1(x) = F_2(x - a)$ für ein a), so können wir etwa den Mittelwert (oder den Median) als Kenngröße nehmen. Auf diese Kenngröße kann das obige Vorgehen analog angewandt werden, um zu entscheiden, ob die Hypothese, dass die Verteilungen sich nicht unterscheiden ($a = 0$), angesichts der Daten haltbar ist.

Aufgabe 3.10	
*	Formulieren Sie die obigen Strategien für Intervalle für einzelne Teststatistiken (Beispiel: Mittelwert) anstelle für Bänder. <p style="text-align: right;">(Fortsetzung)→</p>

¹Vorsicht: es gibt beliebig wilde Definitionen von Bootstrap. Versuchen Sie stets, das Vorgehen mathematisch genau zu formulieren, wenn von Bootstrap die Rede ist.

Aufgabe 3.10	(Fortsetzung)
	<i>Hinweis:</i> Können Sie anstelle der zwei Mittelwerte für beide Gruppen eine eindimensionale zusammenfassende Statistik benutzen?

3.2.1. Kolmogorov Smirnov Tests. In Kapitel 1 haben wir den Kolmogorov-Smirnov-Test zum Vergleich einer Stichprobe $(X_i)_{i=1,\dots,n}$ und der zugehörigen empirischen Verteilung F_n mit einer (festen, vorgegebenen) Verteilung F kennengelernt. Die kritische Testgröße ist dabei

$$\sup |F_n - F|.$$

Wir können diesen Test etwas modifizieren, um zwei empirische Verteilungen zu vergleichen. Anstelle der Modellverteilung F tritt nun eine zweite empirische Verteilung G_m von Beobachtungen $(Y_j)_{j=1,\dots,m}$ mit zu Grunde liegender (unbekannter) Verteilung G . Die kritische Testgröße ist dann

$$\sup |F_n - G_m|.$$

Der darauf basierende Test ist in der Literatur als 2-Stichproben-Kolmogorov-Smirnov-Test zu finden. Dieser Test korrespondiert zum *PP*-Plot und erlaubt es, Bänder zum *PP*-Plot zu konstruieren.

Wir können Bänder auch durch Simulation bestimmen. Im Gegensatz zum 1-Stichproben-Test haben wir jetzt keine vorgegebene Verteilung, aus der wir simulieren können. Unter der Hypothese, dass die Verteilungen F und G sich nicht unterscheiden, verhält sich jedoch bei unabhängigen Beobachtungen der gemeinsame Vektor $(X_1, \dots, X_n, Y_1, \dots, Y_m)$ wie ein Vektor von $n + m$ unabhängigen Zufallszahlen mit identischer Verteilung $F = G$. Bei gegebenen Daten kann diese Beziehung zur Simulation genutzt werden. Durch eine Permutation π der Indizes erzeugt man aus dem Vektor $Z = (X_1, \dots, X_n, Y_1, \dots, Y_m)$ einen neuen Vektor Z' mit $Z'_i = Z_{\pi(i)}$. Die ersten n Komponenten benutzen wir als simulierte Werte $(X'_i)_{i=1,\dots,n}$, die übrigen m Komponenten als simulierte Werte $(Y'_j)_{j=1,\dots,m}$.

Aufgabe 3.11	
*	Programmieren Sie diesen Algorithmus und ergänzen Sie den <i>PP</i> -Plot durch simulierte <i>PP</i> -Plots für eine kleine Anzahl (19?) von Permutation.
	Bestimmen Sie die Permutationsverteilung von $\sup F_n - G_m $ aus den Simulation und berechnen Sie diesen Wert für die ursprünglichen Daten. Können Sie diesen Vergleich benutzen, um ein Testverfahren zu definieren?
	Der implementierte Kolmogorov-Smirnov-Test beinhaltet eine Approximation für den 2-Stichprobenfall. In unserer Simulation wissen wir, dass wir unter der Hypothese simulieren, die Hypothese also zutrifft. Untersuchen Sie die Verteilung des nominellen Niveaus unter den simulierten Bedingungen.

3.3. Tests auf Shift

Wenn wir zusätzliche Verteilungsannahmen machen, können wir andere Entscheidungsverfahren wählen. Für diese Verfahren sind aber die Verteilungsannahmen kritisch. Diese Abhängigkeit von den Verteilungsannahmen kann gemildert oder vermieden werden, wenn

wir geeignete Verteilungsannahmen sicherstellen können. Der F -Test, den wir im letzten Kapitel kennengelernt haben, ist ein Beispiel für ein verteilungsabhängiges Verfahren. Für den Zwei-Stichprobenfall kann dieser Test modifiziert werden zum t -Test, der auch die Richtung des Unterschiedes widerspiegelt. (Das Quadrat der t -Statistik ist eine F -Statistik.)

help(t.test)

`t.test` *Student's t-Test*

Description.

Performs one and two sample t-tests on vectors of data.

Usage.

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula':
```

```
t.test(formula, data, subset, na.action, ...)
```

Arguments.

<code>x</code>	a numeric vector of data values.
<code>y</code>	an optional numeric vector data values.
<code>alternative</code>	a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.
<code>mu</code>	a number indicating the true value of the mean (or difference in means if you are performing a two sample test).
<code>paired</code>	a logical indicating whether you want a paired t-test.
<code>var.equal</code>	a logical variable indicating whether to treat the two variances as being equal. If TRUE then the pooled variance is used to estimate the variance otherwise the Welch (or Satterthwaite) approximation to the degrees of freedom is used.
<code>conf.level</code>	confidence level of the interval.
<code>formula</code>	a formula of the form <code>lhs ~ rhs</code> where <code>lhs</code> is a numeric variable giving the data values and <code>rhs</code> a factor with two levels giving the corresponding groups.
<code>data</code>	an optional matrix or data frame (or similar: see <code>model.frame</code>) containing the variables in the formula <code>formula</code> . By default the variables are taken from <code>environment(formula)</code> .
<code>subset</code>	an optional vector specifying a subset of observations to be used.
<code>na.action</code>	a function which indicates what should happen when the data contain NAs. Defaults to <code>getOption("na.action")</code> .
<code>...</code>	further arguments to be passed to or from methods.

Details.

The formula interface is only applicable for the 2-sample tests.

If `paired` is `TRUE` then both `x` and `y` must be specified and they must be the same length. Missing values are removed (in pairs if `paired` is `TRUE`). If `var.equal` is `TRUE` then the pooled estimate of the variance is used. By default, if `var.equal` is `FALSE` then the variance is estimated separately for both groups and the Welch modification to the degrees of freedom is used.

If the input data are effectively constant (compared to the larger of the two means) an error is generated.

Value.

A list with class `"htest"` containing the following components:

<code>statistic</code>	the value of the t-statistic.
<code>parameter</code>	the degrees of freedom for the t-statistic.
<code>p.value</code>	the p-value for the test.
<code>conf.int</code>	a confidence interval for the mean appropriate to the specified alternative hypothesis.
<code>estimate</code>	the estimated mean or difference in means depending on whether it was a one-sample test or a two-sample test.
<code>null.value</code>	the specified hypothesized value of the mean or mean difference depending on whether it was a one-sample test or a two-sample test.
<code>alternative</code>	a character string describing the alternative hypothesis.
<code>method</code>	a character string indicating what type of t-test was performed.
<code>data.name</code>	a character string giving the name(s) of the data.

See Also.

`prop.test`

Examples.

```
t.test(1:10,y=c(7:20))      # P = .00001855
t.test(1:10,y=c(7:20, 200)) # P = .1245    -- NOT significant anymore

## Classical example: Student's sleep data
plot(extra ~ group, data = sleep)
## Traditional interface
with(sleep, t.test(extra[group == 1], extra[group == 2]))
## Formula interface
t.test(extra ~ group, data = sleep)
```

In seiner einfachsten Form setzt der t -Test voraus, dass wir unabhängig identisch verteilte Stichproben aus Normalverteilungen haben. Tatsächlich reichen schwächere Voraussetzungen. Wenn wir die t -Test-Statistik als

$$(3.1) \quad t = \frac{\widehat{\mu}_1 - \widehat{\mu}_2}{\sqrt{\widehat{Var}(\widehat{\mu}_1 - \widehat{\mu}_2)}}$$

schreiben, so sehen wir, dass t t -verteilt ist, wenn $\widehat{\mu}_1 - \widehat{\mu}_2$ normalverteilt und $(\widehat{Var}(\widehat{\mu}_1 - \widehat{\mu}_2))$ χ^2 verteilt ist, und beide Term unabhängig sind. Der zentrale Grenzwertsatz garantiert, dass $\widehat{\mu}_1 - \widehat{\mu}_2$ unter milden Bedingungen zumindest asymptotisch normalverteilt ist. Analoges gilt oft für $(\widehat{Var}(\widehat{\mu}_1 - \widehat{\mu}_2))$. Gilt die Unabhängigkeit beider Terme, so ist t approximativ t -verteilt.

Aufgabe 3.12	
*	<p>Bestimmen Sie in einer Simulation die Verteilung von \bar{Y}, $\widehat{Var}(Y)$ und der t-Statistik für Y aus der uniformen Verteilung $U[0, 1]$ mit Stichprobenumfang $n = 1, \dots, 10$. Vergleichen Sie die Verteilungen aus der Simulation mit der entsprechenden Normal-, χ^2- bzw. t-Verteilung.</p> <p>Bestimmen Sie in einer Simulation die Verteilung von \bar{Y}, $\widehat{Var}(Y)$ und der t-Statistik für Y aus einer Mischung, die zu 90% aus einer $N(0, 1)$- und zu 10% aus einer $N(0, 10)$-Verteilung stammt, mit Stichprobenumfang $n = 1, \dots, 10$. Vergleichen Sie die Verteilungen aus der Simulation mit der entsprechenden Normal-, χ^2- bzw. t-Verteilung.</p>

Der t -Test hat eine gewisse Robustheit, die ihm eine approximative Gültigkeit geben kann. Man kann sich jedoch ganz von der Normalverteilungs-Voraussetzung befreien. Wenn wir analog zum F -Test bzw. t -Test vorgehen, aber anstelle der Urdaten die Ränge benutzen, gewinnen wir Testverfahren, die verteilungsunabhängig sind (zumindest, solange keine Bindungen auftreten können). Der Wilcoxon-Test ist eine verteilungsunabhängige Variante des t -Tests. Theoretisch entspricht er genau dem t -Test, angewandt auf die (gemeinsam) rangtransformierten Daten. Wie der t -Test ist dieser Test nur darauf ausgelegt, die Nullhypothese (kein Unterschied) gegen eine Shift-Alternative zu testen. Für die praktische Anwendung können arithmetische Vereinfachungen ausgenutzt werden. Deshalb ist die Beziehung zwischen den üblichen Formeln für den t -Test und für den Wilcoxon-Test nicht einfach zu erkennen.

Um den Wilcoxon-Test anzuwenden, muss zum einen die Teststatistik berechnet werden. Zur Bestimmung kritischer Werte, mit denen die Teststatistik zu vergleichen ist, muss zum anderen die Verteilungsfunktion ausgewertet werden. Sind alle Beobachtungen paarweise verschieden, so hängt diese Funktion nur von n_1 und n_2 ab, und relativ einfache Algorithmen stehen zur Verfügung. Diese sind in der Funktion R standardmäßig verfügbar und werden von `wilcox.test()` benutzt. Gibt es Bindungen in den Daten, d.h. gibt es übereinstimmende Werte, so hängt die Verteilung vom speziellen Muster dieser Bindungen ab und die Berechnung ist aufwendiger. `wilcox.test()` greift in diesem Fall auf Approximationen zurück. Zur exakten (im Gegensatz zur approximativen) Auswertung stehen jedoch die entsprechenden Algorithmen ebenfalls zur Verfügung. Dazu benötigt man `library(coin)`. Die exakte Variante des Wilcoxon-Tests findet sich dort etwa als `wilcox_test()`.

Auf den Rängen basierende verteilungsunabhängige Verfahren zu charakterisieren und mit den früher vorgestellten verteilungsunabhängigen Monte-Carlo-Verfahren und deren Varianten zu vergleichen ist ein klassischer Teil der Statistik. Literatur dazu findet man unter den Schlagworten "Rangtests" oder "verteilungsfreie Verfahren". Zusätzliche R-Funktionen finden sich in `library(coin)` sowie in einigen speziellen Paketen.

Natürlich stellt sich die Frage nach dem Informationsverlust. Wenn wir uns auf die Daten beschränken und keine oder geringe Verteilungsannahmen machen, haben wir weniger Information als in einem Modell mit expliziten Verteilungsannahmen. Wenn wir die Daten auf die Ränge reduzieren, verschenken wir zusätzlich möglicherweise Information. Dieser Informationsverlust kann z.B. durch die asymptotische relative Effizienz gemessen werden. Dies ist (asymptotisch) der Stichprobenumfang eines optimalen Tests, der benötigt wird, eine vergleichbare Güte wie ein konkurrierender Test zu erreichen. Beim Wilcoxon-Test

unter Normalverteilung hat dies einen Wert von 94%. Gilt also die Normalverteilungsannahme, so benötigt der (optimale) t -Test nur 94% des Stichprobenumfangs, die der Wilcoxon-Test benötigt. 6% des Stichprobenumfangs sind die Kosten für die Reduzierung auf Ränge. Gilt die Normalverteilungsannahme nicht, so kann der t -Test möglicherweise zusammenbrechen. Der Wilcoxon-Test bleibt ein valider Test auf die Shift-Alternative.

[help\(wilcox.test\)](#)

`wilcox.test` *Wilcoxon Rank Sum and Signed Rank Tests*

Description.

Performs one and two sample Wilcoxon tests on vectors of data; the latter is also known as ‘Mann-Whitney’ test.

Usage.

```
wilcox.test(x, ...)
```

Default S3 method:

```
wilcox.test(x, y = NULL,
            alternative = c("two.sided", "less", "greater"),
            mu = 0, paired = FALSE, exact = NULL, correct = TRUE,
            conf.int = FALSE, conf.level = 0.95, ...)
```

S3 method for class 'formula':

```
wilcox.test(formula, data, subset, na.action, ...)
```

Arguments.

<code>x</code>	numeric vector of data values. Non-finite (e.g. infinite or missing) values will be omitted.
<code>y</code>	an optional numeric vector of data values.
<code>alternative</code>	a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.
<code>mu</code>	a number specifying an optional location parameter.
<code>paired</code>	a logical indicating whether you want a paired test.
<code>exact</code>	a logical indicating whether an exact p-value should be computed.
<code>correct</code>	a logical indicating whether to apply continuity correction in the normal approximation for the p-value.
<code>conf.int</code>	a logical indicating whether a confidence interval should be computed.
<code>conf.level</code>	confidence level of the interval.
<code>formula</code>	a formula of the form <code>lhs ~ rhs</code> where <code>lhs</code> is a numeric variable giving the data values and <code>rhs</code> a factor with two levels giving the corresponding groups.
<code>data</code>	an optional matrix or data frame (or similar: see <code>model.frame</code>) containing the variables in the formula <code>formula</code> . By default the variables are taken from <code>environment(formula)</code> .
<code>subset</code>	an optional vector specifying a subset of observations to be used.
<code>na.action</code>	a function which indicates what should happen when the data contain NAs. Defaults to <code>getOption("na.action")</code> .
<code>...</code>	further arguments to be passed to or from methods.

Details.

The formula interface is only applicable for the 2-sample tests.

If only `x` is given, or if both `x` and `y` are given and `paired` is `TRUE`, a Wilcoxon signed rank test of the null that the distribution of `x` (in the one sample case) or of `x-y` (in the paired two sample case) is symmetric about `mu` is performed.

Otherwise, if both `x` and `y` are given and `paired` is `FALSE`, a Wilcoxon rank sum test (equivalent to the Mann-Whitney test: see the Note) is carried out. In this case, the null hypothesis is that the distributions of `x` and `y` differ by a location shift of `mu` and the alternative is that they differ by some other location shift.

By default (if `exact` is not specified), an exact p-value is computed if the samples contain less than 50 finite values and there are no ties. Otherwise, a normal approximation is used.

Optionally (if argument `conf.int` is true), a nonparametric confidence interval and an estimator for the pseudomedian (one-sample case) or for the difference of the location parameters `x-y` is computed. (The pseudomedian of a distribution F is the median of the distribution of $(u + v)/2$, where u and v are independent, each with distribution F . If F is symmetric, then the pseudomedian and median coincide. See Hollander & Wolfe (1973), page 34.) If exact p-values are available, an exact confidence interval is obtained by the algorithm described in Bauer (1972), and the Hodges-Lehmann estimator is employed. Otherwise, the returned confidence interval and point estimate are based on normal approximations.

With small samples it may not be possible to achieve very high confidence interval coverages. If this happens a warning will be given and an interval with lower coverage will be substituted.

Value.

A list with class `"htest"` containing the following components:

<code>statistic</code>	the value of the test statistic with a name describing it.
<code>parameter</code>	the parameter(s) for the exact distribution of the test statistic.
<code>p.value</code>	the p-value for the test.
<code>null.value</code>	the location parameter <code>mu</code> .
<code>alternative</code>	a character string describing the alternative hypothesis.
<code>method</code>	the type of test applied.
<code>data.name</code>	a character string giving the names of the data.
<code>conf.int</code>	a confidence interval for the location parameter. (Only present if argument <code>conf.int</code> = <code>TRUE</code> .)
<code>estimate</code>	an estimate of the location parameter. (Only present if argument <code>conf.int</code> = <code>TRUE</code> .)

Warning.

This function can use large amounts of memory and stack (and even crash R if the stack limit is exceeded) if `exact` = `TRUE` and one sample is large (several thousands or more).

Note.

The literature is not unanimous about the definitions of the Wilcoxon rank sum and Mann-Whitney tests. The two most common definitions correspond to the sum of the ranks of the first sample with the minimum value subtracted or not: R subtracts and S-PLUS does not, giving a value which is larger by $m(m + 1)/2$ for a first sample of

size m . (It seems Wilcoxon's original paper used the unadjusted sum of the ranks but subsequent tables subtracted the minimum.)

R 's value can also be computed as the number of all pairs $(x[i], y[j])$ for which $y[j]$ is not greater than $x[i]$, the most common definition of the Mann-Whitney test.

References.

Myles Hollander & Douglas A. Wolfe (1973), *Nonparametric statistical inference*. New York: John Wiley & Sons. Pages 27–33 (one-sample), 68–75 (two-sample). Or second edition (1999).

David F. Bauer (1972), Constructing confidence sets using rank statistics. *Journal of the American Statistical Association* **67**, 687–690.

See Also.

`psignrank`, `pwilcox`.

`wilcox.exact` in **exactRankTests** covers much of the same ground, but also produces exact p-values in the presence of ties.

`kruskal.test` for testing homogeneity in location parameters in the case of two or more samples; `t.test` for an alternative under normality assumptions [or large samples]

Examples.

```
## One-sample test.
## Hollander & Wolfe (1973), 29f.
## Hamilton depression scale factor measurements in 9 patients with
## mixed anxiety and depression, taken at the first (x) and second
## (y) visit after initiation of a therapy (administration of a
## tranquilizer).
x <- c(1.83, 0.50, 1.62, 2.48, 1.68, 1.88, 1.55, 3.06, 1.30)
y <- c(0.878, 0.647, 0.598, 2.05, 1.06, 1.29, 1.06, 3.14, 1.29)
wilcox.test(x, y, paired = TRUE, alternative = "greater")
wilcox.test(y - x, alternative = "less") # The same.
wilcox.test(y - x, alternative = "less",
            exact = FALSE, correct = FALSE) # H&W large sample
                                           # approximation

## Two-sample test.
## Hollander & Wolfe (1973), 69f.
## Permeability constants of the human chorioamnion (a placental
## membrane) at term (x) and between 12 to 26 weeks gestational
## age (y). The alternative of interest is greater permeability
## of the human chorioamnion for the term pregnancy.
x <- c(0.80, 0.83, 1.89, 1.04, 1.45, 1.38, 1.91, 1.64, 0.73, 1.46)
y <- c(1.15, 0.88, 0.90, 0.74, 1.21)
wilcox.test(x, y, alternative = "g") # greater
wilcox.test(x, y, alternative = "greater",
            exact = FALSE, correct = FALSE) # H&W large sample
                                           # approximation

wilcox.test(rnorm(10), rnorm(10, 2), conf.int = TRUE)
```



```
## Formula interface.
boxplot(Ozone ~ Month, data = airquality)
wilcox.test(Ozone ~ Month, data = airquality,
            subset = Month %in% c(5, 8))
```

Aufgabe 3.13	
	Benutzen Sie den Wilcoxon-Test, um die Resultate des rechts/links <i>click</i> -Experiments zu vergleichen.

Aufgabe 3.14	
***	<p>Beim rechts/links <i>click</i>-Experiment sind mehrere Effekte vermischt. Einige Probleme:</p> <ul style="list-style-type: none"> • Die Antwortzeit beinhaltet Reaktionszeit, Zeit für die Grob-Bewegung der Maus, Zeit für die Fein-Adjustierung etc. • Für rechts-links-Bewegungen reicht in der Regel ein Schwenken der Hand aus. Für vorwärts-rückwärts-Bewegungen ist in der Regel eine Arm-Bewegung nötig. Es ist nicht zu erwarten, dass beide vergleichbares statistisches Verhalten haben. • Bei aufeinanderfolgenden Registrierungen kann es zum einen Trainings- zum anderen Ermüdungseffekte geben. <p>Können Sie Experiment und Auswertung so modifizieren, dass Unterschiede in der Reaktionszeit untersucht werden können? Können Sie Experiment und Auswertung so modifizieren, dass Unterschiede in der Genauigkeit der Endposition untersucht werden können?</p>
***	<p>Untersuchen und dokumentieren Sie für sich rechts-links-Unterschiede in der Reaktionszeit und in der Genauigkeit. Formulieren Sie ihr Resultat als Bericht.</p>

Aufgabe 3.15	
	<p>Betrachten Sie als Verteilungsfamilien die Shift/Skalenfamilien von $N(0, 1)$ und $t(3)$. Entwerfen Sie ein Szenario, um den Wilcoxon-Test mit dem t-Test jeweils innerhalb dieser Familien zu vergleichen.</p> <p>Führen Sie diesen Test in einer Simulation für Stichprobenumfänge $n_1 = n_2 = 10, 20, 50, 100$ durch und fassen Sie die Resultate zusammen.</p> <p>Führen Sie eine analoge Simulation für die Lognormal-Verteilungen durch.</p>

3.4. Güte

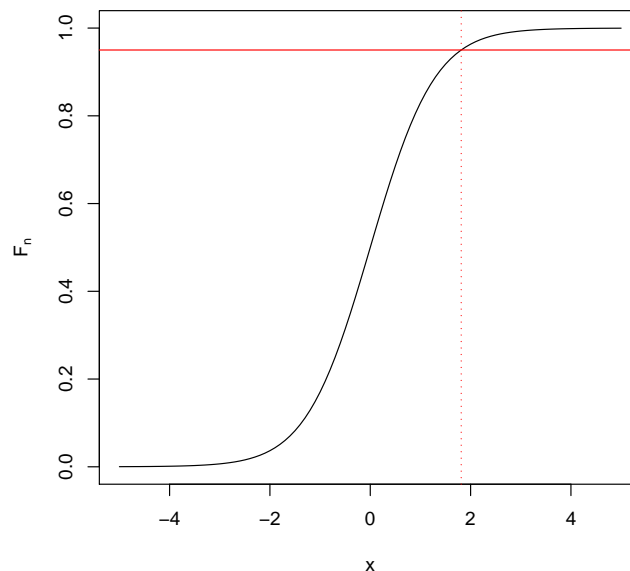
3.4.1. Theoretische Güte. Am Beispiel des t -Tests können wir illustrieren, wie ein Test aufgebaut ist. Der Test benutzt eine Teststatistik, hier die t -Test-Statistik zum

Vergleich zweier Gruppen (3.1). Wir kennen die Verteilung dieser Statistik: für den t -Test ist bei unabhängigen normalverteilten Fehlern und gleicher Varianz die Teststatistik $t(n_1 + n_2 - 2)$ verteilt. Zu gewähltem Niveau α können wir aus der Verteilungsfunktion Grenzen ablesen, die bei dieser Verteilung nur mit einer Wahrscheinlichkeit α unter- bzw. überschritten werden. Benutzen wir beide Grenzen, so erhalten wir einen zweiseitigen Bereich mit der Irrtumswahrscheinlichkeit 2α .

Beispiel 3.4:

Eingabe

```
n1<- 6; n2 <- 6
df <-  n1 + n2 -2
alpha <- 0.05
curve(pt(x,df=df),from=-5, to=5, ylab= expression(F[n]))
abline(h=1-alpha, col="red")      # cut at upper quantile
abline(v=qt(1-alpha, df=df), lty=3, col="red") # get critical value
```



Wollen wir z.B. die Hypothese $\mu_1 = \mu_2$ gegen die Alternative $\mu_1 > \mu_2$ testen, so wählen wir als Verwerfungsbereich den Bereich über der oberen dieser Grenze. Wir wissen, dass wir bei Gültigkeit der Hypothese höchstens mit Wahrscheinlichkeit α zufällig eine Testgröße in diesem Bereich bekommen.

Für den t -Test wissen wir sogar mehr. Unter den Modellvoraussetzungen unabhängig normalverteilter Fehler und gleicher Varianz ist die t -Test-Statistik immer t verteilt. Auf der Hypothese ist sie t verteilt mit Nichzentralitätsparameter 0, folgt also der zentralen t -Verteilung. Auf der Alternative haben wir eine t -Verteilung mit Nichzentralitätsparameter $(\mu_1 - \mu_2)\sigma^{-1}\sqrt{n_1 n_2 / (n_1 + n_2)}$. Damit kann für jede Alternative unter den Modellannahmen die Stärke des Tests abgelesen werden, d.h. die Wahrscheinlichkeit bei Vorliegen dieser Alternative die Hypothese zu verwerfen.

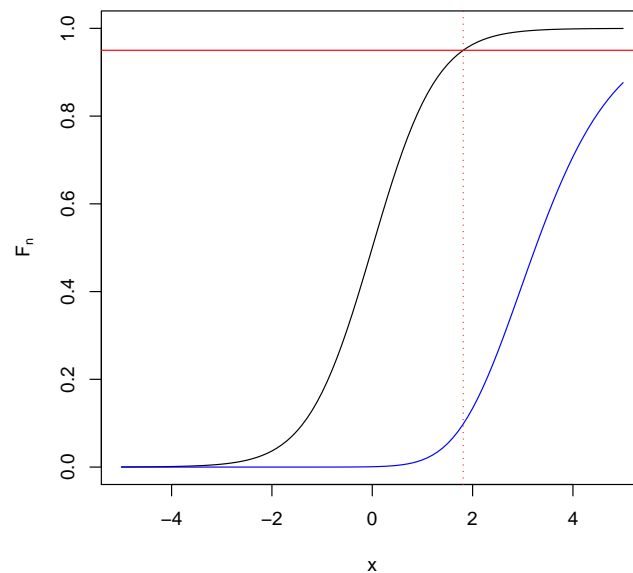
Beispiel 3.5:

Eingabe

```

n1<- 6; n2 <- 6
df <-  n1 + n2 -2
alpha <- 0.05
curve(pt(x,df=df),from=-5, to=5, ylab= expression(F[n]))
abline(h=1-alpha, col="red")      # cut at upper quantile
abline(v=qt(1-alpha, df=df), lty=3, col="red") # get critical value
n1 <- 5
n2 <- 5
n <- n1+n2
theta <- 2
ncp <- theta * sqrt(n1 * n2/(n1+n2))
curve(pt(x,df=df, ncp=ncp),add=TRUE, col="blue")

```



Die Güte des Tests können wir darstellen, indem wir die Verwerfungswahrscheinlichkeit in Abhängigkeit von $(\mu_1 - \mu_2)\sigma^{-1}$ auftragen.²

²Konventionell wird er die Gütefunktion nur auf der Alternative, d.h z.B. für $(\mu_1 - \mu_2) > 0$ betrachtet. Wir setzen sie hier auch auf der Hypothese, d. h. für $(\mu_1 - \mu_2) > 0$ fort.

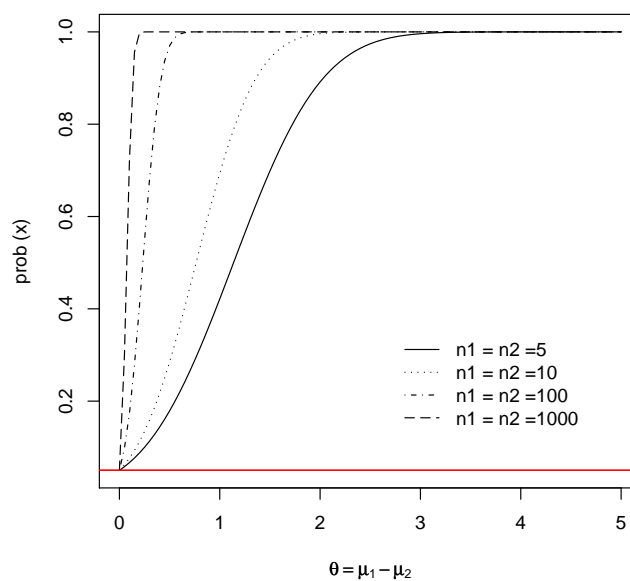
Beispiel 3.6:

```

tpower <- function(n1, n2, alpha, ...) {
  df <- n1 + n2 - 2
  tlim <- qt(1-alpha, df=df)
  prob <- function(theta) {
    pt(tlim, df = df,
       ncp = theta * sqrt(n1 * n2 / (n1+n2)),
       lower.tail=FALSE)
  }
  curve(prob, 0, 5, xlab=expression(theta==mu[1]-mu[2]), ...)
  abline(h=alpha, col="red")
}

tpower(5, 5, 0.05)
tpower(10, 10, 0.05, add =TRUE, lty = 3)
tpower(100, 100, 0.05, add =TRUE, lty = 4)
tpower(1000, 1000, 0.05, add =TRUE, lty = 5)
legend("bottomright",
      lty=c(1,3,4,5),
      legend=c("n1 = n2 =5", "n1 = n2 =10",
              "n1 = n2 =100", "n1 = n2 =1000"),
      inset=0.1, bty="n")

```



Der hier benutzte Zusammenhang kann auch benutzt werden, um zu bestimmen, wie groß der Stichprobenumfang sein muss, um auf der Hypothese höchstens mit einer Wahrscheinlichkeit α fälschlich zu verwerfen, bei Vorliegen einer spezifizierten Alternative jedoch mit einer gewählten Wahrscheinlichkeit die Hypothese richtigerweise zu verwerfen. Dazu gibt es die vorbereitete Funktion `power.t.test()`.

Beispiel 3.7:

<pre>power.t.test(delta=2, power=0.8, sig.level=0.01, type="two.sample", alternative="one.sided")</pre>	<hr style="border: 0.5px solid blue;"/> Eingabe <hr style="border: 0.5px solid blue;"/>
<hr style="border: 0.5px solid green;"/> Ausgabe <hr style="border: 0.5px solid green;"/>	
<pre>Two-sample t test power calculation n = 6.553292 delta = 2 sd = 1 sig.level = 0.01 power = 0.8 alternative = one.sided NOTE: n is number in *each* group</pre>	

3.4.2. Simulation der Güte. Sind die theoretischen Eigenschaften einer Test-Statistik bekannt, so ist dies der beste Weg, die Güte zu analysieren. In einer Umgebung wie **R** haben wir die Möglichkeit, die Güte auch dann zu untersuchen, wenn theoretische Resultate nicht vorliegen oder nicht zugänglich sind. Zu festgelegten Alternativen können wir Zufallsstichproben generieren, Tests durchführen und den relative Anteil der Verwerfungen bestimmen. Generieren wir $nsimul$ unabhängige Zufallsstichproben mit identischer Verteilung, so ist die Anzahl der Verwerfungen binomialverteilt und

$$\hat{p} = \frac{\#Verwerfungen}{nsimul}$$

ein Schätzer für die Verwerfungswahrscheinlichkeit.

Als Beispiel untersuchen wir, wie sich der t -Test verhält, wenn die Daten lognormal verteilt sind. Wir vergleichen zwei Gruppen jeweils mit Stichprobenumfang $n_1 = n_2 = 10$, zunächst auf der Hypothese:

Beispiel (Fortsetzung):

Beispiel 3.8:

	Eingabe
<pre> nsimul <- 300 n1<- 10; n2 <- 10 alpha <- 0.01 x <- 0 for (i in 1:nsimul) { if (t.test(exp(rnorm(n1)),exp(rnorm(n2))), alternative="less", var.equal = TRUE)\$p.value < alpha){ x <- x+1} } p <- x/nsimul cat("estim p", p) </pre>	
estim p 0.006666667	Ausgabe

Die Funktion `prop.test()` berechnet nicht nur diesen Schätzer, sondern auch einen Konfidenzbereich.

Beispiel 3.9:

<code>prop.test(n=nsimul, x=x)</code>	Eingabe
	Ausgabe
1-sample proportions test with continuity correction	
<pre> data: x out of nsimul, null probability 0.5 X-squared = 290.0833, df = 1, p-value < 2.2e-16 alternative hypothesis: true p is not equal to 0.5 95 percent confidence interval: 0.001155561 0.026512753 sample estimates: p 0.006666667 </pre>	

Analog z. B. wenn für die Alternative $\log(x_2)$ nach $N(\mu_2, 1)$ mit $\mu_2 = 1$ verteilt ist:

Beispiel 3.10:

```

Eingabe
-----
nsimul <- 300
n1<- 10; n2 <-10
alpha <- 0.01
x<-0
for (i in 1:nsimul) {
  if (t.test(exp(rnorm(n1)),exp(rnorm(n2, mean = 1)),
            alternative="less",
            var.equal = TRUE)$p.value < alpha){
    x <- x+1}
}
p <- x/nsimul
cat("estim p", p)

```

```

-----
Ausgabe
-----
estim p 0.1866667

```

```

Eingabe
-----
prop.test(n = nsimul, x = x)

```

```

-----
Ausgabe
-----
1-sample proportions test with continuity correction

data:  x out of nsimul, null probability 0.5
X-squared = 116.5633, df = 1, p-value < 2.2e-16
alternative hypothesis: true p is not equal to 0.5
95 percent confidence interval:
 0.1451407 0.2364119
sample estimates:
      p
0.1866667

```

In `library(binom)` finden sich eine Reihe von Werkzeugen für eine differenziertere Analyse der Binomialverteilung.

Die Konfidenzintervalle in diesem Beispiel zeige, dass hier ein Simulationsumfang von $nsimul = 300$ nur grobe Ergebnisse liefert. Für die Simulation wollen wir die Genauigkeit besser kontrollieren. Den Simulationsumfang können wir so wählen, dass eine gewünschte Genauigkeit erreicht werden kann. Eine genaue Planung können wir mit `power.prop.test()` machen. Für Simulationszwecke reicht oft schon eine Abschätzung.

Mit $\hat{p} := Z/n$ als Schätzer einer Wahrscheinlichkeit p haben wir $E(\hat{p}) = p$ und $Var(\hat{p}) = p(1-p)/n$. Ist p tatsächlich der interessierende Parameter, so sind Fehler relativ zum Zielparameter zu messen. Bei $p = 50\%$ ist ein Fehler von $\pm 1\%$ anders zu bewerten, als bei $p = 99\%$. Der relative Fehler, der **Variationskoeffizient**, ist

$$\frac{\sqrt{Var(\hat{p})}}{E(\hat{p})} = \sqrt{\frac{1-p}{np}}.$$

Um einen Variationskoeffizienten von höchstens η zu erhalten, brauchen wir eine Stichprobenumfang

$$n \geq \frac{1-p}{p\eta^2}.$$

Sind n und p in einer Größenordnung, bei der eine Normalapproximation gilt, so haben wir für ein Konfidenzniveau $1 - \alpha$ ein approximatives Konfidenzintervall mit Grenzen

$$\hat{p} \pm \phi_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

Soll das Konfidenzintervall eine Länge von ηp nicht überschreiten, so brauchen wir einen Stichprobenumfang

$$n \geq \frac{\phi_{1-\alpha/2}^2 (1-p)}{p\eta^2}.$$

Wie üblich ist eine Wahl von α zu treffen. Für z. B. $\alpha = 1\%$ mit $\phi_{1-\alpha/2} = 2.575829$ erhalten wir die Werte in Tabelle 3.26. Falls wir mit höheren Quantilen arbeiten, werden wir versuchen, die Fehler relativ zu $1 - p$ zu beschränken. Beispiele sind in Tabelle 3.26.

TABELLE 3.26. Erforderlicher Stichprobenumfang für zweiseitige Konfidenzintervalle mit relativer Länge $\leq \eta$

p	$1 - p$	$n(\alpha = 10\%)$ $\eta = 0.1$	$\eta = 0.01$	$n(\alpha = 1\%)$ $\eta = 0.1$	$\eta = 0.01$
0.500	0.500	271	27 055	663	66 349
0.250	0.750	812	81 166	1990	199 047
0.100	0.900	2435	243 499	5971	597 141
0.010	0.990	26 785	2 678 488	65 685	6 568 547
0.001	0.999	270 283	27 028 379	662 826	662 822 617

Zu merken sind die groben Zahlen: um mit 90% Konfidenz eine Wahrscheinlichkeit im Bereich von $50\% \pm 5\%$ zu schätzen, sind ca. 300 Simulationen notwendig. Um einen Wert bei 99% bis auf $\pm 0.1\%$ genau zu schätzen sind 30000 Simulationen nötig.

3.4.3. Quantilschätzung durch Simulation.

Die andere Seite des Problems ist es, ein Quantil anhand einer Stichprobe zu schätzen. Wir wissen bereits, dass für eine Zufallsvariable X mit stetiger Verteilungsfunktion F die Variable $F(X)$ eine uniforme Verteilung auf $[0, 1]$ hat. Für die Quantilschätzung benötigen wir die Verteilungsfunktion, ausgewertet an den Ordnungsstatistiken. Diese haben wir bereits in Kapitel ?? kennen gelernt. Dort hatten wir als Theorem ??.

THEOREM 3.2. *Sind $X_i, i = 1, \dots, n$ unabhängige Beobachtungen aus einer Verteilung mit stetiger Verteilungsfunktion F und ist $X_{(k:n)}$ die k . Ordnungsstatistik daraus, so ist*

$$F(X_{(k:n)}) \sim P_{beta}(\cdot; k, n - k + 1).$$

Wir wiederholen:

BEMERKUNG 3.3. Im allgemeinen ist die beta-Verteilung schief. Der Erwartungswert der $Beta(k, n - k + 1)$ -Verteilung ist $k/(n + 1)$. Um eine unverzerrte Schätzung des Quantils x_p zu erhalten, benutzt man $X_{(k:n)}$ mit $k/(n + 1) = p$. Die "plug in"-Approximation $k/n = p$ gibt eine verzerrte Schätzung.

Das Theorem kann direkt angewendet werden, um eine obere oder untere Abschätzung für Quantile zu gewinnen. Insbesondere können wir versuchen, das Minimum der beobachteten Wertes $X_{(1:n)}$ als untere Abschätzung für das p -Quantil. zu benutzen Das Konfidenzniveau ist

$$P(X_{(1)} \leq F_p) = P(F(X_{(k)}) \leq p) = I_p(1, n),$$

wobei I das unvollständige Beta-Integral ist. Für die speziellen Parameter $(1, n)$ vereinfacht sich die Beta-Dichte zu $n(1-p)^{n-1}$ und wir bekommen für das unvollständige Beta-Integral $I_p(1, n) = 1 - (1-p)^n$. Daraus folgt

$$P(X_{(1)} \leq F_p) = 1 - (1-p)^n$$

und wir können ein Konfidenzniveau $1 - \alpha$ sicherstellen, wenn

$$n \geq \frac{\ln \alpha}{\ln(1-p)}.$$

Der beobachtete Höchstwert kann als obere Abschätzung für das p -Quantil verwendet werden. Aus Symmetriegründen erhalten wir einen Konfidenzniveau von $1 - \alpha$, wenn

$$n \geq \frac{\ln \alpha}{\ln p}.$$

Beispiele sind in Tabelle 3.27 angegeben.

TABELLE 3.27. Benötigter Stichprobenumfang zur Abschätzung eines Quantils mit Konfidenzniveau $\geq 1 - \alpha$

p		n			
$X_{(1)} \leq F_p$	$X_{(n)} \geq F_p$	$\alpha = 10\%$	$\alpha = 5\%$	$\alpha = 1\%$	$\alpha = 0.5\%$
0.500	0.500	4	5	7	8
0.250	0.750	9	11	17	19
0.100	0.900	22	29	44	51
0.010	0.990	230	299	459	528
0.001	0.999	2302	2995	4603	5296

Zu merken sind wieder die groben Zahlen: um eine einseitige Abschätzung für ein 1% (99%)-Quantil einer stetigen Verteilungsfunktion mit einer Konfidenz von 99% zu erhalten, werden beinahe 500 Simulationen benötigt.

Wir können einseitige Schranken zu Intervallen verknüpfen. Das entsprechende Resultat zur Berechnung der Wahrscheinlichkeit von Intervallen ist in Korollar ??:

KOROLLAR 3.4. *Mit der k_1 -ten und $k_1 + k_2$ -ten Ordnungsstatistik ist das Intervall $(X_{(k_1:n)}, X_{(k_1+k_2:n)})$ ein Konfidenzintervall für das p -Quantil mit der Überdeckungswahrscheinlichkeit*

$$I_p(k_1, n - k_1 + 1) - I_p(k_1 + k_2, n - k_1 - k_2 + 1).$$

Die Simulationsumfänge zur Abschätzung von Quantilen sind drastisch geringer als diejenigen, die zur vergleichbaren Schätzung von Wahrscheinlichkeiten benötigt werden. Im Nachhinein ist dies nicht verwunderlich: die Frage, ob eine Beobachtung über einem bestimmten Quantil liegt, ist einfacher, als die Aufgabe, den p -Wert zu schätzen. In Abschnitt ?? werden wir sehen, dass der notwendige Stichprobenumfang noch einmal drastisch verringert werden kann, wenn die Fragestellung auf ein Testproblem reduziert wird.

Ohne weitere Verteilungsannahmen gibt dies eine erste Möglichkeit, den Umfang einer Simulation festzulegen. In speziellen Situationen können geschickte Einfälle eine bedeutende Reduzierung des Stichprobenumfangs erlauben. Zunächst aber sind die obigen Abschätzungen die Grundlage für Simulationen.

3.5. Qualitative Eigenschaften von Verteilungen

3.6. Ergänzungen

3.7. Statistische Zusammenfassung

Als Leitbeispiel diente in diesem Kapitel der Vergleich von Stichproben. In einfachen Fällen unterscheiden sich Stichproben nur um eine Verschiebung des Mittelwerts. In diesem Fall können die Probleme auf die Ansätze aus Kapitel 2 reduziert werden. In diesem reduzierten Fall stimmen die um den Mittelwert zentrierten Verteilungen überein. Für den allgemeineren Fall, den wir jetzt untersucht haben, gilt diese Vereinfachung nicht. Ein wichtiges Beispiel ist etwa die Untersuchung von Therapie-Studien. Hat eine Behandlung einen homogenen Effekt, so können wir diesen mit den Mitteln von Kapitel 2 untersuchen. Häufig aber gibt es unter einer Behandlung eine Aufspaltung in “Responder” und “Nicht-Responder”. Dies geht über die in Kapitel 2 skizzierten Modelle hinaus, und die allgemeineren Ansätze aus diesem Kapitel 3 werden nötig.

Wir haben uns hier auf den Vergleich von zwei Stichproben beschränkt. Die Praxis führt oft auf andere Probleme. So ist ein typischer Fall, dass eine neue Behandlung mit einer bekannten Referenz-Behandlung verglichen werden soll, wobei für die neue Behandlung nur eine Stichprobeninformation, für die Referenz-Behandlung aber umfassendere Vorinformation bereit steht. Oder eine Referenz-Behandlung soll mit einer Serie von Alternativ-Behandlungen verglichen werden. Diese Probleme gehen über den Rahmen unserer Einführung hinaus. Hier kann nur auf weiterführende Literatur, z.B [Mil81] verwiesen werden.

3.8. Literatur und weitere Hinweise:

[VR02] Venables, W.N.; Ripley, B.D. (2002): Modern Applied Statistics with S. Heidelberg: Springer

[VR00] Venables, W.N.; Ripley, B.D. (2000): S Programming. Heidelberg: Springer

[Mil81] Miller, R. G. (1981): Simultaneous Statistical Inference. Heidelberg: Springer

KAPITEL 4

Dimensionen 1, 2, 3, ..., ∞

4.1. Ergänzungen

In diesem Kapitel beginnen wir Ergänzungen zu R, um uns dann auf statistische Fragen zu konzentrieren. Für werfen einen Blick auf die graphischen Möglichkeiten, die uns zur Verfügung stehen. Die Basis-Graphik von R ist an einem Plotter-Modell orientiert. Die Graphik folgt den Möglichkeiten, die das Zeichnen mit einem Stift bietet. Neben den ein- und zweidimensionalen Möglichkeiten, die wir bis jetzt kennengelernt haben, gibt es Möglichkeiten, eine Funktion darzustellen, die über einem Raster definiert sind. Dazu stehen im wesentlichen drei Funktionen zur Verfügung.

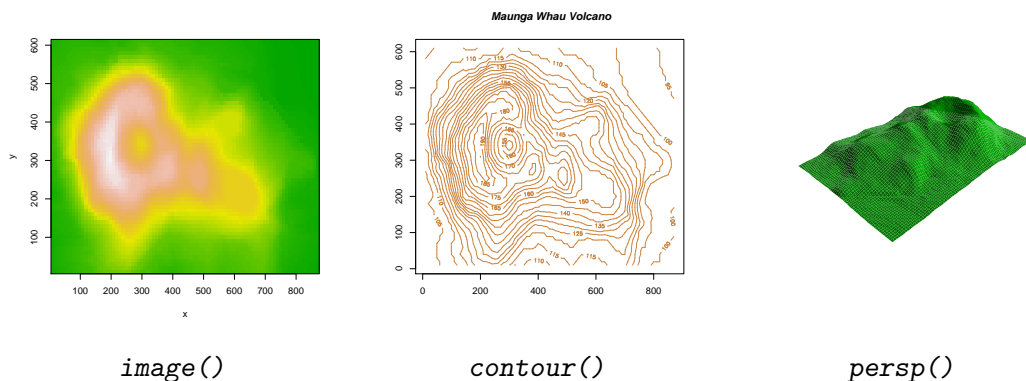
<i>3d-Graphik</i>	
<i>image()</i>	Gibt die Werte einer Variablen z in Graustufen oder Farbcodierung wieder.
<i>contour()</i>	Gibt die Kontouren einer Variablen z .
<i>persp()</i>	Gibt einen perspektivischen Plot einer Variablen z .

image() und *contour()* können auch benutzt werden, um andere Plots zu überlagern.

Beispiel 4.1:

Eingabe

```
#oldpar <- par(mfrow=c(1,3))
x <- 10*(1:nrow(volcano))
y <- 10*(1:ncol(volcano))
image(x, y, volcano, col = terrain.colors(100), axes = FALSE)
axis(1, at = seq(100, 800, by = 100))
axis(2, at = seq(100, 600, by = 100))
box()
title(main = "Maunga Whau Volcano", font.main = 4)
contour(x, y, volcano, levels = seq(90, 200, by = 5),
        col = "peru", main = "Maunga Whau Volcano", font.main = 4)
z <- 2 * volcano      # Exaggerate the relief
x <- 10 * (1:nrow(z)) # 10 meter spacing (S to N)
y <- 10 * (1:ncol(z)) # 10 meter spacing (E to W)
## Don't draw the grid lines : border = NA
#par(bg = "slategray")
persp(x, y, z, theta = 135, phi = 30, col = "green3", scale = FALSE,
      ltheta = -120, shade = 0.75, border = NA, box = FALSE)
```



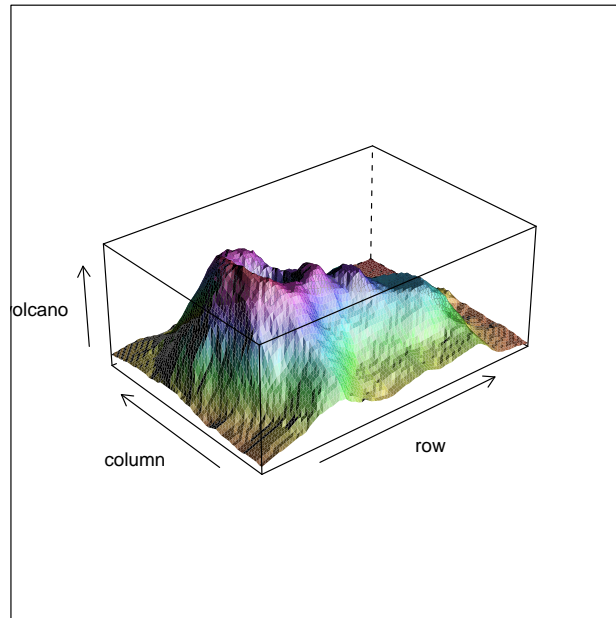
Die Basis-Graphik ist einfach zu handhaben, aber limitiert. Ein neues Grafiksystem arbeitet konzeptuell mit Objekten und einem Kamera-Modell. Die Grafik-Objekte können kombiniert und bearbeitet werden. Die Darstellung erfolgt in einem getrennten Schritt. Einfache 2d-Grafiken können hier nachbearbeitet werden. Für eine 3d-Darstellung können wie bei einer Kamera Abstand, Betrachtungsebene und Brennweite bestimmt werden. Das objektorientierte Graphiksystem besteht aus einer Bibliothek *grid* mit den notwendigen elementaren Operationen, und einer darauf aufbauende Bibliothek *lattice*, die die aus der Basis-Graphik bekannten Darstellungen neu implementiert und durch weitere ergänzt.

Lattice-Objekte werden mit `print()` ausgegeben.

Beispiel 4.2:

Eingabe

```
library(lattice)
## volcano ## 87 x 61 matrix
print(wireframe(volcano, shade = TRUE,
  aspect = c(61/87, 0.4),
  light.source = c(10,0,10)))
```



Basis-Graphik und Lattice-Graphik sind getrennte Graphik-Systeme. Leider benutzen sie auch für vergleichbare Funktionen unterschiedliche Bezeichnungen, und vergleichbare Displays haben unterschiedliche Darstellungen. Eine kleine Übersetzungshilfe ist in Tabelle 4.4 angegeben. Einige Hilfsfunktionen, um beide Graphik-Systeme in Kombination zu nutzen, werden durch die Bibliothek *gridBase* bereitgestellt. Eine ausführliche Einführung in beide Graphik-Systeme ist [Mur06].

Für Visualisierungen im weiten Spektrum von wissenschaftlichen Visualisierungen bis hin zu aufwendigen Spielen wird verbreitet OpenGL benutzt. Dessen Funktionen stehen auch in R durch die Bibliothek *rgl* zur Verfügung. Es gibt jedoch einen deutlichen Unterschied zwischen den üblichen Anforderungen an Graphik, und den speziellen Erfordernissen statistischer Graphik. Wenn es um die Darstellungen von Funktionen geht, ist statistische Graphik noch vergleichbar mit den Anforderungen der üblichen Analysis. Der kleine Unterschied ist, dass Funktionen in der Statistik häufig stückweise konstant oder nur stückweise stetig sind, während z.B. in der Analysis stetige oder sogar in differenzierbare Funktionen die Regel sind. Bei der Darstellung von Daten ändert sich die Situation drastisch. Statistische Daten sind üblicherweise diskret. Glattheitseigenschaften, die die Visualisierung analytischer Daten einfacher machen, fehlen bei statistischen Daten. Deshalb sind spezielle angepasste Visualisierungen nötig.

Basis-Graphik		Lattice
<i>barplot()</i>	bar chart	<i>barchart()</i>
<i>boxplot()</i>	box and whisker plot	<i>bwplot()</i>
	3 dimensional scatter plot	<i>cloud()</i>
<i>contour</i>	contour plot	<i>contourplot()</i>
<i>coplot</i>	conditional scatter plots	<i>contourplot()</i>
<i>curve(density())</i>	density estimator	<i>densityplot()</i>
<i>dotchartt()</i>	dot plot	<i>dotplot()</i>
<i>hist()</i>	dot plot	<i>histogram()</i>
<i>image()</i>	colour map plots	<i>spiom()</i>
	parallel coordinate plots	<i>parallel()</i>
<i>pairs()</i>	scatter plot matrices	<i>wireframe()</i>
<i>persp()</i>	three dimensional surface	<i>wireframe()</i>
<i>plot()</i>	scatter plot	<i>xyplot()</i>
<i>qqnorm()</i>	theoretical $Q - Q$ -plot	<i>qqmath()</i>
<i>qqplot()</i>	empirical $Q - Q$ -plot	<i>qq()</i>
<i>stripchart()</i>	one dimensional scatterplot	<i>stripplot()</i>

TABELLE 4.4. Basis-Graphik und Lattice

4.2. Dimensionen

Wenn wir von einer Dimension zu höheren Dimensionen gehen, gibt es sowohl für die theoretische Untersuchung als auch für die grafische Darstellung neue Herausforderungen. Die linearen Modelle können wieder als leitendes oder warnendes Beispiel dienen.

Die Herausforderungen können von ernsthaften Problemen stammen. So können Verteilungen auf höherdimensionalen Räumen selbst unter Regularitätsvoraussetzungen unüberschaubar komplex sein. Die Klassifikations- und Identifikationsprobleme für Funktionen und Räume aus Analysis und Geometrie geben einen Vorgeschmack davon, was bei der Untersuchung von Wahrscheinlichkeitsverteilungen zu bewältigen ist.

Daneben gibt es hausgemachte Probleme: Eigentore, die durch selbstgetroffene Wahlen erst erzeugt werden.

Ein Beispiel für hausgemachte Probleme kann an linearen Modellen illustriert werden. Die Interpretation des Gauß-Markoff-Schätzers als lineare Projektion zeigt, dass nur scheinbar Koeffizienten für einzelne Regressoren geschätzt werden. Eigentlich wird ein Vektor im von den Regressoren aufgespannten Raum geschätzt; die Zuordnung zu den einzelnen Regressoren ist dann nur lineare Algebra. Diese hängt nicht von dem Einfluss des einzelnen Regressors ab, sondern von der gemeinsamen Geometrie der Regressoren. Nur wenn die Regressoren eine Orthogonalbasis bilden, gibt es eine direkte Interpretation der Koeffizienten. Wird im linearen Modell die Liste der Regressoren z.B. dupliziert, so ändert sich der Raum nicht. Die Rechnungen in Koordinaten werden etwas komplizierter, weil die Regressoren nun auf keinen Fall eine Basis bilden, aber von einem abstrakten Standpunkt bleibt die Situation unverändert. Gibt es aber kein echtes Duplikat, sondern geringfügige Abweichungen (durch minimale "Fehler", Rundungen, Transformationen), so ändert sich

die Situation drastisch. Für den Gauß-Markoff-Schätzer ist nur der von den Regressoren aufgespannte Raum relevant, und selbst durch minimale Änderungen im Duplikat kann sich dessen Dimension verdoppeln. Dies ist ein Beispiel für ein hausgemachtes Problem.

Dies und andere Beispiele sind ein Grund, die Beziehungen zwischen den Variablen genauer zu untersuchen. Bei der Regression etwa betrifft dies nicht nur die Beziehung zwischen Respons und Regressor, sondern, wie durch das letzte Beispiel illustriert, auch die Beziehungen zwischen den Regressoren.

Um die Verbindung zu den Regressionsproblemen zu halten und auf die Erfahrungen in diesem Bereich zurückzugreifen, betten wir formal die Regressionsprobleme in einen allgemeineren Rahmen ein. Bei der Regression hatten wir eine herausgehobene Variable, die Respons Y , deren Verteilung in Abhängigkeit von den Werten der übrigen Variablen, der Regressoren X , modelliert werden sollte. Wir fassen jetzt Respons und Regressor zu einem Datenvektor $Z = (Y; X)$ zusammen und werden auch die gemeinsame Verteilung von Z diskutieren. Wir finden das Regressionsproblem in diesem allgemeineren Rahmen wieder: beim Regressionsproblem suchten wir nach einem Schätzer für die Mittelwertfunktion m im Modell

$$Y = m(X) + \varepsilon.$$

Im allgemeineren Rahmen berücksichtigen wir eine gemeinsame Verteilung von X und Y . Das Regressionsmodell wird damit zum Modell

$$Y = E(Y|X) + \varepsilon$$

und wir haben zunächst die Identifizierung $m(X) = E(Y|X)$.

Wenn wir tatsächlich am ursprünglichen Regressionsmodell interessiert sind, müssen wir weitere Arbeit leisten. Eine Schätzung des bedingten Erwartungswerts $E(Y|X)$ ist nicht dasselbe wie die Schätzung einer Regressionsfunktion $m(X)$. Bei dem Regressionsproblem haben wir keine Annahmen über die Verteilung von X gemacht. Um von $E(Y|X)$ (oder einem Schätzer dafür) auf $m(X)$ zurück zu schließen, müssen wir überprüfen, dass die Schätzung von Verteilungsannahmen über X unabhängig ist. Für unsere augenblicklichen Zwecke ist diese Unterscheidung aber nicht relevant. Wir können uns eine Ignoranz auf Zeit erlauben.

Der allgemeine Rahmen in diesem Kapitel ist also: wir untersuchen Daten $(Z_i)_{i=1,\dots,n}$, wobei die einzelnen Beobachtungen Werte in \mathbb{R}^q annehmen.

Haben wir im wesentlichen lineare Strukturen, so können wir oft auch höher-dimensionale Strukturen mit Methoden analysieren, die für eindimensionale Modelle entwickelt sind. Wir müssen die Methoden evtl. modifizieren oder iteriert anwenden. Sie helfen uns jedoch, die wesentlichen Merkmale zu erkennen. Sie versagen jedoch, wenn sich höhere Dimensionalität mit Nichtlinearität verbindet. Dann sind speziellere Methoden gefragt.

4.3. Selektionen

Ursprünglich bedeutet eine Selektion eine Auswahl von Beobachtungen. Für die grafische Darstellung wird die Selektion mit einer Ausprägung von Attributen (z.B. Farbe, Plot-Zeichen, Dicke) assoziiert. Alle Variablenwerte, die zu Beobachtungen in der Selektion gehören, werden mit diesen Attributen in dieser Ausprägung dargestellt. Dies ermöglicht es, die Verbindung (“*linking*”) der Selektion in verschiedenen Plots zu verfolgen. So können Selektionen helfen, Strukturen in verbundenen Plots, zu erkennen.

In der praktischen Datenanalyse werden die Selektionen variiert (“*brushing*”), um Beobachtungen zu zusammengehörigen Beobachtungsgruppen zusammen zu fassen. Dies

Name	Variable	Einheit, Bem.
rw	relatives Gewicht	
fpg	Plasma-Glukose (nach Fasten)	[mg/100 ml]
ga	Glukosespiegel integriert über 3 Stunden Toleranztest	[mg/100 ml \times h]
ina	Insulinspiegel integriert über 3 Stunden Toleranztest	[μ U/100 ml \times h]
sspg	Plasma-Glukose (steady state)	[mg/100 ml]
cc	Klassifikation	chemische, normale, offene Diabetes

TABELLE 4.5. Diabetes-Datensatz: Variable

ist ein wichtiges Werkzeug der interaktiven Datenanalyse. Selektionen werden, statistisch gesprochen, zur Modellwahl benutzt. Ihnen entspricht das Konzept der lokalen Modelle: anstelle ein für die Daten globales, möglicherweise sehr komplexes Modell zu benutzen, werden für jede Selektion möglicherweise einfachere Modelle bestimmt, die jeweils nur für die Beobachtungen aus dieser Selektion gelten.

Das Linking wird leider von R nicht direkt unterstützt. Wir müssen also jeweils selbst sicher stellen, dass Selektionen mit den entsprechenden Attributen dargestellt werden. Auch die Repräsentation von Selektionen ist in R nicht einheitlich. Bei Funktionsaufrufen können diese durch *selection*-Parameter realisiert sein, oder durch *group*-Variable, oder als Bedingung in einem Formelausdruck. Deshalb müssen wir uns in jedem Fall mit ad-hoc-Lösungen begnügen.

R bleibt weitgehend auf statische Selektionen beschränkt, so dass Brushing nur rudimentär möglich ist.

Selektionen werden im Zusammenhang bei den nachfolgenden Beispielen illustriert.

4.4. Projektionen

Als erstes Beispiel betrachten wir einen Datensatz aus einer Arbeit ([RM79]), in der unterschiedliche Diabetes-Arten untersucht worden. Der Datensatz ist zum Beispiel in *library(locfit)* verfügbar. Die Variablen umfassen Laborwerte zum Glukose-Stoffwechsel und sind in Tabelle 4.5 erklärt.

Eingabe

```
library(locfit)
data(chemdiab)
```

Eine erste Übersicht erhalten wir mit

```
summary(chemdiab)
```

```

      rw          fpg          ga          ina
Min.   :0.7100  Min.   : 70.0  Min.   : 269.0  Min.   : 10.0
1st Qu.:0.8800  1st Qu.: 90.0  1st Qu.: 352.0  1st Qu.:118.0

```

```

Median :0.9800   Median : 97.0   Median : 413.0   Median :156.0
Mean    :0.9773   Mean    :122.0   Mean    : 543.6   Mean    :186.1
3rd Qu.:1.0800   3rd Qu.:112.0   3rd Qu.: 558.0   3rd Qu.:221.0
Max.    :1.2000   Max.    :353.0   Max.    :1568.0   Max.    :748.0
      sspg                                cc
Min.    : 29.0   Chemical_Diabetic:36
1st Qu.:100.0   Normal           :76
Median  :159.0   Overt_Diabetic  :33
Mean    :184.2
3rd Qu.:257.0
Max.    :480.0

```

Wie in der Originalarbeit lassen wir das relative Gewicht außer Betracht. Die chemische Klassifikation *cc* ist aus den Stoffwechseldaten abgeleitet. Sie beinhaltet also keine eigene Information. Zur Orientierung benutzen wir sie dennoch als Markierung, d.h. wir benutzen die Selektion *cc = Chemical_Diabetic, Normal, Overt_Diabetic*. Der Kern des Datensatzes ist vierdimensional mit den Variablen *fpg, ga, ina, sspg*.

4.4.1. Randverteilungen und Scatterplot-Matrix. Wir können versuchen, die mehrdimensionale Verteilung zu untersuchen, indem wir uns die zweidimensionalen *Marginalverteilungen* (Randverteilungen) für alle Variablenpaare ansehen. Die grafische Darstellung dazu heißt *Scatterplot-Matrix*, in R als Funktion *pairs()* implementiert.

[help\(pairs\)](#)

<code>pairs</code>	<i>Scatterplot Matrices</i>
--------------------	-----------------------------

Description.

A matrix of scatterplots is produced.

Usage.

```
pairs(x, ...)
```

```
## S3 method for class 'formula':
pairs(formula, data = NULL, ..., subset,
      na.action = stats::na.pass)
```

Default S3 method:

```
pairs(x, labels, panel = points, ...,
      lower.panel = panel, upper.panel = panel,
      diag.panel = NULL, text.panel = textPanel,
      label.pos = 0.5 + has.diag/3,
      cex.labels = NULL, font.labels = 1,
      rowlattice = TRUE, gap = 1)
```

Arguments.

<code>x</code>	the coordinates of points given as numeric columns of a matrix or dataframe. Logical and factor columns are converted to numeric in the same way that <code>data.matrix</code> does.
<code>formula</code>	a formula, such as <code>~ x + y + z</code> . Each term will give a separate variable in the pairs plot, so terms should be numeric vectors. (A response will be interpreted as another variable, but not treated specially, so it is confusing to use one.)
<code>data</code>	a <code>data.frame</code> (or list) from which the variables in <code>formula</code> should be taken.
<code>subset</code>	an optional vector specifying a subset of observations to be used for plotting.
<code>na.action</code>	a function which indicates what should happen when the data contain NAs. The default is to pass missing values on to the panel functions, but <code>na.action = na.omit</code> will cause cases with missing values in any of the variables to be omitted entirely.
<code>labels</code>	the names of the variables.
<code>panel</code>	<code>function(x,y,...)</code> which is used to plot the contents of each panel of the display.
<code>...</code>	arguments to be passed to or from methods. Also, graphical parameters can be given as can arguments to <code>plot</code> such as <code>main</code> . <code>par("oma")</code> will be set appropriately unless specified.
<code>lower.panel, upper.panel</code>	separate panel functions to be used below and above the diagonal respectively.
<code>diag.panel</code>	optional <code>function(x, ...)</code> to be applied on the diagonals.
<code>text.panel</code>	optional <code>function(x, y, labels, cex, font, ...)</code> to be applied on the diagonals.
<code>label.pos</code>	y position of labels in the text panel.
<code>cex.labels, font.labels</code>	graphics parameters for the text panel.
<code>row1atop</code>	logical. Should the layout be matrix-like with row 1 at the top, or graph-like with row 1 at the bottom?
<code>gap</code>	Distance between subplots, in margin lines.

Details.

The ij th scatterplot contains `x[,i]` plotted against `x[,j]`. The “scatterplot” can be customised by setting panel functions to appear as something completely different. The off-diagonal panel functions are passed the appropriate columns of `x` as `x` and `y`: the diagonal panel function (if any) is passed a single column, and the `text.panel` function is passed a single `(x, y)` location and the column name.

The graphical parameters `pch` and `col` can be used to specify a vector of plotting symbols and colors to be used in the plots.

The graphical parameter `oma` will be set by `pairs.default` unless supplied as an argument.

A panel function should not attempt to start a new plot, but just plot within a given coordinate system: thus `plot` and `boxplot` are not panel functions.

By default, missing values are passed to the panel functions and will often be ignored within a panel. However, for the formula method and `na.action = na.omit`, all cases which contain a missing values for any of the variables are omitted completely

(including when the scales are selected). (The latter was the default behaviour prior to R 2.0.0.)

Author(s).

Enhancements for R 1.0.0 contributed by Dr. Jens Oehlschlaegel-Akiyoshi and R-core members.

References.

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

Examples.

```

pairs(iris[1:4], main = "Anderson's Iris Data -- 3 species",
      pch = 21, bg = c("red", "green3", "blue")[unclass(iris$Species)])

## formula method
pairs(~ Fertility + Education + Catholic, data = swiss,
      subset = Education < 20, main = "Swiss data, Education < 20")

pairs(USJudgeRatings)

## put histograms on the diagonal
panel.hist <- function(x, ...)
{
  usr <- par("usr"); on.exit(par(usr))
  par(usr = c(usr[1:2], 0, 1.5) )
  h <- hist(x, plot = FALSE)
  breaks <- h$breaks; nB <- length(breaks)
  y <- h$counts; y <- y/max(y)
  rect(breaks[-nB], 0, breaks[-1], y, col="cyan", ...)
}
pairs(USJudgeRatings[1:5], panel=panel.smooth,
      cex = 1.5, pch = 24, bg="light blue",
      diag.panel=panel.hist, cex.labels = 2, font.labels=2)

## put (absolute) correlations on the upper panels,
## with size proportional to the correlations.
panel.cor <- function(x, y, digits=2, prefix="", cex.cor)
{
  usr <- par("usr"); on.exit(par(usr))
  par(usr = c(0, 1, 0, 1))
  r <- abs(cor(x, y))
  txt <- format(c(r, 0.123456789), digits=digits)[1]
  txt <- paste(prefix, txt, sep="")
  if(missing(cex.cor)) cex <- 0.8/strwidth(txt)
  text(0.5, 0.5, txt, cex = cex * r)
}
pairs(USJudgeRatings, lower.panel=panel.smooth, upper.panel=panel.cor)

```

Wir benutzen die chemischen Diabetes-Klassen *cc* als Selektionen. Jeder dieser Selektionen wird ein Farbwert zugeordnet; dies ist das verbindende Attribut, das ermöglicht,

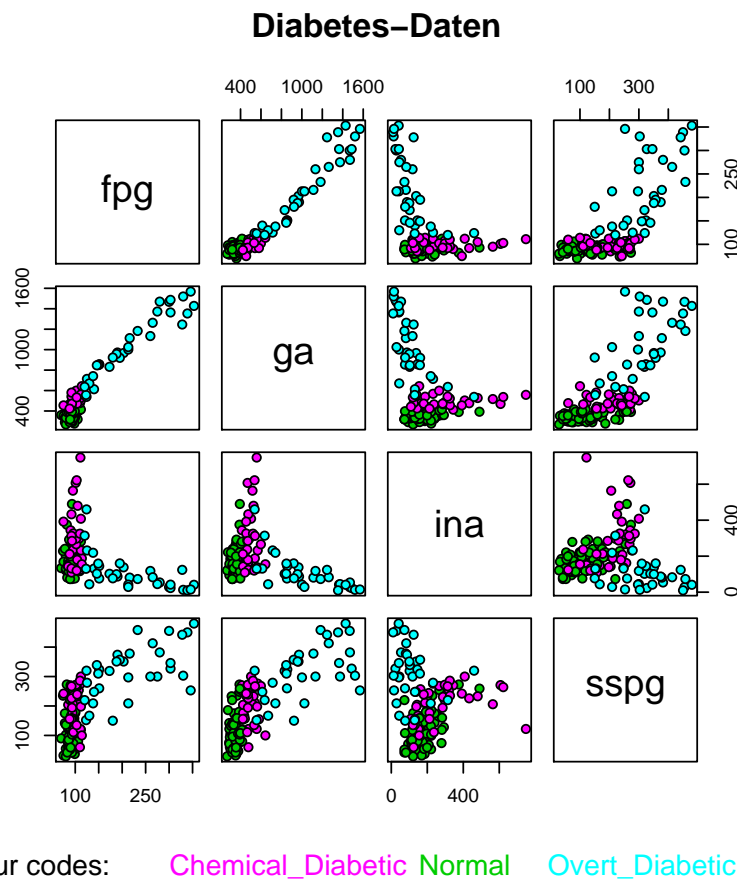
Verbindungen zwischen den Plots zu verfolgen. Um diese Verbindung zu dokumentieren, müssen wir in die Grafiksteuerung eingreifen und die Plots modifizieren. Wir erzeugen mit dem Parameter *oma* einen äusseren Rand, in dem wir eine Legende platzieren.

Beispiel 4.3:

```

pairs(~fpg + ga + ina + sspg, data = chemdiab, pch = 21,
      main = "Diabetes-Daten",
      bg = c("magenta", "green3", "cyan")[unclass(chemdiab$cc)],
      oma = c(8, 8, 8, 8))
mtext(c("Colour codes:", levels(chemdiab$cc)),
      col = c("black", "magenta", "green3", "cyan"),
      at = c(0.1, 0.4, 0.6, 0.8), side = 1, line = 2)

```



Die Funktion *pairs()* kontrolliert nur das "Layout" der Matrix, die Auswahl und Anordnung der Projektionen. Die Darstellung in den Plot-Feldern kann durch den Aufruf gesteuert werden. Die Default-Belegungen führen dazu, dass in der Diagonale die Namen der Variablen und außerhalb der Diagonalen die paarweisen Scatterplots gezeigt werden.

Aufgabe 4.1	
	Generieren Sie eine Scatterplot-Matrix für den Diabetes-Datensatz, die in der Diagonale die Histogramme der jeweiligen Variablen zeigt. <i>Hinweis:</i> Siehe <code>help(pairs)</code> .

Bestimmte Aspekte der Verteilung können aus den Randverteilungen einfach abgelesen werden. Andere geometrische Strukturen sind aus Randverteilungen gar nicht oder nur schwer zu rekonstruieren.

Zwischen dem Glukose-Spiegel bei Fasten `fpg` und dem integrierten Glukose- Spiegel bei Belastung `ga` besteht z.B. ein deutlicher linearer Zusammenhang. Dieser ist in den zweidimensionalen Marginalverteilungen erkennbar und kann mit den Methoden für lineare Modelle untersucht werden.

Diese deutliche Beziehung pflanzt sich auf die Beziehungen zu den anderen Variablen `ina`, `sspg` fort. In der Originalarbeit wird deshalb `fpg` nicht weiter berücksichtigt. Zu untersuchen sind noch die Variablen `ga`, `ina`, `sspg`. Die dreidimensionale Struktur dieses Teils des Datensatzes ist aus den Marginalverteilungen nicht einfach abzulesen.

4.4.2. Projection Pursuit. Geometrische Beziehungen oder stochastische Abhängigkeiten, die nicht parallel zu den Koordinaten-Achsen ausgerichtet sind, werden durch die Randverteilungen nicht ausgedrückt. Wir können die Idee verallgemeinern und anstelle von zweidimensionale Marginalverteilungen beliebige Projektionen benutzen. Dazu greifen wir auf `library(lattice)` zu. Darin ist ein an einer Kamera orientiertes Grafik-Modell implementiert.

Die `grid`-Grafik liefert mit dem Paket `lattice` eine weitgehende Unterstützung für multivariate Darstellungen. `grid` ist dabei die Basis. Das ursprüngliche Grafiksistem von R implementiert ein Modell, dass an der Vorstellung von Stift und Papier orientiert ist. Ein Grafik-Port (Papier) wird eröffnet und darauf werden Linien, Punkte/Symbole gezeichnet. `grid` ist ein zweites Grafiksistem, das an einem Kamera/Objekt-Modell orientiert ist. Grafische Objekte in unterschiedlicher Lage und Richtung werden in einem visuellen Raum abgebildet. Auf der `grid` baut `lattice` auf. In <http://cm.bell-labs.com/cm/ms/departments/sia/project/trellis/> sind die Grundideen zur Visualisierung multi-dimensionaler Daten dokumentiert, die in `lattice` implementiert sind.

Die erzeugte Grafik wird mit `print()` ausgegeben. Mit dem Parameter `split` können wir den Ausgabebereich aufteilen. Leider ist das Linking hier gebrochen: `cloud()` kann zwar eine Legende erzeugen, diese zeigt jedoch die Farbskala bei Beginn der Grafik, nicht die bei der Ausgabe benutzte. Wir müssen deshalb wieder ins System eingreifen und diesmal die Farbtabelle ändern.

```

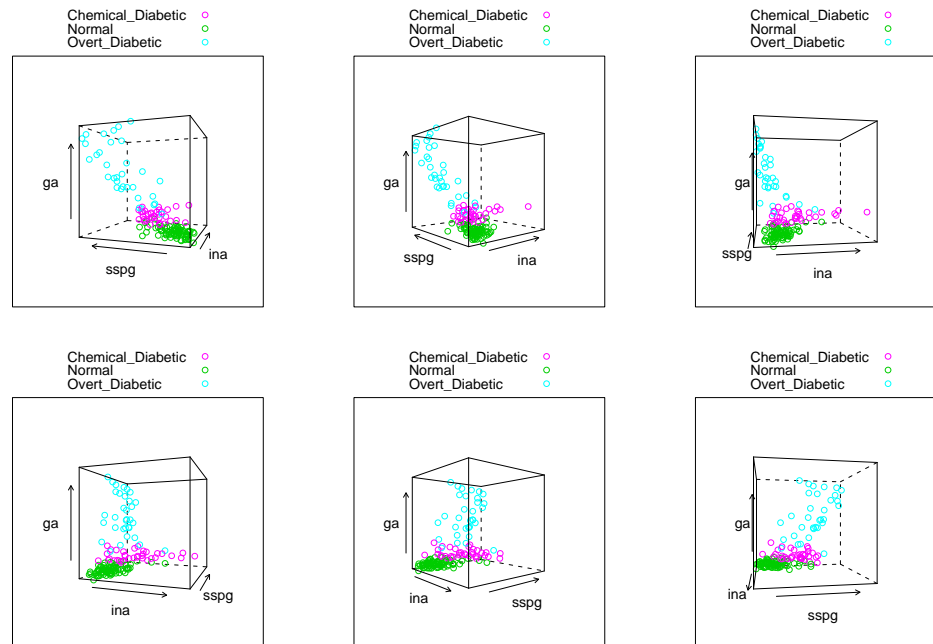
library("lattice")
diabcloud <- function(y, where, more = TRUE, ...) {
  print(cloud(ga ~ ina + sspg, data = chemdiab, groups = cc,
             screen = list(x = -90, y = y), distance = .4, zoom = .6,
             auto.key = TRUE, ...),
        split = c(where, 3, 2), more = more)
}
supsym <- trellis.par.get("superpose.symbol")
supsymold <- supsym
supsym$col = c("magenta", "green3", "cyan")

```

```

trellis.par.set("superpose.symbol" = supsym)
diabcloud(y = 70, where = c(1, 1))
diabcloud(y = 40, where = c(2, 1))
diabcloud(y = 10, where = c(3, 1))
diabcloud(y = -20, where = c(1, 2))
diabcloud(y = -50, where = c(2, 2))
diabcloud(y = -80, where = c(3, 2), more = FALSE)
trellis.par.set("superpose.symbol" = supsymold)
rm(diabcloud, supsymold, supsym)

```



Aufgabe 4.2	
	<p>Modifizieren Sie dieses Beispiel so, dass Sie einen Eindruck der dreidimensionalen Struktur bekommen.</p> <p>Wie unterscheidet sich offene Diabetes von chemischer Diabetes?</p> <p>Wie verhält sich die Normal-Gruppe zu den beiden Diabetes-Gruppen?</p>

Auch mit Serien von Projektionen ist es oft nicht einfach, eine dreidimensionale Struktur zu identifizieren. Mit animierten Folgen kann dies einfacher sein. Unterstützung dazu findet sich in *library(rggobi)*, die allerdings *ggobi*, zu finden in <http://www.ggobi.org/>, als zusätzliche Software voraussetzt.

Was hier ad-hoc gemacht wird, kann auch systematisch durchgeführt werden und für beliebige Dimensionen verallgemeinert werden: man sucht für einen Datensatz im \mathbb{R}^q nach "interessanten" Projektionen. Dazu definiert man einen Index, der messen soll, wie interessant eine Projektion ist, und lässt dann eine Suche laufen, die diesen Index maximiert. Die auf dieser Idee basierende Familie statistischer Verfahren findet man unter dem Stichwort *projection pursuit*. Das System *ggobi* beinhaltet Implementierungen von projection

pursuit für eine Reihe von Indizes, die über die Funktionen in `library(rggobi)` von R aus angesprochen werden können.

4.4.3. Projektionen für dim 1, 2, 3, ... 7. Projektionsmethoden versuchen, in einem höherdimensionalen Datensatz Strukturen niedrigerer Dimension zu identifizieren. Die identifizierbare Dimension ist dabei beschränkt. Projizieren wir Struktur mit einer Dimension, die größer ist als das Projektionsziel, so überdeckt die typische Projektion alles, gibt also keine Information mehr.

Wieviele Dimension können wir erfassen? Die grafische Darstellung in der Ebene gibt zunächst eine niedrig angesetzte Grenze von zwei Dimensionen, d.h. zweidimensionale Strukturen können wir direkt mit cartesischen Koordinaten in der xy -Ebene darstellen. Die Wahrnehmung kann dreidimensionale Strukturen anhand von Hinweisen auf die Raumtiefe (etwa durch Schatten) oder aus Folgen von 2d-Bildern rekonstruieren. Mit Animationen erhalten wir einen Eindruck von veränderlichen 3d-Folgen und sind damit bei vier Dimensionen.

Mit zusätzlichen Informationskanälen wie z.B. mit Farbcodierungen können wir dies leicht erhöhen, bleiben aber effektiv bei vier bis sieben Dimensionen für ein Display.

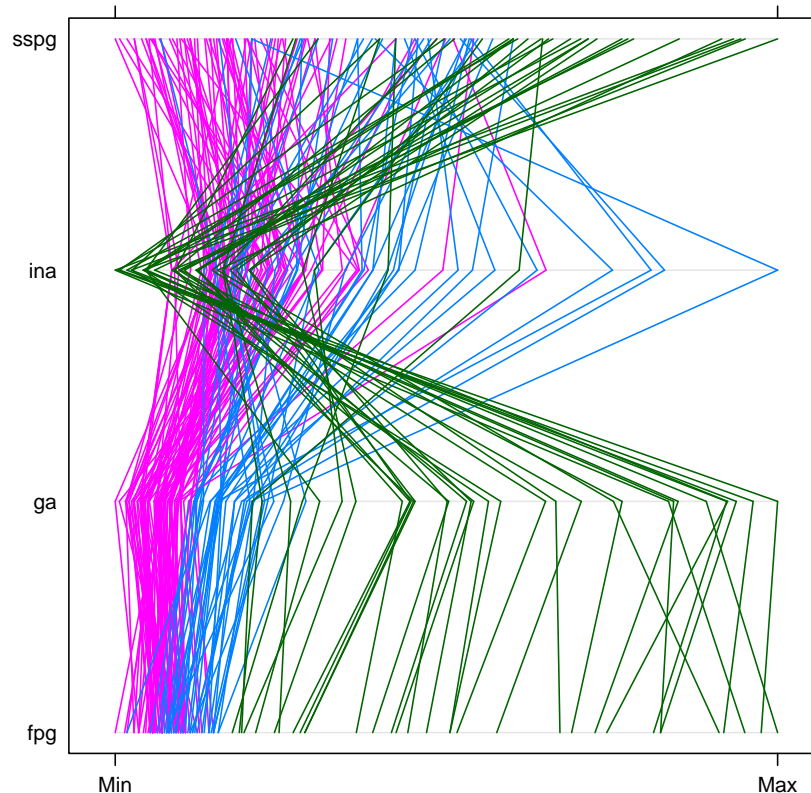
Die Kombination mehrerer Displays hilft kaum über diese Grenze hinaus. Stellen wir mehrere Displays z.B. in einer Scatterplot-Matrix dar, so verlieren wir die Fähigkeit, für die einzelnen Szenen durch die Wahrnehmung komplexere Strukturen zu generieren. Anstelle dessen müssen wir aktiv durch Vergleichen die komplexeren Strukturen aus den zweidimensionalen Displays erarbeiten. Die Fähigkeit, simultane Vergleiche durchzuführen, ist dabei beschränkt. Ebenso die Anzahl von Displays, die simultan auf einem Seiten-Medium wie Bildschirm oder Papier dargestellt werden kann.

4.4.4. Parallel-Koordinaten. Die grafische Darstellung (in kartesischen Koordinaten) sind zunächst auf ein- und zweidimensionale Projektionen beschränkt. Aber selbst bei Darstellungen in der Ebene ist die Beschränkung auf zwei Dimensionen nicht vorgegeben, sondern ist eine Folge unserer Wahl der Darstellung in kartesischen Koordinaten. Plot-Matrizen durchbrechen diese Dimensionsschranke durch Kombination von kartesischen Koordinatensystemen.

Parallel-Koordinaten orientieren die Achsen für die Variablen parallel zueinander. Für Häufigkeiten bei kategorialen Variablen ist dies eine übliche Darstellung: (evtl. überlagerte) Balkendiagramme benutzen Parallel-Koordinaten. Die Prozedur `parallel()` in `library(lattice)` unterstützt Parallel-Koordinaten auch für quantitative Variable. Die zu einem Fall gehörenden Markierungen auf diesen Achsen werden durch einen Linienzug verbunden. Diese Form der Parallel-Koordinaten stammt von A. Inselberg ([ICR87]).

Eingabe

```
library("lattice")
print(parallel(chemdiab[2:5], groups = chemdiab$cc))
```



Die Information ist dieselbe wie in den vorausgehenden Grafiken. Durch die veränderte Darstellung werden die Zusammenhänge auf neue Weise zugänglich.

Aufgabe 4.3	
	Notieren Sie für den <i>chemdiab</i> -Datensatz (schriftlich!) die Beziehungen zwischen den Variablen, die Sie im Parallel-Koordinatenplot erkennen können.
	Anstelle von <i>chemdiab[2:5]</i> können Sie die Variablen auch explizit als <i>chemdiab[c(2, 3, 4, 5)]</i> angeben. Durch dieser Form erhalten Sie Kontrolle über die Reihenfolge der Variablen. Vergleichen Sie zwei unterschiedliche Anordnungen der Variablen und notieren Sie (schriftlich!) ihre Beobachtungen. Welche Variablen-Anordnung gibt die einfachere Darstellung? Welche Beziehungen zwischen den Variablen sind in beiden ablesbar? Welche nur in einer der Anordnungen?

4.5. Schnitte, bedingte Verteilungen und Coplots

Schnitte sind, abstrakt gesehen, bedingte Verteilungen des Typs $P(\cdot \mid X = x)$. Sie sind nur dort zuverlässig, wo der Schnitt eine Bedingung definiert, die ein positives Maß hat. Um die Idee der Reduktion auf bedingte Verteilungen auch auf Daten anwenden zu

können, dicken wir die Schnitte auf. Anstelle bedingter Verteilungen des Typs $P(\cdot \mid X = x)$ zu untersuchen, betrachten wir $P(\cdot \mid \|X - x\| < \varepsilon)$, wobei ε auch mit x variieren kann. In grafischen Darstellungen von Daten verlangt dies eine Serie von Plots, die jeweils nur den durch die Bedingung eingeschränkten Teildatensatz zeigen.

Statistisch führen Projektionen zu Marginalverteilungen und Schnitte zu bedingten Verteilungen. Schnitte und Projektionen sind in gewissem Sinne komplementär: Projektionen zeigen Strukturmerkmale niedriger Dimension. Schnitte sind geeignet, Strukturmerkmale niedriger Co-Dimension zu entdecken. Beide können zur Datenanalyse kombiniert werden. Das Wechselspiel von Projektionen und Schnitten ist in [FB94] untersucht.

Wie die Dimensionsgrenzen bei der Projektion gibt es Grenzen für die Co-Dimension bei den Schnitten. Wir können nur Strukturen kleiner Co-Dimension erfassen. Ist die Co-Dimension zu groß, so ist ein typischer Schnitt leer, gibt also keine Information.

Als erstes Hilfsmittel stellt R die Möglichkeit bereit, zwei Variablen *bedingt* auf eine oder mehrere weitere Variable zu analysieren. Als grafische Darstellung dient dazu der **Coplot**. Er ist eine Variante der Plot-Matrix und zeigt in jedem Feld den Scatterplot zweier Variabler, gegeben die Bedingung.

Der Coplot kann nun auf bestimmte Muster untersucht werden. Sind die dargestellten Variablen stochastisch unabhängig von den bedingenden Variablen, so zeigen alle Plot-Elemente dieselbe Gestalt. Dargestellte Variable und bedingende Variable können dann entkoppelt werden.

Stimmt die Gestalt überein, aber Ort und Größe variieren, so weist dies auf eine (nicht notwendig lineare) Shift/Skalenbeziehung hin. Additive Modelle oder Varianten davon können benutzt werden, um die Beziehung zwischen dargestellten Variablen und bedingenden Variablen zu modellieren.

Verändert sich bei Variation der Bedingung die Gestalt, so liegt eine wesentliche Abhängigkeitsstruktur oder Interaktion vor, die genauerer Modellierung bedarf.

help(coplot)

coplot

Conditioning Plots

Description.

This function produces two variants of the **conditioning plots discussed in the reference below**.

Usage.

```
coplot(formula, data, given.values, panel = points, rows, columns,
       show.given = TRUE, col = par("fg"), pch = par("pch"),
       bar.bg = c(num = gray(0.8), fac = gray(0.95)),
       xlab = c(x.name, paste("Given :", a.name)),
       ylab = c(y.name, paste("Given :", b.name)),
       subscripts = FALSE,
       axlabels = function(f) abbreviate(levels(f)),
       number = 6, overlap = 0.5, xlim, ylim, ...)
co.intervals(x, number = 6, overlap = 0.5)
```

Arguments.

<code>formula</code>	a formula describing the form of conditioning plot. A formula of the form <code>y ~ x a</code> indicates that plots of <code>y</code> versus <code>x</code> should be produced conditional on the variable <code>a</code> . A formula of the form <code>y ~ x a * b</code> indicates that plots of <code>y</code> versus <code>x</code> should be produced conditional on the two variables <code>a</code> and <code>b</code> . All three or four variables may be either numeric or factors. When <code>x</code> or <code>y</code> are factors, the result is almost as if <code>as.numeric()</code> was applied, whereas for factor <code>a</code> or <code>b</code> , the conditioning (and its graphics if <code>show.given</code> is true) are adapted.
<code>data</code>	a data frame containing values for any variables in the formula. By default the environment where <code>coplot</code> was called from is used.
<code>given.values</code>	a value or list of two values which determine how the conditioning on <code>a</code> and <code>b</code> is to take place. When there is no <code>b</code> (i.e., conditioning only on <code>a</code>), usually this is a matrix with two columns each row of which gives an interval, to be conditioned on, but it can also be a single vector of numbers or a set of factor levels (if the variable being conditioned on is a factor). In this case (no <code>b</code>), the result of <code>co.intervals</code> can be used directly as <code>given.values</code> argument.
<code>panel</code>	a <code>function(x, y, col, pch, ...)</code> which gives the action to be carried out in each panel of the display. The default is <code>points</code> .
<code>rows</code>	the panels of the plot are laid out in a <code>rows</code> by <code>columns</code> array. <code>rows</code> gives the number of rows in the array.
<code>columns</code>	the number of columns in the panel layout array.
<code>show.given</code>	logical (possibly of length 2 for 2 conditioning variables): should conditioning plots be shown for the corresponding conditioning variables (default <code>TRUE</code>)
<code>col</code>	a vector of colors to be used to plot the points. If too short, the values are recycled.
<code>pch</code>	a vector of plotting symbols or characters. If too short, the values are recycled.
<code>bar.bg</code>	a named vector with components <code>"num"</code> and <code>"fac"</code> giving the background colors for the (shingle) bars, for numeric and factor conditioning variables respectively .
<code>xlab</code>	character; labels to use for the x axis and the first conditioning variable. If only one label is given, it is used for the x axis and the default label is used for the conditioning variable.
<code>ylab</code>	character; labels to use for the y axis and any second conditioning variable.
<code>subscripts</code>	logical: if true the panel function is given an additional (third) argument <code>subscripts</code> giving the subscripts of the data passed to that panel.
<code>axlabels</code>	function for creating axis (tick) labels when x or y are factors.
<code>number</code>	integer; the number of conditioning intervals, for a and b, possibly of length 2. It is only used if the corresponding conditioning variable is not a factor.
<code>overlap</code>	numeric < 1; the fraction of overlap of the conditioning variables, possibly of length 2 for x and y direction. When <code>overlap < 0</code>, there will be <i>gaps</i> between the data slices.

`xlim` **the range for the x axis.**
`ylim` **the range for the y axis.**
`...` **additional arguments to the panel function.**
`x` **a numeric vector.**

Details.

In the case of a single conditioning variable `a`, when both `rows` and `columns` are unspecified, a “close to square” layout is chosen with `columns >= rows`.

In the case of multiple `rows`, the *order* of the panel plots is from the bottom and from the left (corresponding to increasing `a`, typically).

A panel function should not attempt to start a new plot, but just plot within a given coordinate system: thus `plot` and `boxplot` are not panel functions.

As from R 2.0.0 the rendering of arguments `xlab` and `ylab` is not controlled by `par` arguments `cex.lab` and `font.lab` even though they are plotted by `mtext` rather than `title`.

Value.

`co.intervals(., number, .)` returns a (`number × 2`) matrix, say `ci`, where `ci[k,]` is the range of `x` values for the `k`-th interval.

References.

Chambers, J. M. (1992) *Data for models*. Chapter 3 of *Statistical Models in S* eds J. M. Chambers and T. J. Hastie, Wadsworth & Brooks/Cole.

Cleveland, W. S. (1993) *Visualizing Data*. New Jersey: Summit Press.

See Also.

`pairs`, `panel.smooth`, `points`.

Examples.

```
## Tonga Trench Earthquakes
coplot(lat ~ long | depth, data = quakes)
given.depth <- co.intervals(quakes$depth, number=4, overlap=.1)
coplot(lat ~ long | depth, data = quakes, given.v=given.depth, rows=1)

## Conditioning on 2 variables:
ll.dm <- lat ~ long | depth * mag
coplot(ll.dm, data = quakes)
coplot(ll.dm, data = quakes, number=c(4,7), show.given=c(TRUE,FALSE))
coplot(ll.dm, data = quakes, number=c(3,7),
       overlap=c(-.5,.1)) # negative overlap DROPS values

## given two factors
Index <- seq(length=nrow(warpbreaks)) # to get nicer default labels
coplot(breaks ~ Index | wool * tension, data = warpbreaks, show.given = 0:1)
coplot(breaks ~ Index | wool * tension, data = warpbreaks,
       col = "red", bg = "pink", pch = 21, bar.bg = c(fac = "light blue"))

## Example with empty panels:
with(data.frame(state.x77), {
  coplot(Life.Exp ~ Income | Illiteracy * state.region, number = 3,
```

```

    panel = function(x, y, ...) panel.smooth(x, y, span = .8, ...)
## y ~ factor -- not really sensical, but 'show off':
coplot(Life.Exp ~ state.region | Income * state.division,
       panel = panel.smooth)
})

```

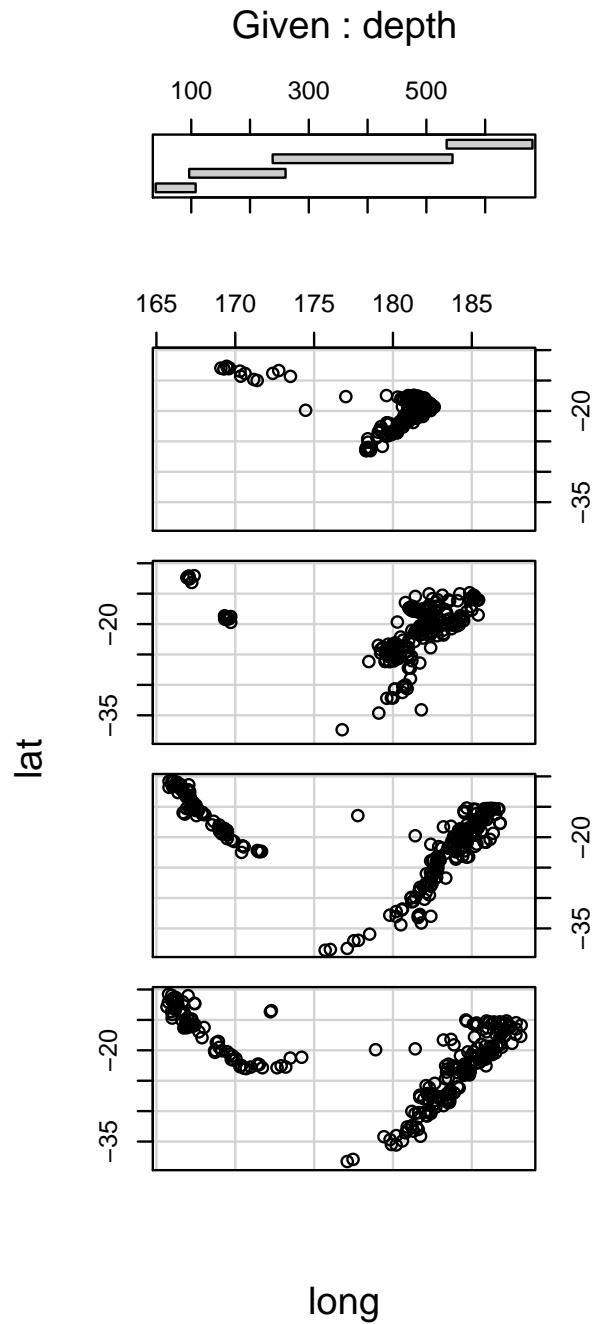
Wir illustrieren die Coplots mit dem “Quakes”-Datensatz. Dieser Datensatz gibt die geografische Länge und Breite einer Reihe von Erdbeben in der Nähe der Fiji-Inseln, zusammen mit der Tiefe des Erdbebenherdes. Wir benutzen die geografische Länge und Breite als Variablen, auf die wir projizieren, und die Tiefe als Covariable, nach der wir Schnitte bilden.

Die Tiefen codieren wir um, damit bei grafischen Darstellungen große Tiefen nach unten zeigen.

```

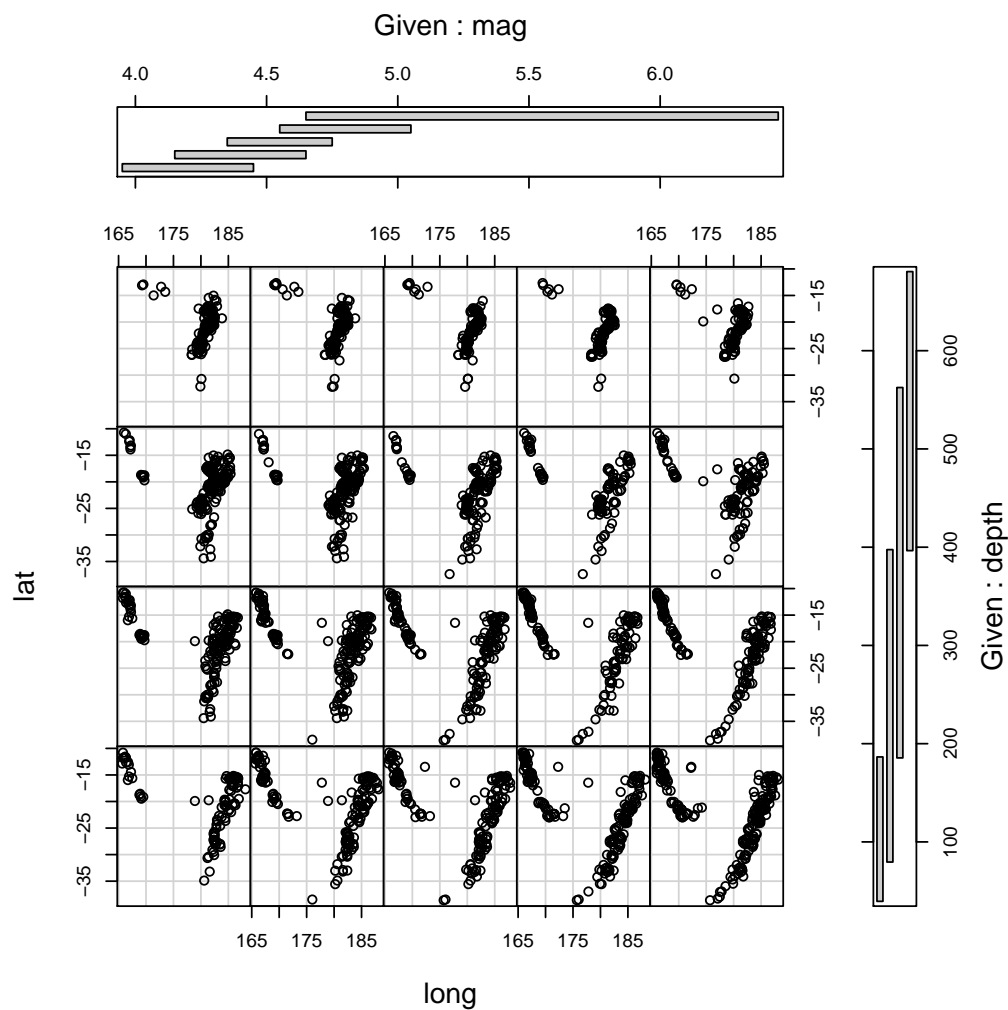
quakes$depth <- -quakes$depth Eingabe
given.depth <- co.intervals(quakes$depth, number = 4, overlap = .1)
coplot(lat ~ long | depth, data = quakes, given.values = given.depth, columns = 1)

```

Analog für zwei Covariable, die Tiefe und die Stärke des Erdbebens.

`coplot(lat ~ long | mag* depth , data = quakes, number = c(5, 4))` *Eingabe*

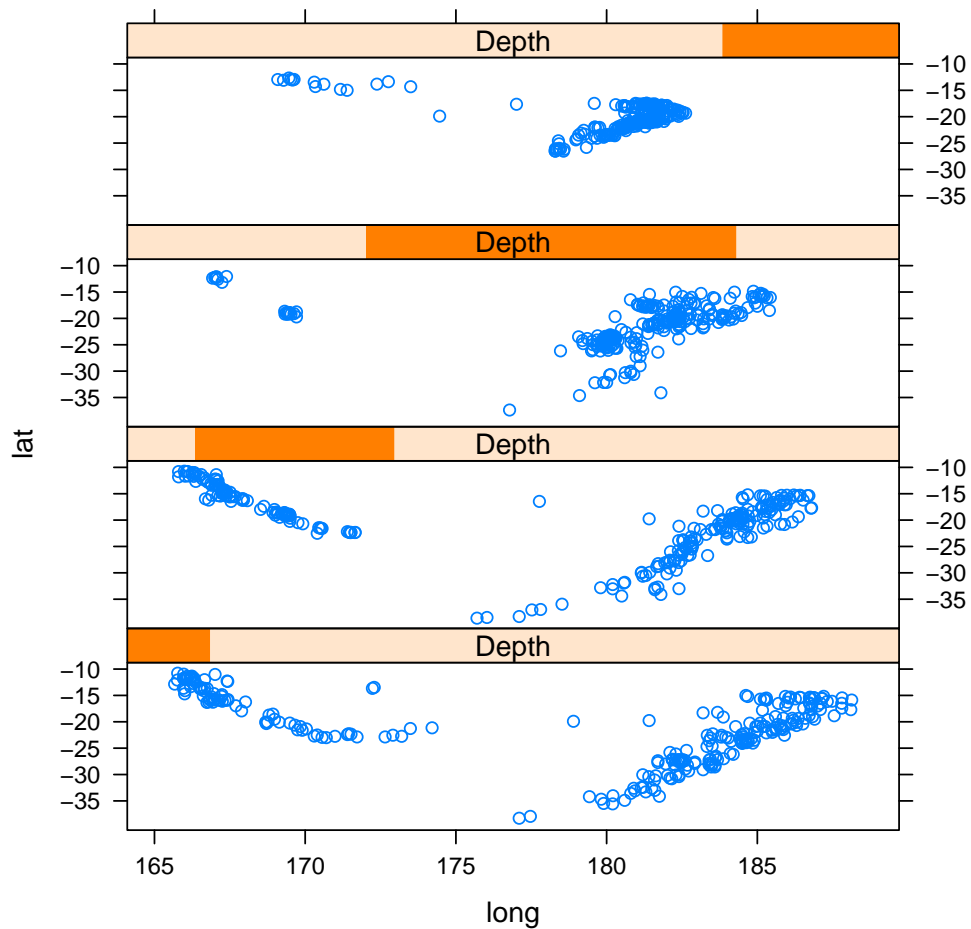


Aufgabe 4.4	
	Analysieren Sie den "quakes"-Datensatz. Fassen Sie Ihre Ergebnisse zusammen. Versuchen Sie, ein formales Modell zu formulieren.
	Wie hängt die geographische Position mit der Tiefe zusammen?
	Ist ein Zusammenhang von Tiefe und Stärke des Erdbebens erkennbar? (Evtl. müssen Sie bei <code>coplot()</code> eine andere Formel wählen.)

Die Idee der Coplots wird generalisiert in den Trellis-Displays (siehe [Cle93]). Trellis-Displays sind in R in `library(lattice)` implementiert.

Eingabe

```
library("lattice")
Depth <- equal.count(quakes$depth, number = 4, overlap = .1)
print(xyplot(lat ~ long | Depth, data = quakes, columns = 1, layout = c(1, 4)))
```



4.6. Transformationen und Dimensionsreduktion

Variable liegen oft in der Form vor, die von den Messprozessen oder fachlichen Konventionen vorgegeben sind. Sie entspricht nicht unbedingt der Form, die von der Sache her vorgegeben ist, oder die für die statistische Modellierung am besten geeignet ist. Diese Form enthält eine gewisse Beliebigkeit:

- Bei einer akustischen Reizbestimmung kann die Stärke der Reizes zum Beispiel durch die Energie beschrieben werden, oder durch den Schalldruck [Phon]. Von der einen zur anderen Skala führt die Logarithmus- Transformation. Das Weber-Fechnersche Gesetz der Psychologie sagt, dass für die menschliche Wahrnehmung die (logarithmische) Phon-Skala die richtige ist.
- Benzinverbrauch wird in den USA als Miles per Gallon angegeben, in Europa als Liter auf 100 km. Bis auf eine Umrechnungskonstante ist die eine Variable das inverse der anderen. Die Angabe in Liter auf 100 km scheint zu einfacheren statistischen Modellen zu führen; Analysen in Miles per Gallon können beliebig kompliziert sein.

Die Wahl der richtigen Variablen kann ein entscheidender Schritt in der Analyse sein. Dabei kann es hilfreich sein, zunächst Transformationen und zusätzliche konstruierte Variablen einzuführen, und dann in einem zweiten Schritt die Dimension wieder zu reduzieren und die effektiven Variablen zu bestimmen.

Koordinatensysteme sind nicht kanonisch vorgegeben. Dies trifft schon auf univariate Probleme zu. Bei univariaten Problemen können wir Koordinatensysteme noch relativ einfach transformieren. Die Modellierung der Fehlerverteilung einerseits und die Transformation der Daten auf eine Standard-Verteilung sind in gewisser Weise austauschbar. In mehrdimensionalen Situationen sind geeignete Transformationsfamilien bisweilen nicht verfügbar oder nicht zugänglich, und die Struktur des Problems kann kritisch von der Wahl geeigneter Koordinaten abhängig sein. Hier hat eine sachorientierte Wahl der Koordinatendarstellung oft den Vorzug vor automatischen Selektionen.

Dieses kann an Anderson's Iris-Datensatz illustriert werden. Der Datensatz hat fünf Dimensionen: vier quantitative Variable (Länge und Breite von Blütenblatt (engl. petal) und Kelchblatt (engl. sepal) von Iris-Blüten) und eine kategoriale Variable (die Spezies: *iris setosa canadensis*, *iris versicolor*, *iris virginica*)¹. Gesucht ist eine Klassifikation der Spezies anhand der vier quantitativen Variablen.

TABELLE 4.11. Iris Spezies.



Die Struktur ist ähnlich der des Diabetes-Datensatzes *chemdiab*. Die Klassifikation nach *iris\$Species* ist hier jedoch eine (extern) gegebene Klassifikation, im Gegensatz zur anhand der anderen Variablen definierten Klassifikation *chemdiab\$cc*. Gesucht ist hier nicht eine allgemeine Beschreibung wie bei *chemdiab*, sondern eine Klassifikationsregel, die *iris\$Species* aus den anderen Variablen ableitet.

Die Spezies definieren die Selektionen, die in diesem Beispiel von Interesse sind.

Um eine erste Übersicht zu bekommen ist es naheliegend, die vier Variablen getrennt nach Spezies zu betrachten. Die Standard-Konventionen von R machen dies umständlich. Die Spezies ist eine kategoriale Variable. Dies veranlasst R, bei der *plot()*-Funktion von einer Punkt-Darstellung zu Box&Whisker-Plots überzugehen.

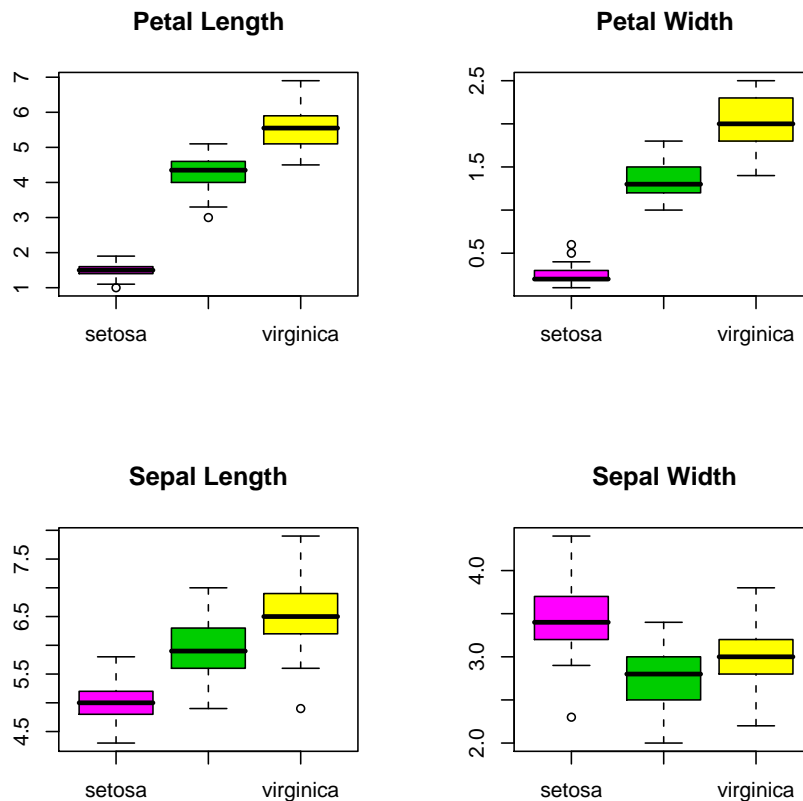
```
oldpar <- par(mfrow = c(2, 2))
plot(iris$Species, iris$Petal.Length,
      ylab = '', main = 'Petal Length', col = c("magenta", "green3", "yellow"))
plot(iris$Species, iris$Petal.Width,
      ylab = '', main = 'Petal Width', col = c("magenta", "green3", "yellow"))
plot(iris$Species, iris$Sepal.Length,
```

¹Photos: The Species Iris Group of North America. Mit freundlicher Genehmigung

```

ylab = '', main = 'Sepal Length', col = c("magenta", "green3", "yellow"))
plot(iris$Species, iris$Sepal.Width,
     ylab = '', main = 'Sepal Width', col = c("magenta", "green3", "yellow"))
par(oldpar)

```



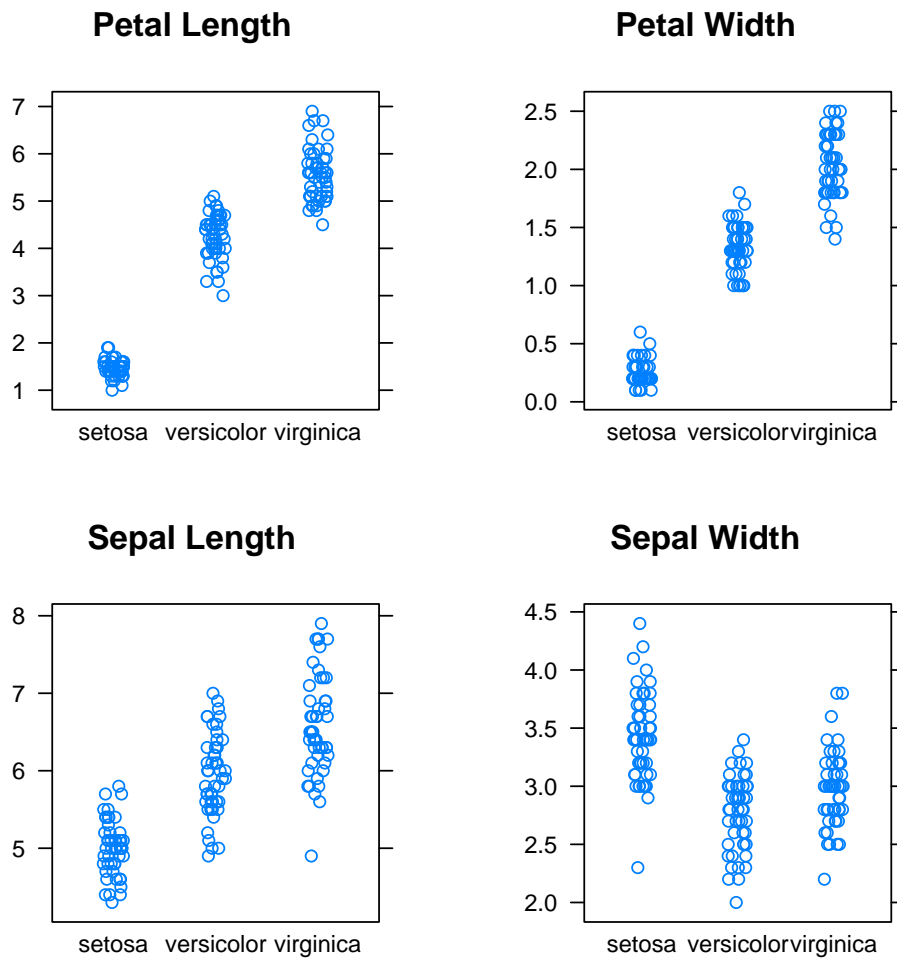
Wir könnten die R-Funktionen modifizieren, um einen Scatterplot der einzelnen Variablen nach Gruppen zu erhalten. Anstelle dessen greifen wir wieder auf `grid` und `lattice` zurück und benutzen die Funktion `stripplot()`. Weil bei der gegebenen Messgenauigkeit Werte vielfach auftreten, benutzen wir ein ‘jitter’: wir ‘verwackeln’ vielfache Werte, um sie getrennt darzustellen.

Eingabe

```

library("lattice")
print(stripplot(Petal.Length ~ Species, data = iris,
               jitter = TRUE, ylab = '', main = 'Petal Length'), split = c(1, 1, 2, 2), more = TRU
print(stripplot(Petal.Width ~ Species, data = iris,
               jitter = TRUE, ylab = '', main = 'Petal Width'), split = c(2, 1, 2, 2), more = TRU
print(stripplot(Sepal.Length ~ Species, data = iris,
               jitter = TRUE, ylab = '', main = 'Sepal Length'), split = c(1, 2, 2, 2), more = TRU
print(stripplot(Sepal.Width ~ Species, data = iris,
               jitter = TRUE, ylab = '', main = 'Sepal Width'), split = c(2, 2, 2, 2))

```



Die eindimensionalen Randverteilungen geben noch wenig Hinweis darauf, wie die drei Gruppen zu trennen sind. Auch die zweidimensionale Darstellung hilft wenig weiter.

Aufgabe 4.5	
	Benutzen Sie die Methoden aus Abschnitt 4.4 und 4.5, um den Datensatz zu untersuchen. Können Sie Klassifikationsregeln erkennen, die die drei Spezies weitgehend richtig klassifizieren?

Mit formalen Methoden wie der Diskriminanzanalyse (z. B. `lda()` in `library(MASS)`) kann die Klassifikation anhand der ursprünglichen Variablen gefunden werden. Die Trennung der Spezies ist nicht trivial.

Die ursprünglichen Variablen repräsentieren jedoch nur den Aspekt der Daten, der technisch am einfachsten erhebbar ist. Biologisch gesehen würde man jedoch anders parametrisieren: die Variablen spiegeln Größe und Form der Blätter wieder. Eine erste Approximation wäre

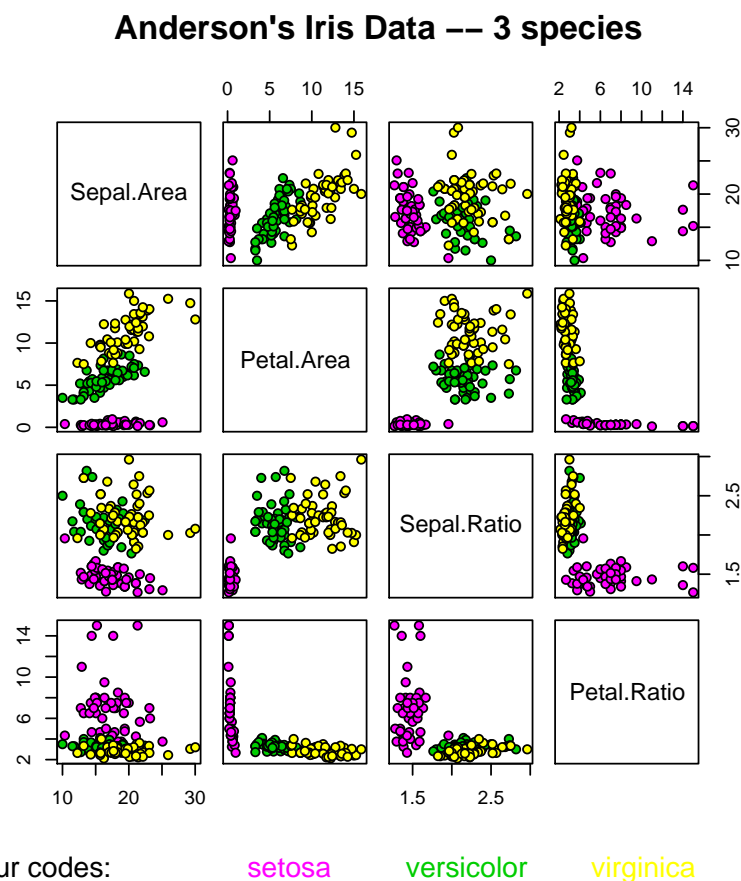
$$(4.1) \quad \text{area} = \text{length} \cdot \text{width}$$

$$(4.2) \quad \text{aspectratio} = \text{length}/\text{width}.$$

Damit erhält man die Darstellung

Eingabe

```
iris$Sepal.Area <- iris$Sepal.Length*iris$Sepal.Width
iris$Petal.Area <- iris$Petal.Length*iris$Petal.Width
iris$Sepal.Ratio <- iris$Sepal.Length/iris$Sepal.Width
iris$Petal.Ratio <- iris$Petal.Length/iris$Petal.Width
pairs(iris[6:9], main = "Anderson's Iris Data -- 3 species",
      pch = 21,
      bg = c("magenta", "green3", "yellow")[unclass(iris$Species)],
      oma = c(8, 8, 8, 8))
mtext(c("Colour codes:", levels(iris$Species)),
      col = c("black", "magenta", "green3", "yellow"),
      at = c(0.1, 0.4, 0.6, 0.8),
      side = 1, line = 2)
```



In der Marginalverteilung sind die Spezies fast vollständig getrennt - mit zwei Grenzfällen. In diesen mehr biologischen Koordinaten sieht man, dass zur Klassifikation Fläche und Längenverhältnis des Blütenblatts allein ausreichen. Jedes kompliziertere formale Verfahren muss sich mit dieser trivialen Klassifikationsregel erst einmal messen.

Selbst eine umfassende Suche, z.B. mit projection pursuit, erfasst nur die Projektionen, also nur spezielle Linearkombinationen der Variablen. Bei den Iris-Daten haben wir

zunächst neue Variablen, die Flächen und Seitenverhältnisse, eingeführt. Dies sind nicht-lineare Transformationen. Erst in einem weiteren Schritt sind dann die klassifizierenden Variablen identifiziert worden, und dabei ist die Dimension drastisch reduziert. Bei echten multivariaten Problemen ist es ganz typisch, dass zunächst eine Dimensionserweiterung notwendig ist, um das Problem zu lösen. Dimensionsreduktion ist erst dann sinnvoll, wenn die beschreibenden Variablen hinreichend komplex sind, um zu einer Lösung zu führen.

4.7. Höhere Dimensionen

4.7.1. Linearer Fall. Haben wir im wesentlichen lineare Strukturen, so können wir oft auch höher-dimensionale Strukturen mit Methoden analysieren, die für eindimensionale Modelle entwickelt sind. Wir müssen die Methoden evtl. modifizieren oder iteriert anwenden. Sie helfen uns jedoch, die wesentlichen Merkmale zu erkennen.

In Kapitel 2 haben wir lineare Modelle bereits allgemein für beliebige Dimension p der Regressoren eingeführt und damit bereits den mehrdimensionalen Fall eingeschlossen. Kapitel 2 setzt voraus, dass das Modell der statistischen Analyse vorab fest steht, d.h. das Information aus dem Datenmaterial die Wahl des Modells nicht beeinflusst, sondern nur die Entscheidung im Rahmen des Modells.

Insbesondere bei höherdimensionalen Problemen ist es jedoch so, dass das Modell erst zu bestimmen ist. Ein wichtiger Spezialfall ist die Auswahl von Regressoren: die Variablen sind Kandidaten, aus denen eine (möglichst kleine) Anzahl von Regressoren zu wählen ist.

Bringen kompliziertere Modelle eine wesentliche Verbesserung gegenüber dem einfachen Modell? Welche Parameter bzw. welche abgeleitete Variable sollten in das Modell einbezogen werden? Die Lehre aus den linearen Modellen ist, dass nicht der Wert des einzelnen Parameters den Beitrag im Modell bestimmt, sondern dass die durch die Parameter bestimmten Räume die wesentlichen Faktoren sind. An dieser Stelle sind angepasste Strategien gefragt. Wir können mit einfachen Modellen beginnen und fragen, ob zusätzliche Parameter einen weiteren Beitrag liefern. Dadurch erreichen wir einen besseren Fit, aber erhöhen die Varianz unserer Schätzungen. Oder wir können mit einem relativ komplexen Modell beginnen, und fragen, ob wir Parameter fortlassen können. Dadurch wird zwar der Restfehler erhöht, wir gewinnen aber an Verlässlichkeit der Schätzungen.

Beide Strategien führen im abstrakten linearen Regressionsmodell zu einem Vergleich von zwei Modellräumen $\mathcal{M}_{X'} \subset \mathcal{M}_X$. Die entsprechenden Schätzer sind $\pi_{\mathcal{M}_{X'}}(Y)$ und $\pi_{\mathcal{M}_X}(Y)$. Die Beziehung zwischen beiden wird klar, wenn wir die orthogonale Zerlegung $\mathcal{M}_X = \mathcal{M}_{X'} \oplus L_X := M_0, L_X := \mathcal{M}_X \ominus \mathcal{M}_{X'}$ von \mathcal{M}_X wählen. Dann ist $\pi_{\mathcal{M}_X}(Y) = \pi_{\mathcal{M}_{X'}}(Y) + \pi_{L_X}(Y)$.

4.7.1.1. *Partielle Residuen und Added-Variable-Plots.* In der Regression sind $\mathcal{M}_{X'}$ und \mathcal{M}_X Räume, die von den Regressor-Variablenvektoren aufgespannt werden. In unserer Situation interessiert uns der Spezialfall

$$X' = \text{span}(X_{1'}, \dots, X_{p'}); X = \text{span}(X_1, \dots, X_p)$$

mit $p > p'$. Dann wird aber L_X aufgespannt von den Vektoren

$$R_{p'+1} = X_{p'+1} - \pi_{\mathcal{M}_{X'}}(X_{p'+1}), \dots, R_p = X_p - \pi_{\mathcal{M}_{X'}}(X_p).$$

Wenn wir also (formal) eine lineare Regression der zusätzlichen Regressoren nach den bereits in X' enthaltenen durchführen, sind die dabei entstehenden Residuen ein Erzeugendensystem für L_X . Eine weitere Regression von Y nach diesen Residuen liefert uns den Term $\pi_{L_X}(Y)$, der den Unterschied zwischen den Modellen beschreibt. Nach Konstruktion wissen wir, dass $\pi_{\mathcal{M}_{X'}}(Y)$ orthogonal zu L_X ist. Bei dieser zweiten Regression wird deshalb

dieser Anteil auf null abgebildet. Wir können diesen Anteil gleich eliminieren und uns auf die Regression von $Y' = Y - \pi_{\mathcal{M}_{X'}}(Y)$ nach $R_{p'+1}, \dots, R_p$ beschränken.

Die Strategiewahl ist einfach: wir untersuchen, ob zusätzliche Parameter in das Modell aufgenommen werden sollten. Anstelle der Scatterplot-Matrix der ursprünglichen Daten betrachten wir die Scatterplots der (formalen) Residuen aus diesem einfachen Modell. Diese Scatterplots werden **Added-Variable-Plots** genannt.

Um den Unterschied zur Scatterplot-Matrix der Ausgangsdaten zu betonen: lineare Strukturen im Scatterplot der Ausgangsdaten sind ein klarer Hinweis auf lineare Abhängigkeiten. Nichtlineare Strukturen, wie z.B. die Dreiecksgestalt in einigen der Scatterplots können eine entsprechende Abhängigkeit widerspiegeln; sie können aber auch Artefakte sein, die als Folge der Verteilungs- und Korrelationsstruktur der Regressoren auftreten. Sie haben in der Regel keine einfache Deutung. Im Gegensatz dazu sind die Darstellungen in der Matrix der Added-Variable-Plots für lineare Effekte der vorausgehenden Variablen adjustiert. Dadurch hängen sie von der Wahl der Reihenfolge ab, in der Variable einbezogen werden. Sie korrigieren aber für lineare Effekte, die aus den Korrelationen zu vorausgehenden Variablen kommen. Dadurch wird eine ganze Reihe von Artefakten vermieden und sie können unter Berücksichtigung des Zusammenhangs unmittelbar interpretiert werden.

Aufgabe 4.6	
	<p>Modifizieren Sie die nachfolgende Prozedur <code>pairslm()</code> so, dass sie für alle Variablen in der ursprünglichen Matrix <code>x</code> die Residuen der Regression nach der neuen Variablen <code>x\$fit</code> berechnet und eine Scatterplot-Matrix dieser Residuen zeigt.</p> <pre data-bbox="531 1108 1268 1182">pairslm <- function(model, x, ...) { x\$fit <- lm(model, x)\$fit; pairs(x, ...)} </pre> <p>Fügen Sie auch Titel, Legenden etc. hinzu. Benutzen Sie den "trees"-Datensatz als Beispiel.</p>

Wir haben den Übergang von p' zu $p' + 1$ Variablen untersucht. Die Scatterplot-Matrix erlaubt uns einen schnellen Überblick über eine (nicht zu) große Zahl von Kandidaten (bei uns drei mögliche zusätzliche Regressoren). Der Übergang von p zu $p - 1$, zur Elimination einer Variablen, ist in gewisser Weise dual dazu. Dies entspricht der zweiten Strategie, der schrittweisen Elimination.

Statt eine einzelne Variable als Leitvariable auszuwählen ist es effizienter, Kombinationen von Variablen als synthetische Leitvariablen zu benutzen. Entsprechende Methoden werden in der Theorie als Hauptkomponentenanalyse behandelt und durch die Funktion `prcomp()` in der `library(mva)` bereitgestellt. Wir kommen daraus in einem späteren Beispiel (Seite 4-40) zurück.

Das Beispiel der linearen Modelle lehrt uns, dass die marginalen Beziehungen nur die halbe Wahrheit sind. Anstelle die einzelnen Regressoren zu betrachten, müssen wir im linearen Modell schrittweise orthogonalisieren. Komponentenweise Interpretationen sind damit fragwürdig - sie sind weitgehend von der Reihenfolge abhängig, in der Variablen einbezogen werden.

In komplexeren Situationen führen formale Methoden oft nur in die Irre. Handwerkliches Geschick ist hier notwendig. Leider sind die Kenntnisse darüber, wie handwerkliche Eingriffe die Gültigkeit formaler Methoden beeinflussen, noch sehr beschränkt. Deshalb ist es gerade hier wichtig, gewählte Strategien anhand von Simulationen kritisch zu beurteilen.

4.7.2. Nichtlinearer Fall. Nichtlineare Beziehungen in höheren Dimension stellen eine Herausforderung dar. Neben den Methoden brauchen wir auch ein Repertoire an Beispielen, die uns zeigen, welche Strukturen auftreten können und worauf wir achten müssen. Das folgende Beispiel, Cusp(Spitzen)-Singularität gehört dazu: es ist mit die einfachste Struktur, die in höheren Dimensionen auftreten kann. Die Basis ist hier ein zweidimensionale Struktur, eine Fläche, die nicht trivial in einem dreidimensionalen Raum eingebettet ist. Das interessante Merkmal ist hier die Aufspaltung von einer unimodalen in eine bimodale Situation.

4.7.2.1. *Beispiel: Spitzen-Nichtlinearität.* Das einfachste Beispiel kann im Hinblick auf physikalische Anwendungen illustriert werden. In physikalischen Systemen hängen Wahrscheinlichkeitsverteilungen oft mit Energiezuständen zusammen; (lokale) Minima der Energie entsprechen dabei den Moden der Verteilung. Ein typischer Zusammenhang ist: verhält sich die Energie wie $\varphi(y)$, so verhält sich die Verteilung nach Standardisierung wie $e^{-\varphi(y)}$. Ist $\varphi(y)$ in der Nähe des Minimums quadratisch, so erhalten wir (bis auf Skalentransformation) Verteilungen aus der Familie der Normalverteilungen.

Die Differentialtopologie lehrt uns, dass auch bei kleinen Störungen oder Variationen dieses qualitative Bild erhalten bleibt. Die Energie bleibt zumindest lokal approximativ quadratisch, und die Normalverteilungen bleiben zumindest approximativ eine geeignete Verteilungsfamilie.

Das Verhalten ändert sich drastisch, wenn das Potential sich lokal wie y^4 verhält. Schon geringe Variationen können dazu führen, dass das Potential lokal quadratisch ist. Aber sie können auch dazu führen, dass das lokale Minimum aufbricht und zu zwei Minima führt. Das typische Bild ist von der Gestalt

$$(4.3) \quad \varphi(y; u, v) = y^4 + u \cdot y^2 + v \cdot y.$$

Dabei sind die Variationen durch die Parameter u, v repräsentiert. Am einfachsten lässt sich die Situation dynamisch interpretieren: wir stellen uns vor, dass u, v äußere Parameter sind, die sich verändern können. Dieses Bild kennen wir von der magnetischen Hysterese: y gibt die Magnetisierung in einer Richtung an, u spielt die Rolle der Temperatur; v die eines äußeren Magnetfelds. Bei hoher Temperatur folgt die Magnetisierung direkt dem äußeren Magnetfeld. Sinkt die Temperatur, so zeigt das Material Gedächtnis: die Magnetisierung hängt nicht nur vom äußeren Magnetfeld ab, sondern auch von der vorhergehenden Magnetisierung.

Ähnliche “Gedächtniseffekte” kennen wir auch in anderen Bereichen. Man stelle sich einen Markt vor mit Preisen y , Kosten v und einem “Konkurrenzdruck” u . Bei ausreichender Konkurrenz folgen die Preise (mehr oder weniger) den Kosten bei sonst gleichen Bedingungen. Bei Monopol-Situationen scheinen die Preise ein Gedächtnis zu haben: sind sie einmal gestiegen, so sinken sie erst, wenn die Kosten drastisch reduziert sind.

Die in Formel 4.3 angegebenen “Entfaltung” des Potentials y^4 hat eine typische Form. Aus

$$(4.4) \quad \varphi'(y; u, v) = 4y^3 + 2u \cdot y + v = 0$$

erhält man die kritischen Punkte (siehe Abb. 4.2). Siehe Abbildung 4.2 auf Seite 4-30.

Projiziert auf die u, v -Ebene gibt dies eine Spitze (engl.: “cusp”, Abb. 4.3). Bei Parametern im inneren dieser Spitze gibt es zwei lokale Minima; außerhalb der Spitze gibt es nur einen Extremalwert.

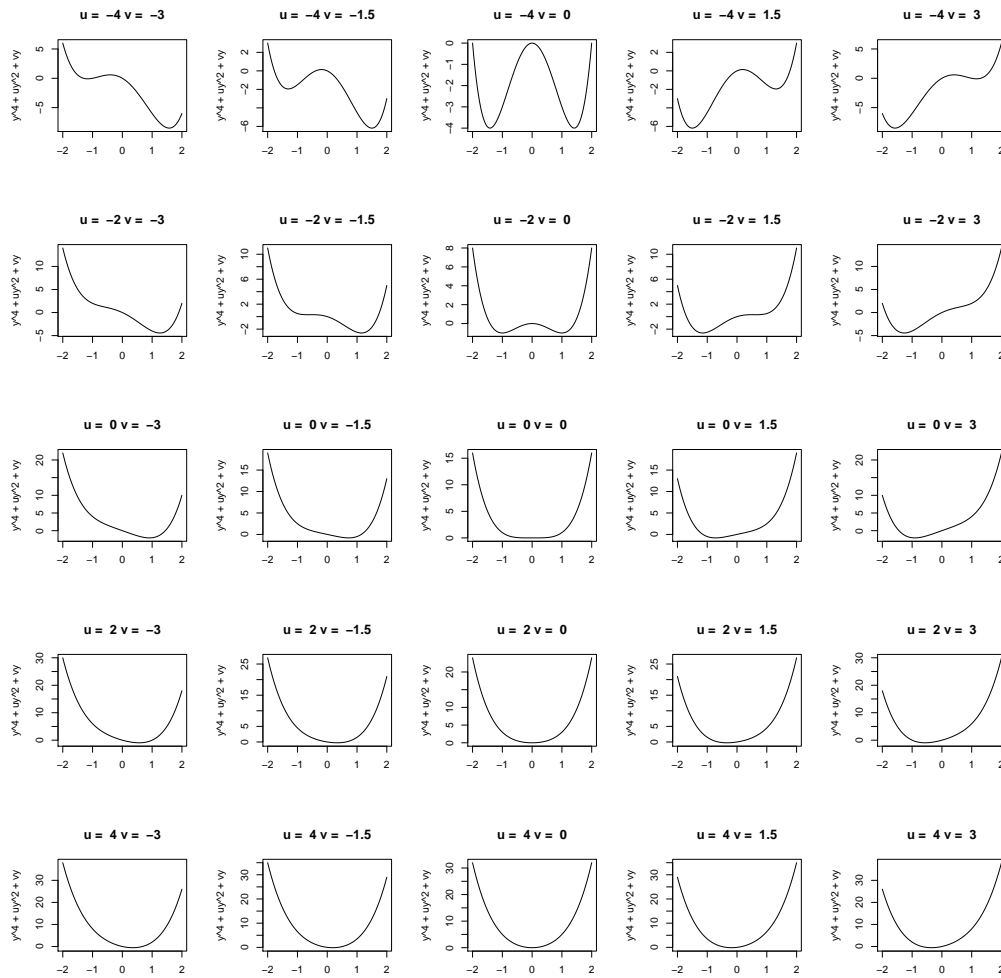


ABBILDUNG 4.1. Entfaltung von y^4 : $\varphi(y; u, v) = y^4 + u \cdot y^2 + v \cdot y$

Die diesen Potentialen entsprechenden Verteilungen sind – bis auf Skalentransformation zur Normalisierung –

$$(4.5) \quad p(y; u, v) \propto e^{-(y^4 + u \cdot y^2 + v \cdot y)}.$$

Die Struktur der Potentiale spiegelt sich auch in den entsprechenden Verteilungen wieder; der exponentielle Abfall macht allerdings die kritische Grenze etwas komplizierter.

Die Situation erscheint hier noch harmlos: der Parameterraum (der Raum der Regressoren) $x = (u, v)$ hat nur zwei Dimensionen. Die Verteilung ist eindimensional mit einer glatten Dichte. Aber die Situation kann mit linearen Methoden nur unzureichend erfasst werden. Der typische nichtlineare Effekt wird nicht erkannt, wenn man darauf nicht vorbereitet ist. Erst das Gesamtbild im Dreidimensionalen vermittelt die eigentliche Struktur.

Dieses einfache Beispiel ist eine Herausforderung. Wie kann eine derartige Struktur diagnostiziert werden?

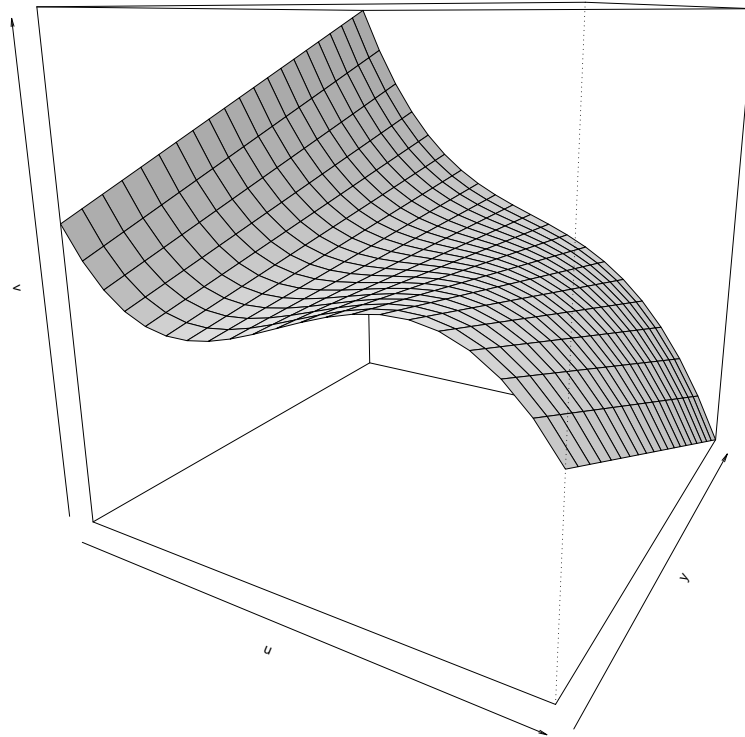


ABBILDUNG 4.2. Kritische Punkte $\varphi'(y; u, v) = 4y^3 + 2u \cdot y + v = 0$

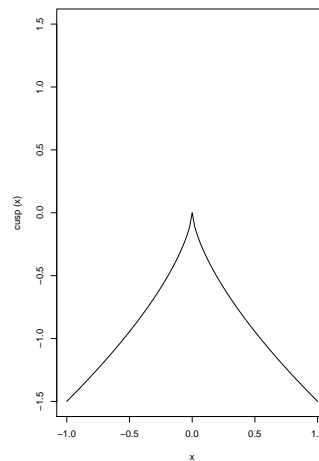


ABBILDUNG 4.3. Grenze zwischen Uni- und Bimodalität im (u, v) -Raum

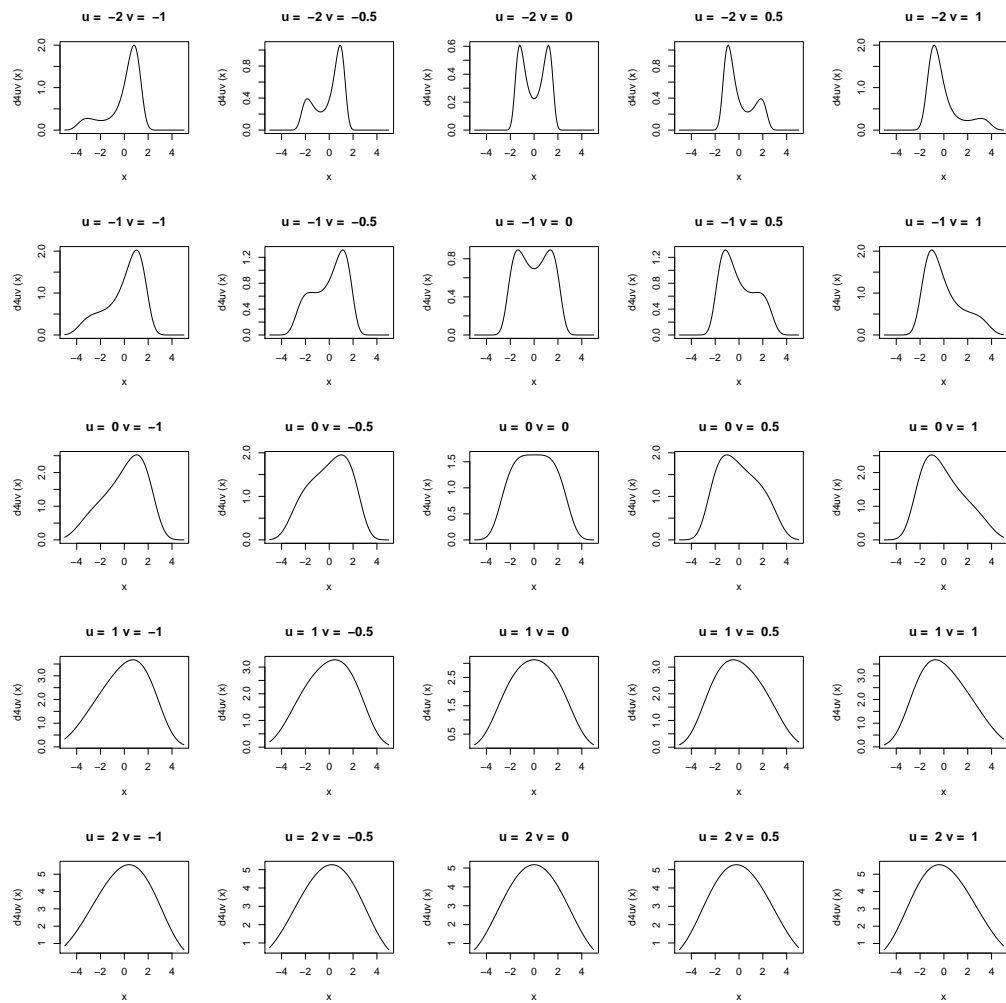


ABBILDUNG 4.4. $p(y; u, v) \propto e^{-(y^4+u \cdot y^2+v \cdot y)}$

Aufgabe 4.7	
	Schreiben Sie eine Funktion <code>dx4exp(x, u, v)</code> , die die zentrierte Wahrscheinlichkeitsdichte zu (4.5) berechnet. Dazu müssen Sie die Dichte aus (4.5) integrieren, um die Normierungskonstante zu bestimmen, und den Erwartungswert berechnen, um die Dichte zu zentrieren. Benutzen Sie für beides eine numerische Integration mit <code>integrate()</code> .
***	<p>Simulieren Sie zu Werten u, v auf einem Gitter in $u = -2 \dots 2$ und $v = -1 \dots 1$ je 100 Zufallszahlen aus <code>dx4exp(x, u, v)</code>. Untersuchen Sie diese mit den Methoden aus Kapitel 2.</p> <p>Können Sie Hinweise auf nicht-lineare Abhängigkeit erkennen? Ist die Bimodalität erkennbar? Wie weit können Sie die Struktur identifizieren?</p>

Bei nichtlinearen Beziehungen können gemeinsame Abhängigkeiten eine große Bedeutung haben. Im allgemeinen erfordert dies Umsicht bei der Modellbildung. Nichtlineare Beziehungen können in Projektionen versteckt sein. Artefakte der (linearen) Projektion können ein Bild vermitteln, das nicht den ursprünglichen Beziehungen entspricht.

4.7.3. “Curse of Dimension”. Ohne angepasste Koordinatensysteme ist eine umfassende Suche nach interessanten Projektionen und Schnitten nötig. Die Anzahl der Möglichkeiten steigt rasch mit der Dimension. Zur Illustration: Um einen Kubus zu identifizieren müssen zumindest die Eckpunkte erkannt werden. In d Dimensionen sind dies 2^d Eckpunkte. Die Anzahl steigt exponentiell mit der Dimension. Dies ist ein Aspekt des als *curse of dimension* bekannten Problems.

Anders betrachtet: bezeichnen wir die Datenpunkte, die in mindestens einer Variablen-dimension extrem sind, so sind dies im eindimensionalen Fall zwei Punkte. in d Dimensionen sind dies typischerweise 2^d Punkte. Betrachten wir nicht nur Koordinatenrichtungen, sondern beliebige Richtungen, so ist typisch jeder Punkt extremal, wenn d sehr groß wird.

Ein dritter Aspekt: im d -dimensionalen Raum ist fast jeder Punkt isoliert. Lokalisierungen, wie wir sie in Abschnitt 2.5 kennengelernt haben, brechen zusammen. Wählen wir um einen Punkt eine Umgebung, die einen Anteil p , z.B. $p = 10\%$ der Variablenspannweite umfasst, so haben wir in einer Dimension typischerweise der Größenordnung nach auch einen Anteil p der Datenpunkte erfasst. In d Dimensionen ist dies nur noch ein Anteil der Größenordnung p^d . Bei zum Beispiel 6 Dimensionen brauchen wir also mehrere Millionen Datenpunkte, damit wir nicht mit leeren Umgebungen arbeiten.

4.7.4. Fallstudie. Als fortlaufendes Beispiel benutzen wir nun den Fat-Datensatz. Dieser Datensatz ist in der Literatur wiederholt veröffentlicht und in R unter anderem im Paket *UsingR* zugänglich.

Ziel der Untersuchung hinter diesem Datensatz ist die Bestimmung des Körperfettanteils. Die verlässlichste Methode ist es, in einem Wasserbad die mittlere Dichte des Gewebes zu bestimmen und daraus auf den Körperfettanteil zurück zu schliessen. Diese Bestimmung ist sehr aufwendig. Kann sich durch einfacher zu messende Körperparameter ersetzt werden? Die zur Verfügung stehenden Parameter sind in Tabelle 4.15 zusammengefasst.

Anhand der Übersicht in Tabelle 4.15 sehen wir gleich, dass metrische Angaben und US-Maße gemischt sind. Damit für uns die Interpretation einfacher ist, stellen wir alle Angaben auf metrische Werte um.

Eingabe

```
library("UsingR")
data(fat)
fat$weightkg <- fat$weight*0.453
fat$heightcm <- fat$height * 2.54
fat$ffweightkg <- fat$ffweight*0.453
```

Die Variablen *body.fat* und *body.fat.siri* sind aus dem gemessenen Wert *density* abgeleitet. Hinter den Formeln stecken Annahmen über die mittlere Dichte von Fett und von fettfreiem Gewebe. Mit diesen Annahmen kann aus *density* der Fettanteil errechnet (oder besser: geschätzt) werden. In beiden Formeln ist der dichtabhängige Faktor $1/\text{density}$. Bis auf (gegebene oder angenommene) Konstanten ist dies also der für uns relevante Term (und nicht *density*).

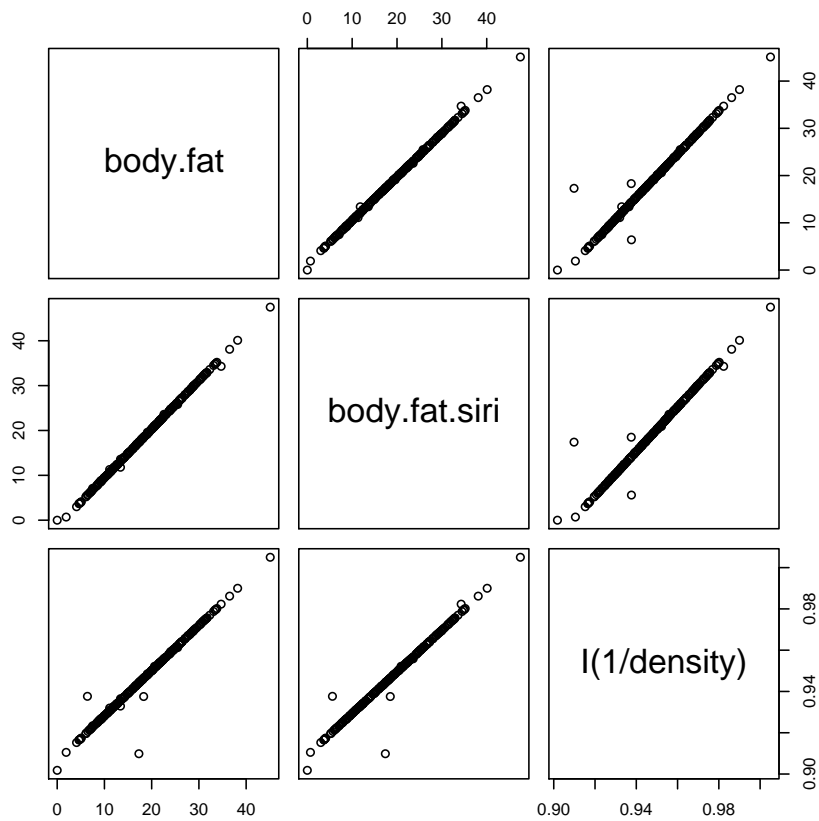
Der erste Schritt ist eine kritische Inspektion und Bereinigung des Datensatzes. Dies ist fast immer nötig, nicht nur bei höherdimensionalen Datensätzen. Bei höherdimensionalen Datensätzen haben wir allerdings oft Redundanzen, die Konsistenzprüfungen und evtl. Korrekturen ermöglichen. In unserem Fall sind *body.fat*, *body.fat.siri*, *ffweight* und *BMI* abgeleitete Größen, die zu anderen Variablen in deterministischer Beziehung stehen.

Name	Variable	Einheit, Bem.
case	Case Number	
body.fat	Percent body fat using Brozek's equation, $457/Density - 414.2$	
body.fat.siri	Percent body fat using Siri's equation, $495/Density - 450$	
density	Density	[g/cm ²]
age	Age	[yrs]
weight	Weight	[lbs]
height	Height	[inches]
BMI	Adiposity index = $Weight/Height^2$	[kg/m ²]
ffweight	Fat Free Weight = $(1 - fractionofbodyfat) * Weight$, using Brozek's formula	[lbs]
neck	Neck circumference	[cm]
chest	Chest circumference	[cm]
abdomen	Abdomen circumference "at the umbilicus and level with the iliac crest"	[cm]
hip	Hip circumference	[cm]
thigh	Thigh circumference	[cm]
knee	Knee circumference	[cm]
ankle	Ankle circumference	[cm]
bicep	Extended biceps circumference	[cm]
forearm	Forearm circumference	[cm]
wrist	Wrist circumference "distal to the styloid pro- cesses"	[cm]

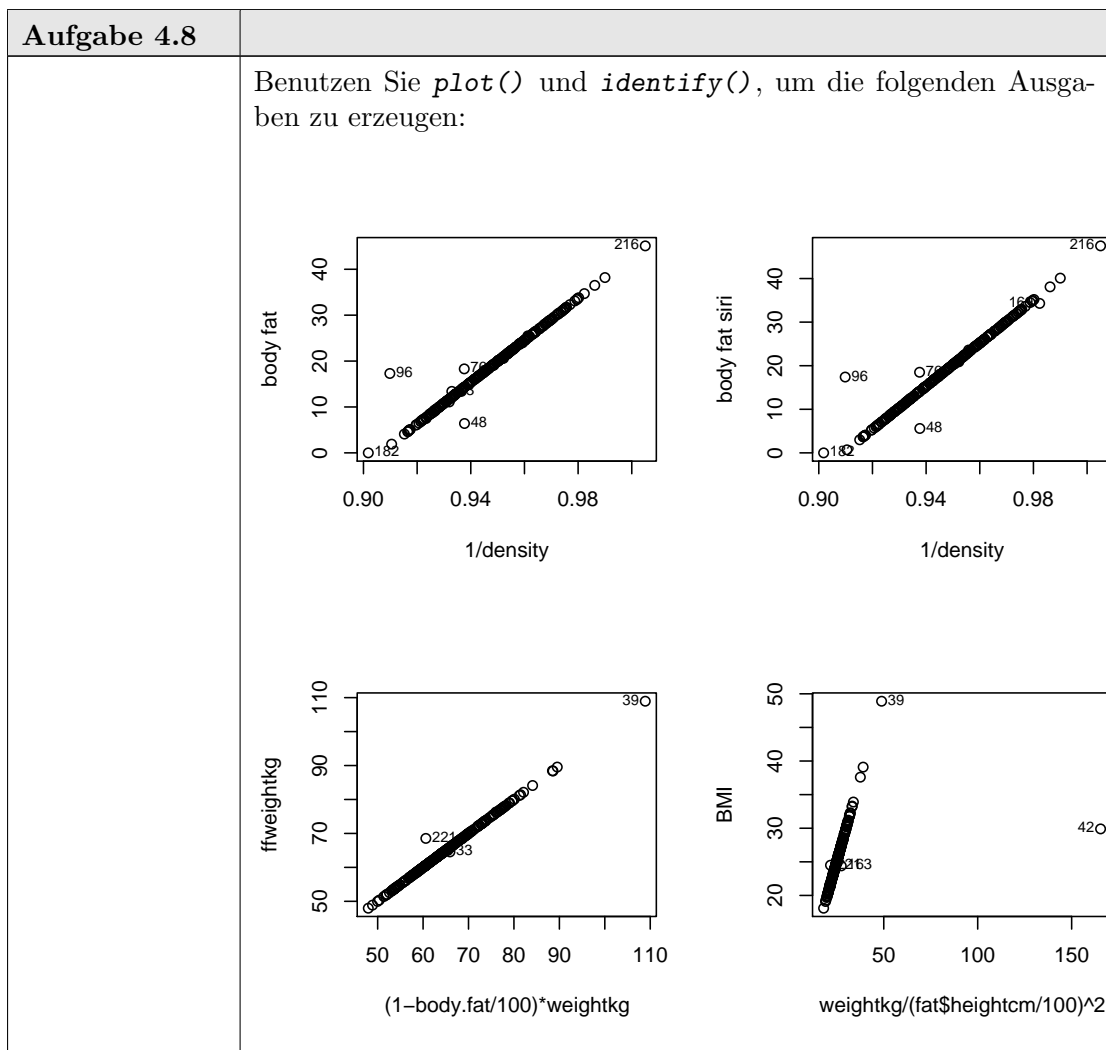
TABELLE 4.15. Fat data set: variables

Wir betrachten zunächst die Gruppe *body.fat*, *body.fat.siri*, $1/density$. Die paarweisen Scatterplots sollten Geraden zeigen. `pairs()` leistet gute Dienste. Wir benutzen es hier in der Formel-Variante. Um zu signalisieren, dass $1/density$ berechnet werden soll, und die Division nicht als Formel-Operator zu verstehen ist, markieren wir den Term entsprechend.

`pairs(~body.fat + body.fat.siri + I(1/density), data = fat)` *Eingabe*



Die inkonsistenten Werte und Ausreißer sind deutlich. Leider ist es in R nicht einfach möglich, Werte in der Scatterplot-Matrix zu markieren.



Wenn eine klare Korrektur vorgenommen werden kann, so sollte es hier getan und im Auswertungsbericht notiert werden. Fall 42 ist einfach: eine Größe von 0.73m bei einem Gewicht von 63.5kg ist unplausibel und inkonsistent zu BMI 29.9. Aus dem BMI lässt sich die Größe rückrechnen. Der eingetragene Wert von 29.5 Zoll sollte wohl 69.5 Zoll sein.

Eingabe

```
fat$height [42] <- 69.5
fat$heightcm[42] <- fat$height[42] * 2.54
```

Fall 216 ist eine Ermessenssache. Die Dichte ist extrem niedrig, der BMI extrem hoch. Andererseits passen die Körpermaße zu diesen Extremen. Diese Fall kann ein Ausreisser sein, der die Auswertungen verzerren kann. Es kann aber auch eine besonders informative Beobachtung sein. Wir notieren ihn als Besonderheit.

Nach dieser Voruntersuchung bereinigen wir den Datensatz. Die Variablen, die keine Information mehr enthalten oder die wir ersetzt haben, löschen wir. Als Zielvariable benutzen wir `body.fat`. Wir behalten jedoch noch die Variable `density` für spätere Zwecke.

Eingabe

```
fat$weight <- NULL
fat$height <- NULL
```

```
fat$ffweight <- NULL
fat$ffweightkg <- NULL
fat$body.fat.siri <- NULL
```

Es gibt eine Reihe von gängigen Indizes (siehe Abb. 4.5). Früher war die Faustformel ‘Idealgewicht = Körpergröße -100’ gängig. Heute ist der “body mass index” BMI = Gewicht/ Körpergröße² gängig. (Handelsübliche Körperfettwaagen bestimmen die elektrische Impedanz. Diese Variable ist im Fat-Datensatz nicht enthalten.)

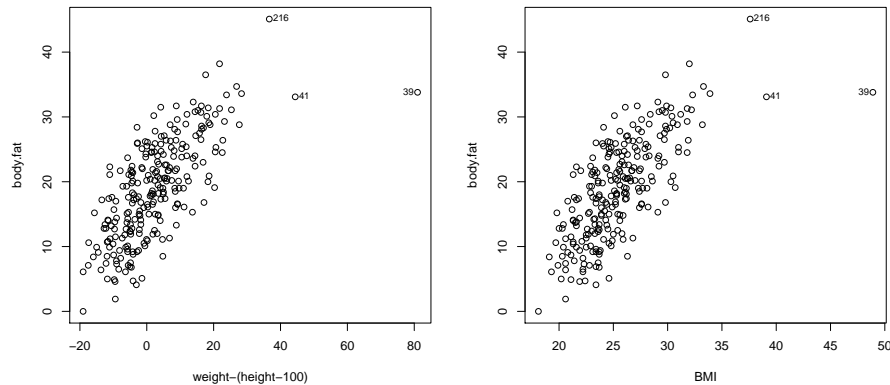


ABBILDUNG 4.5. Fettanteil gegen konventionelle Indizes

Wir können mit den konventionelle Indizes im linearen Modell schätzen. Dabei lassen wir die offensichtlichen Ausreisser und möglichen Hebelpunkte unberücksichtigt. Dazu benutzen wir den `subset`-Parameter der Funktion `lm()`.

Für die Faustformel ‘Idealgewicht = Körpergröße -100’ erhalten wir:

```
lm.height <- lm(body.fat ~ I(weightkg - (heightcm - 100)),
  data = fat,
  subset = -c(39, 41, 216))
summary(lm.height)
```

```
Call:
lm(formula = body.fat ~ I(weightkg - (heightcm - 100)), data = fat,
    subset = -c(39, 41, 216))
```

Residuals:

Min	1Q	Median	3Q	Max
-11.90734	-3.68697	-0.05303	3.65458	12.28000

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	17.70722	0.33296	53.18	<2e-16
I(weightkg - (heightcm - 100))	0.54557	0.03283	16.62	<2e-16

(Intercept) ***

I(weightkg - (heightcm - 100)) ***

```

---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.166 on 247 degrees of freedom
Multiple R-squared:  0.5279,    Adjusted R-squared:  0.526
F-statistic: 276.2 on 1 and 247 DF,  p-value: < 2.2e-16

```

Die Regression von *body.fat* nach *BMI* ergibt:

```

lm.BMI <- lm(body.fat~BMI,      Eingabe
             data = fat,
             subset = -c(39, 41, 216))
summary(lm.BMI)

```

```

Call:
lm(formula = body.fat ~ BMI, data = fat, subset = -c(39, 41,
216))

```

```

Residuals:
    Min       1Q   Median       3Q      Max
-12.49460  -3.53561  -0.05228   3.69129  11.72720

```

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -25.6130     2.6212  -9.772  <2e-16 ***
BMI           1.7564     0.1031  17.042  <2e-16 ***

```

```

---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.097 on 247 degrees of freedom
Multiple R-squared:  0.5404,    Adjusted R-squared:  0.5385
F-statistic: 290.4 on 1 and 247 DF,  p-value: < 2.2e-16

```

Der Fit ist jedoch mit $R^2 = 0.53$ bzw. $R^2 = 0.54$ in beiden Fällen nur mäßig.

Selbst mit allen Datenpunkten und allen Regressoren wird maximal $R^2 = 0.75$ erreicht:

```

lm.fullres <- lm(body.fat ~ age + BMI + neck + chest +
                 abdomen + hip + thigh + knee + ankle +
                 bicep + forearm + wrist + weightkg + heightcm,
                 data = fat)
summary(lm.fullres)

```

```

Call:
lm(formula = body.fat ~ age + BMI + neck + chest + abdomen +
    hip + thigh + knee + ankle + bicep + forearm + wrist + weightkg +
    heightcm, data = fat)

```

```

Residuals:
    Min       1Q   Median       3Q      Max
-10.0761  -2.6118  -0.1055   2.8993   9.2691

```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-50.804727	36.489198	-1.392	0.16513	
age	0.061005	0.029862	2.043	0.04217	*
BMI	0.782993	0.733562	1.067	0.28688	
neck	-0.439082	0.218157	-2.013	0.04528	*
chest	-0.040915	0.098266	-0.416	0.67751	
abdomen	0.866361	0.085550	10.127	< 2e-16	***
hip	-0.206231	0.136298	-1.513	0.13159	
thigh	0.246127	0.135373	1.818	0.07031	.
knee	-0.005706	0.229564	-0.025	0.98019	
ankle	0.135779	0.208314	0.652	0.51516	
bicep	0.149100	0.159807	0.933	0.35177	
forearm	0.409032	0.186022	2.199	0.02886	*
wrist	-1.514111	0.493759	-3.066	0.00242	**
weightkg	-0.389753	0.221592	-1.759	0.07989	.
heightcm	0.187196	0.199854	0.937	0.34989	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.991 on 237 degrees of freedom

Multiple R-squared: 0.7497, Adjusted R-squared: 0.7349

F-statistic: 50.7 on 14 and 237 DF, p-value: < 2.2e-16

Dies ist ein Modell mit 15 Koeffizienten. Das Modell ist so komplex, das es kaum zu interpretieren ist, und man wird versuchen, das Modell zu reduzieren. Anstelle "von Hand" nach einfacheren Modellen zu suchen, kann dieser Prozess automatisiert werden. Dazu dient die Funktion `regsubsets()` in `library(leaps)`. Der quadratische Fehler (bzw. das Bestimmtheitsmaß R^2) muss dabei modifiziert werden: der quadratische Fehler wird minimiert, wenn wir alle Regressoren ins Modell aufnehmen, also immer im vollen Modell. Zur Modellwahl benutzt man Varianten des quadratischen Fehlers (bzw. des Bestimmtheitsmaßes R^2), die für die Anzahl der Parameter adjustiert sind.

Aufgabe 4.9	
*	Benutzen Sie <code>library(leaps)</code> <code>lm.reg <- regsubsets(body.fat ~ age + BMI + neck + chest + abdomen + hip + thigh + knee + ankle + bicep + forearm + wrist + weightkg + heightcm, data = fat)</code> und inspizieren Sie das Resultat mit <code>summary(lm.reg)</code> <code>plot(lm.reg, scale = r2)</code> <code>plot(lm.reg, scale = "bic")</code> <code>plot(lm.reg, scale = Cp)</code> Hinweis: siehe <code>help(plot.regsubsets)</code>
*	Benutzen Sie die Funktion <code>leaps()</code> zur Modellselektion.

Allerdings sind nun die Werkzeuge, die wir in Kapitel 2 kennengelernt haben, unbrauchbar geworden. Die statistischen Aussagen in der Zusammenfassung sind nur gültig,

wenn Modell bzw. Hypothesen unabhängig vom Datenmaterial festgelegt sind. Wenn anhand des Datenmaterials das Modell erst bestimmt wird, ist unklar, wie die geschätzten Koeffizienten verteilt, d.h. wie Konfidenzintervalle zu bestimmen sind bzw. wie zu testen ist. Die Software hat keine Information darüber, dass wir uns in einem Modellwahlprozess befinden und gibt die Wahrscheinlichkeiten aus, die bei festem Modell unter Normalverteilungsannahme gelten würden.

Auch die Diagnostik wird unbrauchbar: der zentrale Grenzwertsatz sorgt dafür, dass unter schwachen Unabhängigkeitsannahmen die Residuen bei der großen Anzahl von Termen approximativ normalverteilt sind, selbst wenn dies für die Fehler nicht zutrifft.

Wir sind in einer Sackgasse.

Wir illustrieren nun einen anderen Zugang, der etwas weiter führt. Dazu versetzen wir uns an den Anfang der Analyse, nach der ersten Inspektion und Datenkorrektur. Damit wir nicht in das oben gesehen Problem laufen, dass die statistischen Verteilungen durch vorhergehende Modellwahl schritte beeinflusst werden, teilen wir den Datensatz auf. Einen Teil benutzen wir als Trainingsteil, an dem wir das Modell wählen und verschiedenen Alternativmöglichkeiten durchspielen können. Der Rest wird als Auswertungsteil reserviert. Dessen Information wird erst nach Modellwahl für die statistische Analyse benutzt.

Bei genauerer Überlegung zeigt sich, dass der Modellwahl schritt nur für die Abschätzung der Fehler kritisch ist, nicht für die Parameter-Schätzung. Wird der Fehler anhand der Daten geschätzt, die zur Modellwahl benutzt sind, so unterschätzen wir tendenziell die Fehler. Der Auswertungsteil dient der verlässlichen Fehlerabschätzung und Residuen-diagnostik. Dies ist eine eingeschränkte Aufgabe. Deshalb reservieren wir dafür nur einen kleineren Teil.

Eingabe

```
sel <- runif(dim(fat)[1])
fat$train <- sel < 2/3
rm(sel)
```

Die Ausreisser eliminieren wir aus dem Trainingsteil

Eingabe

```
fat$train[c(39, 41, 216)] <- FALSE
summary(fat$train)
```

Ausgabe

Mode	FALSE	TRUE
logical	93	159

Unsere Zielvariable ist `body.fat`, oder, proportional dazu, `1/density`.

Wir versuchen zunächst, die Variablen inhaltlich zu sortieren. Für die Dichte habe wir eine physikalische Definition

$$Dichte = \frac{Gewicht}{Volumen}.$$

Unter den Variablen, die als Regressoren in Betracht kommen, haben wir eine Variable, die direkt das Gewicht angibt (`weight` bzw. `weightkg`), eine ganze Reihe von Variablen, die Körpermaße widerspiegeln, sowie das Alter `age`.

Aus der gemessenen Dichte und dem gemessenen Gewicht lässt sich das Volumen er rechnen. Wir erweitern dadurch die Variablen. Da wir hier nur eine Gewichtsmessung pro Person haben, bleibt kein Platz für personenbezogene Statistik.

Eingabe

```
fat$vol <- fat$weightkg/fat$density
```

Wir versuchen nun, das Volumen `fat$vol` zu schätzen. Die Körpermaße sind lineare Werte. In einer groben Approximation können wir daraus Volumen-Werte ableiten. Die einzige Längeninformation, die wir haben, steckt in `height`. Mangels besserer Information nehmen wir an, dass alle Körperteile eine Länge haben, die proportional zur Körpergröße ist. Wir zielen auf ein lineares Modell. Deshalb können wir lineare Faktoren vernachlässigen. Approximieren wir die Körperteile durch Zylinder, so erhalten wir, bis auf lineare Faktoren

Eingabe

```
fat$neckvol <- fat$neck^2 * fat$heightcm
fat$chestvol <- fat$chest^2 * fat$heightcm
fat$abdomenvol <- fat$abdomen^2 * fat$heightcm
fat$hipvol <- fat$hip^2 * fat$heightcm
fat$thighvol <- fat$thigh^2 * fat$heightcm
fat$kneevol <- fat$knee^2 * fat$heightcm
fat$anklevol <- fat$ankle^2 * fat$heightcm
fat$bicepvol <- fat$bicep^2 * fat$heightcm
fat$forearmvol <- fat$forearm^2 * fat$heightcm
fat$wristvol <- fat$wrist^2 * fat$heightcm
```

Als nächstes untersuchen wir die interne Struktur der Regressor-Kandidaten im Trainings-teil. Wir tun dies getrennt für die linearen Variablen und für die Volumen- Variablen. Dazu benutzen wir die Funktion `prcomp()`, die zu gegebenen Variablen schrittweise beste lineare Prediktoren liefert.

Für die approximativen Körperteil-Volumen sind die Hauptkomponenten:

Eingabe

```
pcfatvol <- prcomp(fat[, 20:29], subset = fat$train)
round(pcfatvol$rotation, 3)
```

	Ausgabe							
	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8
neckvol	0.054	-0.017	0.063	-0.070	0.570	0.045	0.787	-0.169
chestvol	0.548	0.457	0.691	0.012	-0.113	-0.004	-0.003	-0.012
abdomenvol	0.650	0.290	-0.697	-0.059	0.049	-0.011	-0.017	0.016
hipvol	0.483	-0.748	0.099	0.437	-0.031	0.069	0.003	0.018
thighvol	0.185	-0.370	0.074	-0.870	-0.244	-0.017	0.079	0.012
kneevol	0.056	-0.082	0.055	-0.061	0.336	-0.822	-0.278	-0.159
anklevol	0.016	-0.032	0.034	-0.005	0.140	-0.236	-0.047	-0.120
bicepvol	0.051	-0.047	0.079	-0.180	0.552	0.508	-0.531	-0.333
forearmvol	0.024	-0.019	0.074	-0.086	0.388	0.019	-0.104	0.906
wristvol	0.009	-0.005	0.017	0.002	0.110	-0.053	0.044	0.003
	PC9	PC10						
neckvol	0.067	0.093						
chestvol	0.005	0.002						
abdomenvol	-0.015	-0.002						
hipvol	0.011	0.004						
thighvol	-0.015	-0.019						
kneevol	0.299	0.056						
anklevol	-0.949	0.073						
bicepvol	0.030	0.010						

```
forearmvol -0.049 0.044
wristvol -0.048 -0.990
```

Das Muster der Vorzeichen bei den Ladungen gibt Hinweise auf die interne Struktur. Die erste Hauptkomponente $PC1$ ist eine Linearkombination von Variablen, die im wesentlichen den Torso beschreiben. Die zweite Hauptkomponente kontrastiert den Oberkörper bis zum Bauch mit den unteren Teil des Torsos. Die dritte unterscheidet das Bauchvolumen vom Rest des Torsos.

Aufgabe 4.10	
	Skizzieren Sie für die nachfolgenden Komponenten $PC4, \dots, PC10$, welche Körpergeometrie durch sie beschrieben wird.

Der Versuch, das errechnete Volumen durch die approximativen Körperteil-Volumen darzustellen, ergibt für den Trainingsteil ein hohes Bestimmtheitsmaß.

```
Eingabe
```

```
lm.vol <- lm(vol ~ neckvol + chestvol + abdomenvol +
  hipvol + thighvol + kneevol +
  anklevol + bicepvol + forearmvol +
  wristvol,
  data = fat, subset = fat$train)
summary(lm.vol)
```

```
Ausgabe
```

```
Call:
lm(formula = vol ~ neckvol + chestvol + abdomenvol + hipvol +
    thighvol + kneevol + anklevol + bicepvol + forearmvol + wristvol,
    data = fat, subset = fat$train)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-5.7799 -1.1548  0.1726  1.1230  4.4839
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.869e-01  1.519e+00  0.386 0.699789
neckvol     1.450e-05  8.860e-06  1.637 0.103805
chestvol    9.449e-06  1.320e-06  7.156 3.58e-11 ***
abdomenvol  1.213e-05  1.154e-06 10.514 < 2e-16 ***
hipvol      7.830e-06  2.009e-06  3.897 0.000147 ***
thighvol    1.497e-05  3.425e-06  4.373 2.30e-05 ***
kneevol     -5.814e-06  9.550e-06 -0.609 0.543548
anklevol    5.387e-05  1.405e-05  3.834 0.000186 ***
bicepvol    2.246e-05  8.691e-06  2.584 0.010744 *
forearmvol  3.047e-05  9.577e-06  3.182 0.001783 **
wristvol    9.623e-06  4.106e-05  0.234 0.815029
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1.867 on 148 degrees of freedom
```

Multiple R-squared: 0.9777, Adjusted R-squared: 0.9762
 F-statistic: 649.4 on 10 and 148 DF, p-value: < 2.2e-16

Bei Einbeziehung auch der linearen Variablen können wir im Trainingsteil des Bestimmtheitsmaß nur geringfügig erhöhen.

Mit den Funktionen aus `library(leap)` kann wieder automatisch nach “optimalen” Modellen gesucht werden.

Eingabe

```
library(leaps)
l1 <- leaps(x = fat[, c(6:15, 20:29)], y = fat$vol)
```

Wenn wir versuchen wollen, nicht das Volumen zu schätzen, sondern den Fettanteil als Linearkombination der entsprechenden Komponenten darzustellen, können wir die entsprechenden Hilfsvariablen konstruieren.

Eingabe

```
fat$neckvol <- fat$neckvol / fat$weightkg
fat$chestvol <- fat$chestvol / fat$weightkg
fat$abdomenvol <- fat$abdomenvol / fat$weightkg
fat$hipvol <- fat$hipvol / fat$weightkg
fat$thighvol <- fat$thighvol / fat$weightkg
fat$kneevol <- fat$kneevol / fat$weightkg
fat$anklevol <- fat$anklevol / fat$weightkg
fat$bicepvolf <- fat$bicepvolf / fat$weightkg
fat$forearmvol <- fat$forearmvol / fat$weightkg
fat$wristvol <- fat$wristvol / fat$weightkg
```

Wir beginnen mit einem einfachen Modell. Wir benutzen nur eine Variable (`abdomenvolf`) aus der Gruppe der Variablen, die den Torso beschreibt, und eine der Variablen (`wristvol`) aus den höheren Hauptkomponenten. Damit erreichen wir fast die Genauigkeit des ersten Modells mit dem vollen Variablensatz.

Eingabe

```
lm.volf <- lm(body.fat ~ abdomenvolf + wristvol, data = fat, subset = fat$train)
summary(lm.volf)
```

Ausgabe

Call:
`lm(formula = body.fat ~ abdomenvolf + wristvol, data = fat, subset = fat$train)`

Residuals:

Min	1Q	Median	3Q	Max
-10.4925	-2.8068	0.2003	3.3089	8.5725

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-4.1738092	6.4631443	-0.646	0.519
abdomenvolf	0.0024661	0.0001843	13.378	< 2e-16 ***
wristvol	-0.0313941	0.0055045	-5.703	5.73e-08 ***

 Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 4.282 on 156 degrees of freedom
 Multiple R-squared: 0.6888, Adjusted R-squared: 0.6848
 F-statistic: 172.7 on 2 and 156 DF, p-value: < 2.2e-16

Aufgabe 4.11	
*	Ergänzen Sie die Variablen durch andere volumenbezogene Variable in dem obigen Modell. gewinnen Sie an Schätzgenauigkeit?
**	Versuchen Sie, die Variable <code>age</code> in die Modellierung mit einzubeziehen. Wie berücksichtigen Sie <code>age</code> im Modell?
**	Die Funktion <code>mvr()</code> in <code>library(pls)</code> steht bereit, um Regressionen auf der Basis der Hauptkomponenten durchzuführen. Benutzen Sie die Funktion zur Regression. Wie unterscheidet sie sich von der gewöhnlichen Kleinste-Quadrate-Regression?

Zur Konstruktion des Modells haben wir den Trainingsteil benutzt. Die Genauigkeit des so gewonnenen Modells überprüfen wir nun am Auswertungsteil. Dazu benutzen wir die Funktion `predict.lm()`, die ein mit `lm()` geschätztes lineares Modell auf einen neuen Datensatz anwendet. Z.B.

Eingabe

```
fat.eval <- fat[fat$train == FALSE, ]
pred <- predict.lm(lm.volf, fat.eval, se.fit = TRUE)
```

Aufgabe 4.12	
*	Schätzen Sie die Genauigkeit des Modells durch die Daten des Auswertungsteils.
*	Führen Sie die eine Diagnostik des gewonnen Modells anhand der Daten des Auswertungsteils durch.

4.8. Hohe Dimensionen

Probleme in kleinen Dimensionen können wir umfassend darstellen und analysieren. Höhere Dimensionen erfordern es oft, eine spezielle Analyse-Strategie zu entwerfen. Die formale Anwendung von Standard-Methoden kommt hier schnell an ihre Grenzen.

Höhere Dimensionen, etwa von 10 bis 100, sind in vielen Anwendungsbereichen üblich. Aber auch Probleme in großen Dimensionen sind alltäglich. Wir müssen uns darüber im Klaren sein, dass die Dimension eine Frage der Modellierung ist, nicht nur eine Frage des Problems. Digitales Video (DV PAL) zum Beispiel zeichnet Bilder im Format 720×576 auf. Ein einzelnes Bild mit drei Farben gibt also einen Vektor im $720 \times 576 \times 3 = 1244160$ -dimensionalen Raum, jede Sekunde Video das 25-fache. Haben wir Bilddaten zu bearbeiten, so ist es unsere Wahl, ob wir die Bildbearbeitung als Problem mit Dimension $d = 1244160$, betrachten, oder als Folge von 1244160 (nicht unabhängigen!) Beobachtungen mit $d = 1$.

Beim Übergang von $d = 1244160$ zu $d = 1$ verlagern wir Information, die implizit in den Dimensionen steckt, in Strukturinformation, die folglich modelliert werden muss.

Als Anmerkung: in der Praxis geht man einen Mittelweg. Man zerlegt das Bild in Blöcke, z.B. der Größe 64×64 . Pixel innerhalb eines Blockes werden simultan behandelt. Die Blöcke werden sequentiell behandelt - sichtbar bei der nächsten Störung im Fernsehen.

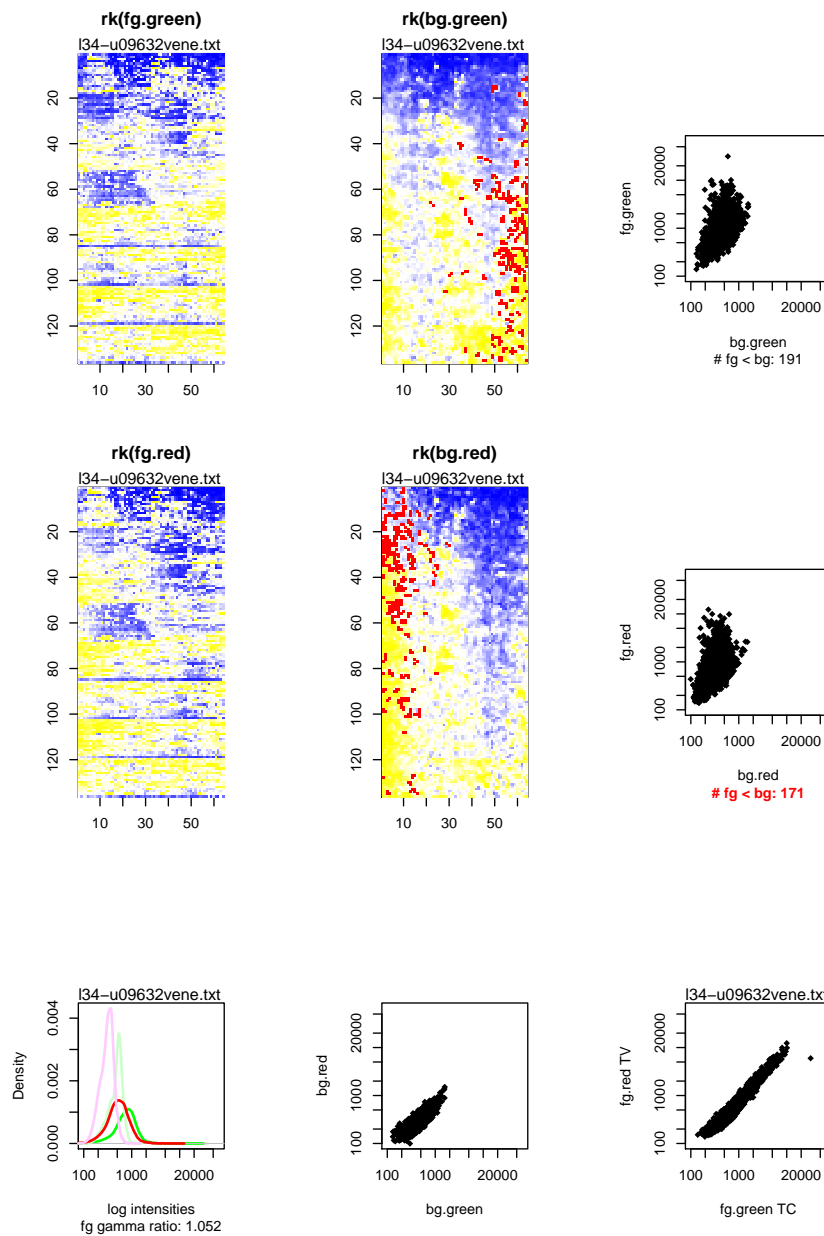


ABBILDUNG 4.6. Ein 4227×4 Datum aus einem Microarray-Experiment

Bei hochdimensionalen Problem ist die Statistik oft gar nicht sichtbar - sie ist versteckt in der Hardware als “imbedded system”.

Die Abbildung 4.6 ist ein Beispiel aus einer Analyse mit R für einen hochdimensionalen Datensatz vom Mikroarray-cDNA-Daten ([Saw02]). Ein einzelnes Datum in diesem Datensatz besteht aus Messungen an 4227 Proben mit jeweils 4 Teilmessungen (*fg.green*, *fg.red*, *bg.green*, *bg.red*). Die wesentliche Funktion, die hier zur Visualisierung benutzt wird, ist *image()*, mit der eine Variable *z* anhand eine Farbtabelle gegen zwei Koordinaten *x*, *y* dargestellt werden kann. Dargestellt ist eine Beobachtung. Die vier Kanäle der Teilmessungen sind nebeneinander gestellt.

Die Farben codieren das Ergebnis einer Voranalyse - die roten Punkte signalisieren Problemzonen auf dem cDNA-Chip. In diesem Fall kann aus dem Muster der Selektion ein spezifisches Problem in der Fertigung identifiziert werden.

Themenorientierte Übersichten über R-Pakete, insbesondere auch zu multivariaten Problemen, sind in <http://cran.at.r-project.org/src/contrib/Views/> zu finden.

4.9. Statistische Zusammenfassung

Die Analyse multivariater Daten konnte in diesem Zusammenhang nur gestreift werden. Multivariate Probleme tauchen implizit schon bei Regressionsproblemen auf (siehe Kapitel 2). Bei den einfachen Regressionsproblemen bezogen sich die multivariaten Aspekte aber nur auf deterministische Parameter. Im allgemeinen Fall haben wir aber eine multivariate statistische Verteilung zu analysieren. An dieser Stelle muss die Einführung abbrechen, und weiteres bleibt weiterführenden Vorlesungen vorbehalten.

4.10. Literatur und weitere Hinweise:

R als Programmiersprache: Übersicht

R ist eine interpretierte Ausdruckssprache. Ausdrücke sind zusammengesetzt aus Objekten und Operatoren.

A.1. Hilfe und Information

R <i>Hilfe</i>	
<code>help()</code>	Information über ein Objekt/eine Funktion <i>Beispiel:</i> <code>help(help)</code>
<code>args()</code>	Zeigt Argumente einer Funktion
<code>example()</code>	Führt evtl. vorhandene Beispiele aus <i>Beispiel:</i> <code>example(plot)</code>
<code>help.search()</code>	Sucht Information über ein Objekt/eine Funktion
<code>apropos()</code>	Lokalisiert nach Stichwort
<code>demo()</code>	Führt Demos zu einem Themenbereich aus <i>Beispiel:</i> <code>demo(graphics)</code> <code>demo()</code> listet die zur Verfügung stehenden Themenbereiche

A.2. Namen und Suchpfade

Objekte werden durch Namen identifiziert. Anhand des Namens werden Objekte in einer Kette von Suchbereichen identifiziert. Die aktuellen Suchbereiche können mit `search()` inspiziert werden.

R <i>Suchpfade</i>	
<code>search()</code>	Liste der aktuellen Suchbereiche, beginnend mit <code>.GlobalEnv</code> bis hinab zum Basis-Paket <code>package:base</code> . <i>Beispiel:</i> <code>search()</code>
<code>searchpaths()</code>	Liste der Zugriffspfade zu aktuellen Suchbereichen <i>Beispiel:</i> <code>searchpaths()</code>
<code>objects()</code>	Liste der Objekte in einem Suchbereich <i>Beispiele:</i> <code>objects()</code> <code>objects("package:base")</code>

(Fortsetzung)→

R <i>Suchpfade</i> (Fortsetzung)	
<code>ls()</code>	Liste der Objekte in einem Suchbereich <i>Beispiele:</i> <code>ls()</code> <code>ls("package:base")</code>
<code>ls.str()</code>	Liste der Objekte und ihrer Struktur in einem Suchbereich <i>Beispiele:</i> <code>ls.str()</code> <code>lsf.str("package:base")</code>
<code>find()</code>	Lokalisiert nach Stichwort. Findet auch überlagerte Einträge <i>Aufruf:</i> <code>find(what, mode = "any", numeric = FALSE, simple.words = TRUE)</code>
<code>apropos()</code>	Lokalisiert nach Stichwort. Findet auch überlagerte Einträge <i>Aufruf:</i> <code>apropos(what, where = FALSE, ignore.case = TRUE, mode = "any")</code>

Funktionen können sowohl bei Definition als auch bei Aufruf geschachtelt sein. Dies macht eine Erweiterung der Suchpfade nötig. Die dynamische Identifikation von Objekten benutzt Umgebungen (environments), um in Funktionen lokale Variable oder globale Variablen aufzulösen.

R <i>Suchpfade</i> (Fortsetzung)	
<code>environment()</code>	Aktuelle Auswertungsumgebung <i>Beispiel:</i> <code>environment()</code>
<code>sys.parent()</code>	Vorausgehende Auswertungsumgebungen <i>Beispiel:</i> <code>sys.parent(1)</code>

Objekte haben zwei implizite Attribute, die erfragt werden mit `mode()` und `length()`. Die Funktion `typeof()` gibt den (internen) Speichermodus eines Objektes an.

Ein `class`-Attribut benennt die Klasse eines Objektes.

A.3. Anpassung

R bietet eine Reihe von Möglichkeiten, das System zu konfigurieren, so dass beim Start und beim Ende bestimmte Kommandos ausgeführt werden. Falls vorhanden, werden beim Start die Dateien `.Rprofile` und `.RData` eingelesen und ausgewertet. Details können system-spezifisch sein. Die jeweils spezifische Information erhält man mit `help(Startup)`.

A.4. Basis-Datentypen

R <i>Basis-Datentypen</i>	
<i>numeric</i>	<i>real</i> oder <i>integer</i> . In R: <i>real</i> ist stets doppelt-genau. Einfache Genauigkeit wird für externe Aufrufe zu anderen Sprachen mit <i>.C</i> oder <i>.FORTRAN</i> unterstützt. Funktionen wie <i>mode()</i> und <i>typedef()</i> können je nach Implementierung auch den Speicherungsmodus (<i>single</i> , <i>double</i> ...) melden. <i>Beispiele:</i> 1.0 2 3.14E0
<i>complex</i>	komplex, in cartesischen Koordinaten <i>Beispiel:</i> 1.0+0i
<i>logical</i>	TRUE, FALSE. In R: auch vordefinierte Variable T, F. In S-Plus sind T und F Basis-Objekte.
<i>character</i>	Zeichenketten. Delimiter sind alternativ " oder '. <i>Beispiel:</i> T", 'klm'
<i>list</i>	Allgemeine Liste. Die Listenelemente können auch von unterschiedlichem Typ sein. <i>Beispiel:</i> list(1:10, "Hello")
<i>function</i>	R-Funktion <i>Beispiel:</i> sin
<i>NULL</i>	Spezialfall: leeres Objekt <i>Beispiel:</i> NULL

Zusätzlich zu den Konstanten TRUE und FALSE gibt es drei spezielle Werte für Ausnahmesituationen:

<i>spezielle Konstanten</i>	
<i>TRUE</i>	Alternativ: T. Typ: logical.
<i>FALSE</i>	Alternativ: F. Typ: logical.
<i>NA</i>	“not available”. Typ: logical. NA ist von TRUE und FALSE verschieden
<i>NaN</i>	“not a valid numeric value”. Implementationsabhängig. Sollte dem IEEE Standard 754 entsprechen. Typ: numeric. <i>Beispiel:</i> 0/0

(Fortsetzung)→

<i>spezielle Konstanten</i> (Fortsetzung)	
<i>Inf</i>	unendlich. Implementationsabhängig. Sollte dem IEEE Standard 754 entsprechen. Typ: numeric. <i>Beispiel:</i> 1/0

A.5. Ausgabe von Objekten

Die Objekt-Attribute und weitere Eigenschaften können abgefragt oder mit Ausgaberroutinen angefordert werden. Die Ausgaberroutinen sind in der Regel *polymorph*, d.h. sie erscheinen in Varianten, die den jeweiligen Objekten angepasst werden.

R <i>Inspektion</i>	
<code>print()</code>	Standard-Ausgabe
<code>structure()</code>	Ausgabe, optional mit Attributen
<code>summary()</code>	Standard-Ausgabe als Übersicht, insbesondere für Modellanpassungen
<code>plot()</code>	Standard-Grafikausgabe

A.6. Inspektion von Objekten

Die folgende Tabelle fasst die wichtigsten Informationsmöglichkeiten über Objekte zusammen.

<i>Inspektion von Objekten</i>	
<code>str()</code>	Stellt die interne Struktur eines Objekts in kompakter Form dar. <i>Aufruf:</i> <code>str(<object>)</code>
<code>structure()</code>	Stellt die interne Struktur eines Objekts dar. Dabei können Attribute für die Darstellung als Parameter übergeben werden. <i>Beispiel:</i> <code>structure(1:6, dim = 2:3)</code> <i>Aufruf:</i> <code>structure(<object>, ...)</code>
<code>class()</code>	Objekt-Klasse. Bei neueren Objekten ist die Klasse als Attribut gespeichert. In älteren S oder R-Versionen ist sie durch Typ und andere Attribute implizit bestimmt.
<code>mode()</code>	Modus (Typ) eines Objekts.
<code>storage.mode()</code>	Speichermodus eines Objekts.
<code>typeof()</code>	Modus eines Objekts. Kann vom Speichermodus abweichen. Je nach Implementierung kann etwa eine numerische Variable standardmäßig doppelt- oder einfach genau abgespeichert werden.
<code>length()</code>	Länge = Anzahl der Elemente
<code>attributes()</code>	Liest/setzt Attribute eines Objekts, wie z.B. Namen, Dimensionen, Klassen.
<code>names()</code>	Namen-Attribut für Elemente eines Objekts, z.B. eines Vektors. <i>Aufruf:</i> <code>names(<obj>)</code> gibt das Namen-Attribut von <code><obj></code> . <code>names(<obj>)<-<charvec></code> setzt es. <i>Beispiel:</i> <code>x<-values</code> <code>names(x)<- <charvec></code>

A.7. Inspektion des Systems

Die folgende Tabelle fasst die wichtigsten Informationsmöglichkeiten über die allgemeine Systemumgebung zusammen.

<i>System-Inspektion</i>	
<code>search()</code>	aktueller Suchpfad
<code>ls()</code>	aktuelle Objekte
<code>methods()</code>	generische Methoden <i>Aufruf:</i> <code>methods(<fun>)</code> zeigt spezialisierte Funktionen zu <code><fun></code> , <code>methods(class = <c>)</code> die klassenspezifischen Funktionen zu class <code><c></code> . <i>Beispiele:</i> <code>methods(plot)</code> <code>methods(class = lm)</code>
<code>data()</code>	zugreifbare Daten
<code>library()</code>	zugreifbare Bibliotheken
<code>help()</code>	allgemeines Hilfe-System
<code>options()</code>	globale Optionen
<code>par()</code>	Parameter-Einstellungen des Grafik-Systems

Die Optionen des `lattice`-Systems können mit `trellis.par.set()` bzw. `lattice.options()` kontrolliert werden.

R ist im umgebenden Betriebssystem verankert. Einige Variable, wie z.B. Zugriffspfade, Zeichencodierung etc. werden von dort übernommen.

<i>System-Umgebung</i>	
<code>getwd()</code>	aktuelles Arbeitsverzeichnis
<code>setwd()</code>	setzt aktuelles Arbeitsverzeichnis
<code>dir()</code>	listet Dateien im aktuellen Arbeitsverzeichnis
<code>system()</code>	ruft System-Funktionen auf

A.8. Komplexe Datentypen

Die Interpretation von Basistypen oder abgeleiteten Typen kann durch ein oder mehrere `class`-Attribute spezifiziert werden. Polymorphe Funktionen wie `print` oder `plot` werten dieses Attribut aus und rufen nach Möglichkeit entsprechend der Klasse spezialisierte Varianten auf (Siehe 2.6.5Seite 2-39).

Zur Speicherung von Datumsangaben und Zeiten stehen entsprechende Klassen bereit. Nähere Information zu diesen Datentypen erhält man mit

`help(DateTimeClasses)`.

R ist vektor-basiert. Einzelne Konstanten oder Werte sind nur Vektoren der speziellen Länge 1. Sie genießen keine Sonderbehandlung.

<i>Zusammengesetzte Objekttypen</i>	
Vektoren	R Basis-Datentypen
Matrizen	Vektoren mit zwei-dimensionalem Layout
Arrays	<p>Vektoren mit höherdimensionalem Layout</p> <p><code>dim()</code> definiert Dimensionsvektor</p> <p><i>Beispiel:</i> <code>x <- runif(100)</code> <code>dim(x) <- c(5, 5, 4)</code></p> <p><code>array()</code> konstruiert neuen Vektor mit gegebener Dimensionsstruktur</p> <p><i>Beispiel:</i> <code>z <- array(0, c(4, 3, 2))</code></p> <p><code>rbind()</code> kettet Reihen an</p> <p><code>cbind()</code> kettet Spalten an</p>
Faktoren	<p>Sonderfall für kategorielle Daten</p> <p><code>factor()</code> wandelt Vektor in Faktor um</p> <p><i>Siehe auch</i> Abschnitt 2.2.1</p> <p><code>ordered()</code> wandelt Vektor im Faktor mit geordneten Stufen um. Dies ist eine Abkürzung für <code>factor(x, ..., ordered = TRUE)</code></p> <p><code>levels()</code> gibt die Stufen eines Faktors an</p> <p><i>Beispiel:</i> <code>x <- c("a", "b", "a", "c", "a")</code> <code>xf <- factor(x)</code> <code>levels(xf)</code> ergibt <code>[1] "abc"</code></p> <p><code>tapply()</code> wendet eine Funktion getrennt für alle Stufen von Faktoren einer Faktorliste an</p>
Listen	Analog Vektoren, mit Elementen auch unterschiedlichen Typs

(Fortsetzung)→

Zusammengesetzte Objekttypen (Fortsetzung)	
	<pre>list() erzeugt Liste Aufruf: list(<Komponenten>) [[]] Indexweiser Zugriff auf Komponenten Liste\$Komponente Zugriff nach Namen Beispiel: l <- list(name = "xyz", age = 22, fak = math") > l[[2]] 22 > l\$age 22</pre>
Datenrahmen	<p>data frames Analog Arrays bzw. Listen, mit spaltenweise einheitlichem Typ und einheitlicher Spaltenlänge</p> <pre>data.frame() analog list(), aber Restriktionen müssen erfüllt sein attach() fügt Datenrahmen in die aktuelle Suchliste ein, d.h. für Komponenten reicht der Komponentename. detach()</pre>

A.9. Zugriff auf Komponenten

Die Länge von Vektoren ist ein dynamisches Attribut. Sie wird bei Bedarf erweitert und gekürzt. Insbesondere gilt implizit eine "Recycling-Regel": Hat ein Vektor nicht die erforderliche Länge für eine Operation, so wird er periodisch bis zur erforderlichen Länge wiederholt.

Auf Vektor-Komponenten kann über Indizes zugegriffen werden. Die Indizes können explizit oder als Regel-Ausdruck angegeben werden.

R <i>Index-Zugriff</i>	
<code>x[⟨indices⟩]</code>	Indizierte Komponenten von x <i>Beispiel:</i> <code>x[1:3]</code>
<code>x[-⟨indices⟩]</code>	x ohne indizierte Komponenten <i>Beispiel:</i> <code>x[-3]</code> x ohne 3. Komponente
<code>x[⟨condition⟩]</code>	Komponenten von x , für die $\langle \text{condition} \rangle$ gilt. <i>Beispiel:</i> <code>x[x<0.5]</code>

Vektoren (und andere Objekte) können auf höherdimensionale Konstrukte abgebildet werden. Die Abbildung wird durch zusätzliche Dimensions-Attribute beschrieben. Nach Konvention erfolgt eine spaltenweise Einbettung, d.h. der erste Index variiert zuerst (FORTRAN-Konvention). Operatoren und Funktionen können die Dimensions-Attribute auswerten.

R <i>Index-Zugriff</i>	
<code>dim()</code>	Setzt oder liest die Dimensionen eines Objekts <i>Beispiel:</i> <code>x <- 1:12 ; dim(x) <- c(3, 4)</code>
<code>dimnames()</code>	Setzt oder liest Namen für die Dimensionen eines Objekts
<code>nrow()</code>	Gibt die Anzahl der Zeilen = Dimension 1
<code>ncol()</code>	Gibt die Anzahl der Spalten = Dimension 2
<code>matrix()</code>	Erzeugt eine Matrix mit vorgegebenen Spezifikationen <i>Aufruf:</i> <code>matrix(data = NA, nrow = 1, ncol = 1, byrow = FALSE, dimnames = NULL)</code> <i>Siehe auch</i> Beispiel 1.8 (Seite 1-21)
<code>array()</code>	Erzeugt eine evtl. höherdimensionale Matrix <i>Beispiel:</i> <code>array(x, dim = length(x), dimnames = NULL)</code>

R <i>Array-Zugriffe</i>	
<code>cbind()</code> <code>rbind()</code>	Verkettet Zeilen bzw. Spalten

(Fortsetzung)→

R <i>Array-Zugriffe</i> (Fortsetzung)	
<code>split()</code>	Teilt einen Vektor nach Faktoren auf
<code>table()</code>	Erzeugt eine Tabelle von Besetzungszahlen

R <i>Iteratoren</i>	
<code>apply()</code>	wendet eine Funktion auf die Zeilen oder Spalten einer Matrix an <i>Aufruf:</i> <code>apply(x, MARGIN, FUNCTION, ...)</code> Margin = 1: Zeilen, Margin = 2: Spalten. <i>Siehe auch</i> Beispiel 1.8 (Seite 1-21)
<code>lapply()</code>	wendet eine Funktion auf die Elemente einer Liste an <i>Aufruf:</i> <code>lapply(X, FUN, ...)</code>
<code>sapply()</code>	wendet eine Funktion auf die Elemente einer Liste, eines Vektors oder einer Matrix an. Falls mögliche werden Dimensionsnamen übernommen. <i>Aufruf:</i> <code>sapply(X, FUN, ..., simplify = TRUE, USE.NAMES = TRUE)</code>
<code>tapply()</code>	wendet eine Funktion auf Komponenten eines Objekts in Abhängigkeit von einer Liste von kontrollierenden Faktoren an.
<code>by()</code>	Objekt-orientierte Variante von <code>tapply</code> <i>Aufruf:</i> <code>by(data, INDICES, FUN, ...)</code>
<code>aggregate()</code>	Berechnet Statistiken für Teilmengen <i>Aufruf:</i> <code>aggregate(x, ...)</code>
<code>replicate()</code>	Wertet eine Ausdruck wiederholt aus (z. Bsp. mit Erzeugung von Zufallszahlen zur Simulation). <i>Aufruf:</i> <code>replicate(n, expr, simplify = TRUE)</code>
<code>outer()</code>	erzeugt eine Matrix mit allen Paar-Kombinationen aus zwei Vektoren, und wendet eine Funktion auf jedes Paar an. <i>Aufruf:</i> <code>outer(vec1, vec2, FUNCTION, ...)</code>

A.10. Tabellen-Transformationen

<i>Transformationen</i>	
<code>seq()</code>	Erzeugt eine Sequenz
<code>abbreviate()</code>	

<i>Transformationen</i>	
<code>duplicated()</code>	Prüft auf mehrfach auftretende Werte
<code>unique()</code>	Erzeugt Vektor ohne mehrfach auftretende Werte
<code>match()</code>	Gibt Position eines Werts in einem Vektor
<code>pmatch()</code>	Partielles Matching

<i>Zeichenketten-Transformationen</i>	
<code>casefold()</code>	Wandelt in Klein- oder Großbuchstaben um
<code>tolower()</code>	Wandelt in Kleinbuchstaben um
<code>toupper()</code>	Wandelt in Großbuchstaben um
<code>chartr()</code>	Übersetzt Zeichen in einem Zeichen-Vektor
<code>substring()</code>	

<i>Transformationen</i>	
<code>table()</code>	Erzeugt eine Kreuztabelle
<code>expand.grid()</code>	Erzeugt einen Datenrahmen mit allen Kombinationen gegebener Faktoren
<code>reshape()</code>	Wandelt zwischen einer Kreuztabelle (Spalte pro Variable) und einer langen Tabelle (Variablen in Zeilen, mit zusätzlicher Indikator-Spalte) um
<code>merge()</code>	Kombiniert Datenrahmen

<i>Transformationen</i>	
-------------------------	--

<i>Transformationen</i>	
<code>t()</code>	Transponiert Zeilen und Spalten Aufruf: <code>t(x)</code>

(Fortsetzung)→

<i>Transformationen</i> (Fortsetzung)	
<i>aperm()</i>	Generalisierte Permutation <i>Aufruf:</i> <code>aperm(x, perm)</code> Dabei ist <i>perm</i> eine Permutation der Indizes von <i>x</i> .
<i>split()</i>	Teilt einen Vektor nach einem Faktor auf
<i>unsplit()</i>	Kombiniert Komponenten zu einem Vektor

A.11. Operatoren

Ausdrücke in R können aus Objekten und Operatoren zusammengesetzt sein. Die folgende Tabelle ist nach Vorrang geordnet (höchster Rang oben).

R <i>Basisoperatoren</i>	
\$	Komponenten-Selektion <i>Beispiel:</i> <code>list\$item</code>
[[[Indizierung, Elementzugriff <i>Beispiel:</i> <code>x[i]</code>
^	Potenzierung <i>Beispiel:</i> <code>x^3</code>
-	unitäres Minus
:	Folge-Generierung <i>Beispiele:</i> <code>1:5</code> <code>5:1</code>
%<name>%	spezielle Operatoren. Können auch benutzer-definiert sein. <i>Beispiele:</i> <code>"%deg2%"<-function(a, b) a + b^2</code> <code>2 %deg2% 4</code>
* /	Multiplikation, Division
+ -	Addition, Subtraktion
< > < = > = == ! =	Vergleichsoperatoren
!	Negation
& &&	und, oder && , sind "Shortcut"-Operatoren
<- ->	Zuweisung

Haben die Operanden nicht die gleiche Länge, so wird der kürzere Operand zyklisch wiederholt.

Operatoren der Form %<name>% können vom Benutzer definiert werden. Die Definition folgt den Regeln für Funktionen.

Ausdrücke können als Folge mit trennendem Semikolon geschrieben werden. Ausdrucksgruppen können durch {...} zusammengefasst werden.

A.12. Funktionen

Funktionen sind spezielle Objekte. Funktionen können Resultat-Objekte übergeben.

R <i>Funktions-</i> <i>deklarationen</i>	
Deklaration	<code>function (<formale Parameterliste>)</code> <Ausdruck> <i>Beispiel:</i> <code>fak <- function(n) prod(1:n)</code>
Formale Parameter	<Parametername> <Parametername> = <Default-Wert>
Formale Parameterliste	Liste von formalen Parametern, durch Komma getrennt <i>Beispiele:</i> <code>n, mean = 0, sd = 1</code>
...	Variable Parameterliste. Variable Parameterlisten können innerhalb von Prozeduren weitergegeben werden. <i>Beispiel:</i> <code>mean.of.all <- function (...)mean(c(...))</code>
Funktions-Resultate	<code>return <Wert></code> bricht Funktionsauswertung ab und übergibt Wert
	<Wert> als letzter Ausdruck in einer Funktionsdeklaration: übergibt Wert
Funktions-Resultate	<Variable><<-<Wert> übergibt Wert. Normalerweise wirken Zuweisungen nur auf lokale Kopien der Variablen. Die Zuweisung mit <<- jedoch sucht die Zielvariable in der gesamten Umgebungshierarchie.

R <i>Funktions-</i> <i>aufruf</i>	
Funktionsaufruf	<Name>(<Aktuelle Parameterliste>) <i>Beispiel:</i> <code>fak(3)</code>
Aktuelle Parameterliste	Werte werden zunächst der Position nach zugeordnet. Abweichend davon können Namen benutzt werden, um Werte gezielt zuzuordnen. Dabei reichen die Anfangsteile der Namen (Ausnahme: nach einer variablen Parameterliste müssen die Namen vollständig angegeben werden). Mit der Funktion <code>missing()</code> kann überprüft werden, ob für einen formalen Parameter ein entsprechender aktueller Parameter fehlt. <i>Aufruf:</i> <code><Werteliste></code> <code><Parametername> = <Werte></code> <i>Beispiel:</i> <code>rnorm(10, sd = 2)</code>

Parameter bei Funktionen werden dem Wert nach übergeben. Soll der damit verbundene Aufwand vermieden werden, so kann mit Hilfe der *environment*-Information direkt auf Variable zugegriffen werden. Entsprechende Techniken sind in [GI00] beschrieben.

Spezialfall: Funktionen mit Namen der Form `xxx<-` erweitern die Zuweisungsfunktion.

Beispiel:

```
"inc<-" <-function (x, value) x+value
x <- 10
inc(x)<- 3
x
```

In R-Zuweisungsfunktionen **muss** das Wert-Argument "value" heißen.

A.13. Debugging und Profiling

R bietet eine Reihe von Werkzeugen zur Identifizierung von Fehlern. Diese sind besonders im Zusammenhang mit Funktionen hilfreich. Mit `browser()` kann in einen Browser-Modus geschaltet werden. In diesem Modus sind die üblichen R-Anweisungen möglich. Daneben gibt es eine kleine Zahl von speziellen Anweisungen. Der Browser-Modus kann mit `debug()` automatisch bei Eintritt in eine Funktion aktiviert werden. Durch den speziellen Prompt `Browse[xx]>` ist der Browser-Modus erkennbar.

`<return>`: geht zur nächsten Anweisung, falls die Funktion unter `debug`-Kontrolle steht. Fährt mit der Anweisungsausführung fort, falls `browser` direkt aufgerufen wurde.

`n`: geht zur nächsten Anweisung (auch falls `browser` direkt aufgerufen wurde).

`cont`: Fährt mit der Anweisungsausführung fort.

`c`: Kurzform für `cont`. Fährt mit der Anweisungsausführung fort.

`where`: Zeigt Aufrufverschachtelung.

`Q`: Stoppt Ausführung und springt in Grundzustand zurück.

<i>Debug-Hilfen</i>	
<code>browser()</code>	Hält die Ausführung an und geht in den Browser-Modus. <i>Aufruf:</i> <code>browser()</code>
<code>recover()</code>	<code>recover()</code> zeigt eine Liste der aktuellen Aufrufe, aus der einer zur <code>browser()</code> -Inspektion gewählt werden kann. Mit <code>c</code> kehrt man aus dem <code>browser</code> zu <code>recover</code> zurück. Mit <code>0</code> verlässt man <code>recover()</code> <i>Aufruf:</i> <code>recover()</code> <i>Hinweis:</i> Mit <code>options(error = recover)</code> kann die Fehlerbehandlung so konfiguriert werden, dass im Fehlerfalle automatisch <code>browser()</code> aufgerufen wird.
<code>debug()</code>	Markiert eine Funktion zur Debugger-Kontrolle. Bei nachfolgenden Aufrufen der Funktion wird der Debugger aktiviert und schaltet in den Browser-Modus. <i>Aufruf:</i> <code>debug(<Funktion>)</code>
<code>undebug()</code>	Löscht Debugger-Kontrolle für eine Funktion. <i>Aufruf:</i> <code>undebug(<Funktion>)</code>
<code>trace()</code>	Markiert eine Funktion zur Trace-Kontrolle. Bei nachfolgenden Aufrufen der Funktion wird der Aufruf mit seinen Argumenten angezeigt. <i>Aufruf:</i> <code>trace(<Funktion>)</code>

(Fortsetzung)→

Debug-Hilfen	
(Fortsetzung)	
<code>untrace()</code>	Löscht Trace-Kontrolle für eine Funktion. <i>Aufruf:</i> <code>untrace(<Funktion>)</code>
<code>traceback()</code>	Im Fehlerfall innerhalb einer Funktion wird die aktuelle Aufrufverschachtelung in einer Variablen <code>.Traceback</code> gespeichert. <code>traceback()</code> wertet diese Variable aus und zeigt den Inhalt an. <i>Aufruf:</i> <code>traceback()</code>
<code>try()</code>	Erlaubt benutzer-definierte Fehlerbehandlung. <i>Aufruf:</i> <code>try(<Ausdruck>)</code>

Um die Laufzeit in einzelnen Bereichen zu messen, bietet R ein “profiling”, das jedoch nur verfügbar ist, wenn R mit den entsprechenden Optionen kompiliert worden ist. Die beim Compilieren benutzten Informationen können mit `capabilities()` erfragt werden.

Profiling-Hilfen	
<code>system.time()</code>	Misst die Ausführungszeit einer Anweisung. Diese Funktion ist stets verfügbar. <i>Aufruf:</i> <code>system.time(<expr>, <gcFirst>)</code>
<code>Rprof()</code>	Registriert periodisch die jeweils aktiven Funktionen. Diese Funktion ist nur verfügbar, wenn R für “profiling” kompiliert ist. Mit <code>memory.profiling = TRUE</code> wird außer der Zeit auch periodisch die Speicherplatznutzung protokolliert. Diese Option ist nur verfügbar, wenn R entsprechend kompiliert ist. <i>Aufruf:</i> <code>Rprof(filename = "Rprof.out", append = FALSE, interval = 0.02, memory.profiling = FALSE)</code>
<code>Rprofmem()</code>	Registriert Speicherplatz-Anforderungen im Anforderungsfall. Diese Funktion ist nur verfügbar, wenn R für “memory profiling” kompiliert ist. <i>Aufruf:</i> <code>Rprofmem(filename = "Rprofmem.out", append = FALSE, threshold = 0)</code>
<code>summaryRprof()</code>	Fasst die Ausgabe von <code>Rprof()</code> zusammen und berichtet den Zeitbedarf je Funktion. <i>Aufruf:</i> <code>summaryRprof(filename = "Rprof.out", chunksize = 5000, memory = c("none", "both", "tseries", "stats"), index = 2, diff = TRUE, exclude = NULL)</code>

A.14. Kontrollstrukturen

R <i>Kontrollstrukturen</i>	
<i>if</i>	<p>Bedingte Ausführung</p> <p><i>Aufruf:</i> <code>if (<log. Ausdruck 1>) <Ausdruck2></code> Der logische Ausdruck 1 darf nur einen logischen Wert ergeben. Für vektorisierten Zugriff benutze man <i>ifelse</i>.</p> <p><i>Aufruf:</i> <code>if (<log. Ausdruck1>) <Ausdruck2> else <Ausdruck3></code></p>
<i>ifelse</i>	<p>Elementweise bedingte Ausführung</p> <p><i>Aufruf:</i> <code>ifelse(<log. Ausdruck1>, <Ausdruck2>, <Ausdruck3>)</code> Wertet den logischen Ausdruck 1 elementweise auf einen Vektor an, und übergibt bei wahren Resultat den elementweisen Wert von Ausdruck2, sonst von Ausdruck3)</p> <p><i>Beispiel:</i> <code>trimmedX <- ifelse (abs(x)<2, X, 2)</code></p>
<i>switch</i>	<p>Auswahl aus einer Liste von Alternativen</p> <p><i>Aufruf:</i> <code>switch(<Ausdruck1>, ...)</code> Ausdruck1 muss einen numerischen Wert oder eine Zeichenkette ergeben. ... ist eine explizite Liste der Alternativen.</p> <p><i>Beispiel:</i> <code>centre <- function (x , type) { switch(type, mean = mean(x), median = median(x), trimmed = mean(x, trim = .1)) }</code></p>
<i>for</i>	<p>Iteration (Schleife)</p> <p><i>Aufruf:</i> <code>for (<name> in <Ausdruck1>) <Ausdruck2></code></p>
<i>repeat</i>	<p>Wiederholung. Muss z.B. mit <i>break</i> verlassen werden.</p> <p><i>Aufruf:</i> <code>repeat <Ausdruck></code></p> <p><i>Beispiel:</i> <code>pars<-init repeat { res<- get.resid (data, pars) if (converged(res)) break pars<-new.fit (data, pars) }</code></p>
<i>while</i>	<p>Bedingte Wiederholung</p> <p><i>Aufruf:</i> <code>while (<log. Ausdruck>) <Ausdruck></code></p> <p><i>Beispiel:</i> <code>pars<-init; res <- get.resid (data, pars) while (!converged(res)) { pars<-new.fit(data, pars) res<- get.resid }</code></p>
<i>break</i>	verlässt die aktuelle Schleife
<i>next</i>	verlässt einen Schleifenzyklus und springt zum nächsten

A.15. Verwaltung und Anpassung

<i>objects()</i> <i>ls()</i>	Liste der aktuellen Objekte
<i>rm()</i>	Löscht die angegebenen Objekte <i>Aufruf:</i> <i>rm</i> (<Objektliste>)

A.16. Ein- und Ausgabe in Ausgabeströme

R <i>Ein/Ausgabe</i>	
<code>write()</code>	Schreibt Daten in eine Datei. <i>Aufruf:</i> <code>write(val, file)</code> <i>Beispiel:</i> <code>write(x, file = "data")</code>
<code>source()</code>	Führt die R-Anweisungen aus der angegebenen Datei aus. <i>Aufruf:</i> <code>source("<Dateiname> ")</code> <i>Beispiel:</i> <code>source(cmnds.R)</code>
<code>sink()</code>	Lenkt Ausgaben in die angegebene Datei. <i>Aufruf:</i> <code>sink("<Dateiname>")</code> <i>Beispiel:</i> <code>sink()</code> lenkt die Ausgabe wieder auf die Konsole.
<code>dump()</code>	Schreibt für ein Objekt die definierenden Kommandos. Mit <code>source()</code> kann aus der Ausgabe das Objekt regeneriert werden <i>Aufruf:</i> <code>dump(list, file = "<dumpdata.R>", append = FALSE)</code>

A.17. Externe Daten

Zum Editieren und für die Eingabe nach Spreadsheet-Art innerhalb von R gibt es `edit()` (früherer Name: `data.entry()`).

Für den Austausch müssen die Datenformate zwischen allen Beteiligten abgestimmt sein. Zum Import aus Datenbanken und anderen Paketen steht eine Reihe von Bibliotheken zur Verfügung, z.B. `stataread` für Stata, `foreign` für SAS, Minitab und SPSS, `RODBC` für SQL. Weitere Information findet sich im Manual "Data Import/Export" ([R D07b]).

Innerhalb von R werden vorbereitete Daten üblicherweise als `data frames` bereitgestellt. Sind zusätzliche Objekte wie Funktionen oder Parameter nötig, so können sie gebündelt als Paket bereit gestellt werden (siehe Aufgabe A.18 (Seite A-31)).

Für den Austausch zu R kann ein spezielles Austauschformat benutzt werden. Dateien in diesem Format können mit `save()` generiert werden und haben konventionell die Namensendung `.Rda`. Diese Dateien werden mit `load()` wieder geladen.

Daten werden allgemeiner mit der Funktion `data()` geladen. Abhängig von der Endung des Dateinamens der Eingabedatei verzweigt `data()` in mehreren Spezialfällen. Neben den `.Rda` sind übliche Endungen für reine Daten-Eingabedateien `.tab` oder `.txt`. Die online-help-Funktion `help(data)` gibt weitere Auskunft.

<i>Ein- Ausgabe von Daten für R</i>	
<code>save()</code>	Speichert Daten in externe Datei. <i>Aufruf:</i> <code>save(<Namen der zu speichernden Objekte>, file = <Dateiname>, ...)</code>
<code>load()</code>	Lädt Daten aus exterener Datei. <i>Aufruf:</i> <code>load(file = <Dateiname>, ...)</code>
<code>data()</code>	Lädt Daten. <code>data()</code> kann unterschiedliche Formate verarbeiten, wenn die Zugriffspfade und Datei-Namen den R-Konventionen folgen. <i>Aufruf:</i> <code>data(..., list = character(0), package = c(.packages(), .Autoloaded), lib.loc = .lib.loc)</code> <i>Beispiel:</i> <code>data(crimes) # lädt den Datensatz 'crimes'</code>

Für den flexiblen Austausch mit anderen Programmen werden Daten in der Regel als Text-Dateien bereitgestellt, nach Möglichkeit

- in Tabellenform,
- nur ASCII-Zeichen (z.B. keine Umlaute!)
- Variablen spaltenweise angeordnet
- Spalten durch Tabulator-Sprünge getrennt.
- evtl. Spaltenüberschriften in Zeile 1
- evtl. Zeilennr. in Spalte 1.

Dafür wird zum Lesen die Funktion `read.table()` und zum Schreiben die Funktion `write.table()` bereitgestellt. Neben `read.table()` gibt es eine Reihe von Varianten, die auf andere gebräuchliche Datenformate abgestimmt sind. Diese sind unter `help(read.table)` aufgeführt.

<i>Ein- Ausgabe von Daten zum Austausch</i>	
<code>read.table()</code>	Liest Daten-Tabelle Aufruf: <code>read.table(file, header = FALSE, sep = "\t", ...)</code> Beispiele: <code>read.table(<Dateiname>, header = TRUE, sep = '\t')</code> Überschriften in Zeile 1, Zeilennr. in Spalte 1 <code>read.table(<Dateiname>, header = TRUE, sep = '\t')</code> keine Zeilennr., Überschriften in Zeile 1,
<code>write.table()</code>	Schreibt Daten-Tabelle Aufruf: <code>write.table(file, header = FALSE, sep = '\t', ...)</code> Beispiele: <code>write.table(<data frame>, <Dateiname>, header = TRUE, sep = '\t')</code> Überschriften in Zeile 1, Zeilennr. in Spalte 1 <code>write.table(<data frame>, <Dateiname>, header = TRUE, sep = '\t')</code> keine Zeilennr., Überschriften in Zeile 1,

Defaultmäßig konvertiert `read.table()` Daten in *factor*-Variable, falls möglich. Dieses Verhalten kann mit dem Parameter `as.is` beim Aufruf von `read.table()` modifiziert werden. Diese Modifikation ist z. B. nötig, um Datums- und Zeitangaben einzulesen, wie in dem folgen Beispiel aus [GP04]:

```
# date col in all numeric format yyyyymmdd
df <- read.table("laketemp.txt", header = TRUE)
as.Date(as.character(df$date), "%Y-%m-%d")
# first two cols in format mm/dd/yy hh:mm:ss
# Note as.is = in read.table to force character
library("chron")
df <- read.table("oxygen.txt", header = TRUE,
as.is = 1:2)
chron(df$date, df$time)
```

Für sequentielles Lesen steht `scan()` zur Verfügung. Dateien mit stellengenau fest vorgegebenem Format können mit `read.fwf()` gelesen werden.

A.18. Libraries, Pakete

Externe Information kann in (Text)-Dateien und Paketen(Packages) gespeichert sein. Bibliotheken und Pakete sind dabei nach speziellen R-Konventionen strukturiert. “Bibliotheken” sind Sammlungen von “Paketen”.

Zusätzliche Funktionen werden in der Regel als Pakete bereitgestellt. Pakete werden mit `library()`

geladen. Im Paket enthaltene Datensätze sind dann direkt auffindbar und werden mit `data()`

(ohne Argument) aufgelistet.

Beispiel:

```
library(nls)
data()
data(Puromycin)
```

<i>Pakete</i>	
<code>library()</code>	Lädt Zusatzpaket <i>Aufruf:</i> <code>library(package, ...)</code> <i>Siehe auch</i> Abschnitt 1.5.6
<code>require()</code>	Lädt Zusatzpaket; gibt Warnung bei Fehler. <i>Aufruf:</i> <code>require(package, ...)</code>
<code>detach()</code>	Gibt Zusatzpaket frei und entfernt es aus dem Suchpfad. <i>Aufruf:</i> <code>detach(<name>)</code>
<code>install.packages()</code>	Installiert Pakete in <code><lib></code> , lädt sie bei Bedarf aus dem Archiv <i>CRAN</i> <i>Aufruf:</i> <code>install.packages(pkgs, lib, CRAN = getOption("CRAN"), ...)</code>
<code>package.manager()</code>	Falls implementiert: Interface zur Verwaltung installierter Pakete. <i>Aufruf:</i> <code>package.manager()</code>
<code>package.skeleton()</code>	Erstellt das Gerüst für ein neues Paket. <i>Aufruf:</i> <code>package.skeleton(name = "<anRpackage>", list, ...)</code>

Detailinformation zur Erstellung von R-Paketen findet man in “Writing R Extensions” ([R D08]).

A.19. Lineare Algebraoperatoren

Für die lineare Algebra sind die wichtigsten Funktionen weitgehend standardisiert und in C-Bibliotheken wie BLAS/ATLAS und Lapack verfügbar. R benutzt diese Bibliotheken und bietet für die wichtigsten Funktionen einen direkten Zugang.

<i>Lineare Algebra</i>	
<i>eigen()</i>	Berechnet Eigenwerte und Eigenvektoren von reellen oder komplexen Matrizen
<i>svd()</i>	Eigenwertzerlegung einer Matrix
<i>qr()</i>	QR-Zerlegung einer Matrix
<i>determinant()</i>	Determinante einer Matrix
<i>solve()</i>	Löst lineare Gleichung

Falls möglich sollten jedoch statistische Funktionen benutzt und der direkte Zugriff auf Funktionen der linearen Algebra vermieden werden.

A.20. Modell-Beschreibungen

Lineare statistische Modelle können durch Angabe einer Design-Matrix X spezifiziert werden und in der allgemeinen Form

$$Y = X\beta + \varepsilon$$

dargestellt werden, wobei die Matrix X jeweils genauer bestimmt werden muß.

R erlaubt es, Modelle auch dadurch zu spezifizieren, dass die Regeln angegeben werden, nach denen die Design-Matrix gebildet wird.

Operator	Syntax	Bedeutung	Beispiel
\sim	$Y \sim M$	Y hängt von M ab	$Y \sim X$ ergibt $E(Y) = a + bX$
$+$	$M_1 + M_2$	M_1 und M_2 einschliessen	$Y \sim X + Z$ $E(Y) = a + bX + cZ$
$-$	$M_1 - M_2$	M_1 einschliessen, aber M_2 ausschliessen	$Y \sim X - 1$ $E(Y) = bX$
$:$	$M_1 : M_2$	Tensorprodukt, d.h. alle Kombinationen von Stufen von M_1 und M_2	
% in %	$M_1 \% \text{ in } \% M_2$	modifiziertes Tensorprodukt	$a + b \% \text{ in } \% a$ entspricht $a + a : b$
$*$	$M_1 * M_2$	“gekreuzt”	$M_1 + M_2$ entspricht $M_1 + M_2 + M_1 : M_2$
$/$	M_1 / M_2	“geschachtelt”: $M_1 + M_2 \% \text{ in } \% M_1$	
\wedge	$M \wedge n$	M mit allen “Interaktionen” bis Stufe n	
$I()$	$I(M)$	Interpretiere M . Terme in M behalten ihre ursprüngliche Bedeutung; das Resultat bestimmt das Modell.	$Y \sim (1 + I(X^2))$ entspricht $E(Y) = a + bX^2$

TABELLE A.53. Wilkinson-Rogers-Notation für lineare Modelle

Die Modell-Spezifikation ist auch für allgemeinere, nicht lineare Modelle möglich.

Beispiele

$$y \sim 1 + x \quad \text{entspricht } y_i = (1 \ x_i)(\beta_1 \ \beta_2)^\top + \varepsilon$$

$$y \sim x \quad \text{Kurzschreibweise für } y \sim 1 + x$$

(Konstanter Term wird implizit angenommen)

$y \sim 0 + x$	entspricht $y_i = x_i \cdot \beta + \varepsilon$
$\log(y) \sim x1 + x2$	entspricht $\log(y_i) = (1 \ x_{i1} \ x_{i2})(\beta_1 \ \beta_2 \ \beta_3)^T + \varepsilon$ (Konstanter Term wird implizit angenommen)
$y \sim A$	Einweg-Varianzanalyse mit Faktor A
$y \sim A + x$	Covarianzanalyse mit Faktor A und Covariable x
$y \sim A * B$	Zwei-Faktor-Kreuz-Layout mit Faktoren A und B
$y \sim A/B$	Zwei-Faktor hierarchisches Layout mit Faktor A und Subfaktor B

Um zwischen verschiedenen Modellen ökonomisch wechseln zu können, steht die Funktion `update()` zur Verfügung.

Modell- Verwaltung	
<code>formula()</code>	extrahiert Modellformel aus einem Objekt
<code>terms()</code>	extrahiert Terme der Modell-Formal aus einem Objekt
<code>contrasts()</code>	spezifiziert Kontraste
<code>update()</code>	Wechsel zwischen Modellen
<code>model.matrix()</code>	Generiert die Design-Matrix zu einem Modell

Anwendungsbeispiel:

```
lm(y ~ poly(x, 4), data = experiment)
```

analysiert den Datensatz "experiment" mit einem linearen Modell für polynomiale Regression vom Grade 4.

Standard- Analysen	
<code>lm()</code>	lineares Modell <i>Siehe auch</i> Kapitel 2
<code>glm()</code>	generalisiertes lineares Modell
<code>nls()</code>	nicht-lineare kleinste Quadrate
<code>nlm()</code>	allgemeine nicht-lineare Minimierung
<code>update()</code>	Wechsel zwischen Modellen
<code>anova()</code>	Varianz-Analyse

A.21. Grafik-Funktionen

R bietet zwei Grafik-Systeme: Das Basis-Grafiksystem von R implementiert ein Modell, das an der Vorstellung von Stift und Papier orientiert ist. Das Lattice-Grafiksystem ist ein zusätzliches zweites Grafiksystem, das an einem Kamera/Objekt-Modell orientiert ist. Information über Lattice erhält man mit `help(lattice)`, eine Übersicht über die Funktionen in Lattice mit `library(help = lattice)`. Informationen über das Basis-Grafiksystem folgen hier.

Grafik-Funktionen fallen im wesentlichen in drei Gruppen:

“high level“-Funktionen. Diese definieren eine neue Ausgabe.

“low level“-Funktionen. Diese modifizieren eine vorhandene Ausgabe.

Parametrisierungen. Diese modifizieren die Voreinstellungen des Grafik-Systems.

A.21.1. high level Grafik.

“high level”	
<code>plot()</code>	Generische Grafikfunktion
<code>pairs()</code>	paarweise Scatterplots
<code>coplot()</code>	Scatterplots, bedingt auf Covariable
<code>qqplot()</code>	Quantil-Quantil-Plot
<code>qqnorm()</code>	Gauß-Quantil-Quantil-Plot
<code>qqline()</code>	fügt eine Linie zu einem Gauß-Quantil-Quantil-Plot hinzu, die durch das erste und dritte Quantil verläuft.
<code>hist()</code>	Histogramm <i>Siehe auch</i> Abschnitt 1.3.2, Seite 1-27
<code>boxplot()</code>	Box&Whisker-Plot
<code>dotplot()</code>	
<code>curve()</code>	Wertet eine Funktion oder einen Ausdruck nach Bedarf aus und zeichnet eine Kurve. <i>Beispiel:</i> <code>curve(dnorm, from = -3, to = 3)</code>
<code>image()</code>	farbcodiertes z gegen x, y
<code>contour()</code>	Contourplot von z gegen x, y
<code>persp()</code>	3D-Fläche

A.21.2. low level Grafik. Die high-level-Funktionen haben in der Regel einen Parameter `add`. Wird beim Aufruf `add = FALSE` gesetzt, so können sie auch benutzt werden, um zu einem vorhandenen Plot Elemente hinzu zu fügen. Daneben gibt es eine Reihe von low-level-Funktionen, die voraussetzen, dass bereits eine Plot-Umgebung geschaffen ist.

“low level”	
-------------	--

(Fortsetzung)→

<i>“low level”</i> (Fortsetzung)	
<code>points()</code>	Generische Funktion. Markiert Punkte an angegebenen Koordinaten. <i>Aufruf:</i> <code>points(x, ...)</code>
<code>lines()</code>	Generische Funktion. Verbindet Punkte an angegebenen Koordinaten. <i>Aufruf:</i> <code>lines(x, ...)</code>
<code>abline</code>	Fügt Linie (in mehreren Darstellungen) zum Plot hinzu. <i>Aufruf:</i> <code>abline(a, b, ...)</code>
<code>polygon()</code>	Fügt Polygon mit spezifizierten Ecken hinzu.
<code>axis()</code>	Fügt Achsen hinzu.

Daneben hat R rudimentäre Möglichkeiten für Interaktion mit Grafik.

<i>Interaktionen</i>	
<code>locator()</code>	bestimmt die Position von Mausklicks. Eine aktuelle Grafik muss definiert sein, bevor <code>locator()</code> benutzt wird. <i>Beispiel:</i> <code>plot(runif(19))</code> <code>locator(n = 3, type = "l")</code>

A.21.3. Annotationen und Legenden. Die high-level-Funktion bieten in der Regel die Möglichkeiten, Standard-Beschriftungen durch geeignete Parameter zu kontrollieren.

`main` = Haupt-Überschrift, über dem Plot

`sub` = Plot-Unterschrift

`xlab` = Beschriftung der x-Achse

`ylab` = Beschriftung der y-Achse

Beschreibungen erhält man mit `help(plot.default)`.

Zur Ergänzung stehen low-level-Funktionen bereit.

<i>“low level”</i>	
<code>title()</code>	Setzt Überschrift, analog high-level-Parametern. <i>Aufruf:</i> <code>title(main = NULL, sub = NULL, xlab = NULL, ylab = NULL, ...)</code>
<code>text</code>	Fügt Text an spezifizierten Koordinaten hinzu. <i>Aufruf:</i> <code>text(x, y = NULL, text, ...)</code>
<code>legend()</code>	Fügt einen Block mit einer Legende hinzu. <i>Aufruf:</i> <code>legend(x, y = NULL, text, ...)</code>

(Fortsetzung)→

“low level” (Fortsetzung)	
<code>mtext()</code>	Fügt Randbeschriftung hinzu. <i>Aufruf:</i> <code>mtext(text, side = 3, ...)</code> . Die Ränder werden bezeichnet durch 1 = unten, 2 = links, 3 = oben, 4 = rechts)

R gibt auch (eingeschränkte) Möglichkeiten zum Formelsatz. Ist der Text-Parameter eine Zeichenkette, so wird sie direkt übernommen. Ist der Text-Parameter ein (unausgewerteter) R-Ausdruck, so wird versucht, die mathematisch übliche Darstellung zu geben. R-Ausdrücke können mit den Funktionen `expression()` oder `bquote()` erzeugt werden.

Beispiel:

```
text(x, y, expression(paste(bquote("(" , atop(n, x), ")"),
. (p)^x, . (q)^{ n-x})))
```

Ausgabe-Beispiele erhält man mit `demo(plotmath)`.

A.21.4. Grafik-Parameter und Layout.

Parametrisierungen	
<code>par()</code>	Setzt Parameter des Basis-Grafiksystems <i>Aufruf:</i> siehe <code>help(par)</code> <i>Beispiel:</i> <code>par(mfrow = c(m, n))</code> unterteilt den Grafikbereich in m Zeilen und n Spalten, die Zeile für Zeile gefüllt werden. <code>par(mfcol = c(m, n))</code> füllt den Bereich Spalte für Spalte.
<code>split.screen()</code>	Teilt den Grafik-Bereich in Teile <i>Aufruf:</i> <code>split.screen(figs, screen, erase = TRUE)</code> . Hat <code>figs</code> zwei Einträge, so werden damit die Anzahl der Zeilen und Spalten festgelegt. Ist <code>figs</code> eine Matrix, so gibt jede Zeile die Koordinaten eines Grafikbereichs in relativen Koordinaten $[0 \dots 1]$ an. <code>split.screen()</code> kann auch geschachtelt werden.
<code>screen()</code>	Wählt Grafik-Bereich für die nächste Ausgabe. <i>Aufruf:</i> <code>screen(n = cur.screen , new = TRUE)</code> .
<code>layout()</code>	Unterteilt den Grafik-Bereich. Diese Funktion ist mit anderen Layout-Funktionen nicht verträglich.

A.22. Einfache Statistische Funktionen

<i>Statistik-Funktionen</i>	
<i>sum()</i>	summiert Komponenten eines Vektors
<i>cumsum()</i>	bildet kumulierte Summen
<i>prod()</i>	multipliziert Komponenten eines Vektors
<i>cumprod()</i>	bildet kumulierte Produkte
<i>length()</i>	Länge eines Objekts, z.B. Vektors
<i>max()</i> <i>min()</i>	Maximum, Minimum. Siehe auch <i>pmax</i> , <i>pmin</i>
<i>range()</i>	Minimum und Maximum
<i>cummax()</i> <i>cummin()</i>	Kumulatives Maximum, Minimum
<i>quantile()</i>	Stichprobenquantile. Für theoretische Verteilungen: <i>qxxxx</i> , z.B. <i>qnorm</i>
<i>median()</i>	Median
<i>mean()</i>	Mittelwert auch getrimmte Mittel
<i>var()</i>	Varianz, Varianz / Kovarianzmatrix
<i>sort()</i> <i>rev()</i>	Sortierung
<i>order()</i>	Sortierung nach Leit-Element, auch für mehrere Variable
<i>rev()</i>	Umgekehrte Sortierung
<i>rank()</i>	Stichprobenränge

A.23. Verteilungen, Zufallszahlen, Dichten...

Der Basis-Generator für uniforme Zufallszahlen wird von *Random* verwaltet. Verschiedene mögliche Basis-Generatoren stehen zur Verfügung. **Für ernsthafte Simulation wird eine Lektüre der Empfehlungen von Marsaglia et al. dringend empfohlen.** (Siehe `help(.Random.seed)`). Alle nicht-uniformen Zufallszahlengeneratoren sind vom aktuellen Basisgenerator abgeleitet. Eine Übersicht über die wichtigsten nicht-uniformen Zufallszahlengeneratoren, ihre Verteilungsfunktionen und ihre Quantile findet sich am Ende dieses Abschnitts.

R Zufallszahlen	
<code>.Random.seed</code>	<code>.Random.seed</code> ist eine globale Variable, die den augenblicklichen Zustand des Zufallszahlengenerators speichert. Diese Variable kann gesichert und mit <code>set.seed()</code> wieder restauriert werden.
<code>set.seed()</code>	initialisiert den Zufallszahlengenerator <i>Aufruf:</i> <code>set.seed(seed, kind = NULL)</code>
<code>RngKind()</code>	<code>RngKind()</code> gibt den Namen des aktuellen Basisgenerators. <code>RngKind (<name>)</code> setzt einen Basisgenerator. <i>Aufruf:</i> <code>RngKind()</code> <code>RngKind(<name>)</code> <i>Beispiel:</i> <code>RngKind("Wichmann-Hill")</code> <code>RngKind("Marsaglia-Multicarry")</code> <code>RngKind("Super-Duper")</code>
<code>sample()</code>	<code>sample()</code> zieht eine Zufallsstichprobe aus den im Vektor <i>x</i> angegebenen Werten, mit oder ohne Zurücklegen (je nach Wert von <code>replace</code>). Size ist defaultmäßig die Länge von <i>x</i> . Optional kann <code>prob</code> ein Vektor von Wahrscheinlichkeiten für die Werte von <i>x</i> sein. <i>Aufruf:</i> <code>sample(x, size, replace = FALSE, prob)</code> <i>Beispiel:</i> Zufällige Permutation: <code>sample(x)</code> <code>val<-c("H", T)</code> <code>prob<-c(0.3, 0.7)</code> <code>sample(val, 10,</code> <code>replace = T, prob)</code>

Sollen Simulationen reproduzierbar sein, so muss der Zufallszahlengenerator in einen kontrollierten Zustand gesetzt sein. Ein Beispiel dafür ist die folgende Anweisungsfolge:

```
save.seed <- .Random.seed
save.kind <- RNGkind()
```

Mit `set.seed(save.seed, save.kind)` wird dann der Zustand des Genarators bei Bedarf restauriert.

Die einzelnen Funktionsnamen für die wichtigsten nicht-uniformen Generatoren und Funktionen setzen sich aus einem Präfix und dem Kurznamen zusammen. Allgemeiner Schlüssel: *xxxx* ist der Kurzname

rxxxx erzeugt Zufallszahlen

dxxxx Dichte oder Wahrscheinlichkeit

pxxxx Verteilungsfunktion

qxxxx Quantile

Beispiel:

`x<-runif(100)` erzeugt 100 U(0,1)-verteilte Zufallsvariable

`qf(0.95, 10, 2)` berechnet das 95%-Quantil der F(10,2)-Verteilung.

<i>Verteilungen</i>	<i>Kurzname</i>	<i>Parameter und Default-Werte</i>
Beta	<i>beta</i>	<i>shape1, shape2, ncp = 0</i>
Binomial	<i>binom</i>	<i>size, prob</i>
Cauchy	<i>cauchy</i>	<i>location = 0, scale = 1</i>
χ^2	<i>chisq</i>	<i>df, ncp = 0</i>
Exponential	<i>exp</i>	<i>rate = 1</i>
F	<i>f</i>	<i>df1, df2 (ncp = 0)</i>
Gamma	<i>gamma</i>	<i>shape, scale = 1</i>
Gauß	<i>norm</i>	<i>mean = 0, sd = 1</i>
Geometrisch	<i>geom</i>	<i>prob</i>
Hypergeometrisch	<i>hyper</i>	<i>m, n, k</i>
Lognormal	<i>lnorm</i>	<i>meanlog = 0, sdlog = 1</i>
Logistisch	<i>logis</i>	<i>location = 0, scale = 1</i>
Negativ-Binomial	<i>nbinom</i>	<i>size, prob</i>
Poisson	<i>pois</i>	<i>lambda</i>
Student's t	<i>t</i>	<i>df</i>
Tukey Studentised Range	<i>tukey</i>	
Uniform	<i>unif</i>	<i>min = 0, max = 1</i>
Wilcoxon Signed Rank	<i>signrank</i>	<i>n</i>
Wilcoxon Rank Sum	<i>wilcox</i>	<i>m, n</i>
Weibull	<i>weibull</i>	<i>shape, scale = 1</i>

A.24. Verarbeitung von Ausdrücken

Die Sprachausdrücke von R sind genau so Objekte wie Daten oder Funktionen. Wie diese können sie gelesen oder verändert werden.

<i>Umwandlungen</i>	
<code>parse()</code>	Wandelt Eingabe in eine Liste von R-Ausdrücken um. <code>parse</code> führt den Parse-Schritt durch, wertet die Ausdrücke aber nicht aus.
<code>deparse()</code>	Wandelt einen R-Ausdruck in interner Darstellung in eine Zeichen-darstellung um.
<code>expression()</code>	erzeugt einen R-Ausdruck in interner Darstellung. <i>Beispiel:</i> <code>integrate <- expression(integral(fun, lims))</code> <i>Siehe auch</i> 1.3.1: Mathematischer Formelsatz in Plot-Beschriftungen
<code>substitute()</code>	R-Ausdrücke mit Auswertung aller definierten Terme.
<code>bquote()</code>	R-Ausdrücke mit selektiver Auswertung. Terme in <code>.()</code> werden ausgewertet. <i>Beispiele:</i> <code>n<-10; bquote(n^2 == .(n*n))</code>

<i>Auswertung</i>	
<code>eval()</code>	wertet einen Ausdruck aus.

Literaturverzeichnis

- [BCW88] BECKER, Richard A. ; CHAMBERS, John M. ; WILKS, Allan R.: *The New S Language*. London : Chapman & Hall, 1988
- [CH92] CHAMBERS, John M. ; HASTIE, Trevor J.: *Statistical Models in S*. London : Chapman & Hall, 1992
- [Cha98] CHAMBERS, John M.: *Programming with Data*. New York : Springer, 1998. – ISBN 0-387-98503-4
- [Cle93] CLEVELAND, William S.: *Visualizing Data*. Murray Hill : AT&T Bell Laboratories, 1993
- [FB94] FURNAS, George W. ; BUJA, Andreas: Prosection views: dimensional inference through sections and projections. In: *J. Comput. Graph. Statist.* 3 (1994), Nr. 4, S. 323–385. – ISSN 1061–8600
- [GI00] GENTLEMAN, Robert ; IHAKA, Ross: Lexical Scope and Statistical Computing. In: *Journal of Computational and Graphical Statistics* 9 (2000), S. 491–508
- [GP04] GROTHENDIECK, Gabor ; PETZOLDT, Thomas: R Help Desk: Date and Time Classes in R. In: *R News* 4 (2004), June, Nr. 1, S. 29–32
- [GS77] GÄNSSLER, Peter ; STUTE, Winfried: *Wahrscheinlichkeitstheorie*. Springer, 1977
- [ICR87] INSELBERG, Alfred ; CHOMUT, Tuval ; REIF, Mordechai: Convexity algorithms in parallel coordinates. In: *J. Assoc. Comput. Mach.* 34 (1987), Nr. 4, S. 765–801. – ISSN 0004–5411
- [JK70] JOHNSON, N.L. ; KOTZ, S.: *Discrete Distributions*. New York : Wiley, 1970. – ISBN 0–471–44360–3
- [Jør93] JØRGENSEN, Bent: *The Theory of Linear Models*. New York-London : Chapman & Hall, 1993. – ISBN 0–412–04261–1
- [Mil81] MILLER, R. G.: *Simultaneous Statistical Inference*. New York : Springer, 1981. – ISBN 0–387–90584–0
- [Mur06] MURRELL, Paul: *R Graphics*. Boca Raton, Fla. [u.a.] : Chapman & Hall/CRC, 2006. – XIX, 301 S. S. – ISBN 1–58488–486–X
- [R D07a] R DEVELOPMENT CORE TEAM: An Introduction to R / R Project. 2007. – Forschungsbericht
- [R D07b] R DEVELOPMENT CORE TEAM: R Data Import/Export / R Project. 2007. – Forschungsbericht
- [R D07c] R DEVELOPMENT CORE TEAM: The R language definition / R Project. 2007. – Forschungsbericht
- [R D07d] R DEVELOPMENT CORE TEAM: The R Reference Index / R Project. 2007. – Forschungsbericht
- [R D08] R DEVELOPMENT CORE TEAM: Writing R Extensions / R Project. 2008. – Forschungsbericht
- [Rao73] RAO, C. R.: *Linear Statistical Inference and Its Applications*. 2. Wiley, 1973
- [RM79] REAVEN, G. M. ; MILLER, R. G.: An Attempt to Define the Nature of Chemical Diabetes Using a Multidimensional Analysis. In: *Diabetologia* 16 (1979), S. 17–24
- [Saw94a] SAWITZKI, Günther: Numerical Reliability of Data Analysis Systems. In: *Computational Statistics & Data Analysis* 18 (1994), Nr. 2, S. 269–286
- [Saw94b] SAWITZKI, Günther: Report on the Numerical Reliability of Data Analysis Systems. In: *Computational Statistics & Data Analysis* 18 (1994), Nr. 2, S. 289 – 301
- [Saw02] SAWITZKI, Günther: Quality Control and Early Diagnostics for cDNA Microarrays. In: *R News* 2 (2002), March, Nr. 1, S. 6–10
- [VR00] VENABLES, William N. ; RIPLEY, Brian D.: *S Programming*. Springer, 2000. – ISBN 0-387-98966-8

- [VR02] VENABLES, William N. ; RIPLEY, Brian D.: *Modern Applied Statistics with S*. 4. Heidelberg : Springer, 2002. – ISBN 0-387-95457-0

Index

- *Topic **aplot**
 - coplot, 4-15
- *Topic **debugging**
 - browser, A-21
 - debug, A-21
 - recover, A-21
 - traceback, A-21
- *Topic **distribution**
 - qqnorm, 3-6
 - Uniform, 1-4
- *Topic **hplot**
 - coplot, 4-15
 - pairs, 4-7
 - qqnorm, 3-6
- *Topic **htest**
 - t.test, 3-11
 - wilcox.test, 3-14
- *Topic **loess**
 - loess, 2-34
- *Topic **models**
 - anova, 2-20
- *Topic **regression**
 - anova, 2-20
 - lm, 2-11
- *Topic **smooth**
 - loess, 2-34
- PP*-Plot, **1-38**
- QQ*-Plot, **1-38**
- .Random.seed, 1-5

- abbreviate, A-15
- abline, 1-12
- add1, 2-21
- Added-Variable-Plots, **4-27**
- aggregate, A-14
- Annotation, A-38
- anova, 2-12, **2-20**, 2-29, A-36
- anova.lm, 2-13
- aov, 2-11–2-13, 2-21
- aperm, A-16
- apply, 1-23, A-14
- apropos, A-1, A-2
- args, A-1

- array, A-11, A-13
- as.data.frame, 2-11, 2-35
- attach, A-12
- attr, 2-39
- attributes, A-7
- axis, A-38

- Bandbreite, **1-10**
- barchart, 4-4
- barplot, 4-4
- bedingt, **4-15**
- Beschriftung, A-38
- Bindung, 3-13
- Bootstrap, **3-9**
- Box-Cox-Transformation, **2-33**
- boxcox, 2-33
- boxplot, 1-35, 4-4, A-37
- bquote, 1-24, A-39, A-45
- browser, A-21
- brushing, **4-5**
- bwplot, 4-4
- by, A-14

- c, 1-8
- casefold, A-15
- cbind, A-11, A-13
- chartr, A-15
- chisq.test, 1-28
- citation, 1-49
- class, 2-12, 2-39, A-7
- cloud, 4-4, 4-11
- co.intervals (*coplot*), 4-15
- coef, 2-13, 2-40
- coefficients, 2-21
- confint, 2-13, 2-41
- contour, 4-1, 4-2, A-37
- contourplot, 4-4
- contrasts, A-36
- Coplot, **4-15**
- coplot, **4-15**, A-37
- cummax, A-41
- cummin, A-41
- cumprod, A-41
- cumsum, A-41

- curse of dimension, **4-32**
- curve, *1-21, A-37*
- data, *1-48, 2-38, A-9, A-29, A-31*
- data.entry, *A-29*
- data.frame, *A-12*
- data.matrix, *4-8*
- Datenstrukturen, *1-17, A-11*
- DateTimeClasses, *A-30*
- Datum
 - see DateTimeClasses, *A-30*
- debug, *A-21*
- Debugging, *1-46, A-21*
- demo, *A-1*
- density, *1-15*
- densityplot, *4-4*
- deparse, *A-45*
- Design-Matrix, **2-2, 2-14**
- detach, *A-12, A-31*
- determinant, *A-33*
- dim, *A-11, A-13*
- dimnames, *A-13*
- dir, *A-9*
- dotchartt, *4-4*
- dotplot, *4-4, A-37*
- drop1, *2-21*
- dump, *A-27*
- dunif (*Uniform*), *1-4*
- duplicated, *A-15*
- edit, *A-29*
- effects, *2-13, 2-21, 2-40*
- eigen, *A-33*
- environment, *1-47, A-2*
- Erwartungswert, **1-33**
- eval, *1-46, A-45*
- exakter Test, *3-13*
- example, *A-1*
- expand.grid, *A-15*
- expression, *1-13, A-39, A-45*
- factor, *2-4, 4-16, A-11*
- Faktor, **2-4**
 - Stufen, *2-5*
- find, *A-2*
- Fit, **2-6**
- fitted, *2-13, 2-41*
- fitted.values, *2-21*
- formula, *2-12, A-36*
- function, **1-43**, *A-19–A-23*
- function, *4-16*
- Funktion
 - polymorph, *siehe polymorph*
- Güte, **1-31**
- Gauß-Markov-Schätzer, **2-5**
- getwd, *A-9*
- glht, *2-29*
- glm, *2-13, A-36*
- help, *A-1, A-9*
- help.search, *A-1*
- hist, *1-15, 1-18, 4-4, A-37*
- histogram, *4-4*
- Hut-Matrix, **2-7**
- identify, *4-35*
- image, *4-1, 4-2, 4-4, 4-44, A-37*
- influence, *2-41*
- inherits, *2-39*
- install.packages, *1-47, A-31*
- integrate, *4-31*
- Kern, **1-10**
- Kleinste-Quadrate-Schätzer, **2-5**
- Kontrast, **2-21, 2-27**
- kruskal.test, *3-16*
- ks.test, *1-27*
- lapply, *A-14*
- lattice.options, *A-9*
- lda, *4-24*
- leaps, *4-38*
- legend, *1-42, A-38*
- Legende, *A-38*
- length, *1-12, A-2, A-7, A-41*
- levels, *A-11*
- library, *A-9, A-31*
- lines, *A-38*
- linking, **4-5**
- list, *A-12*
- lm, *2-5, 2-11, 2-21, 2-40, 4-36, 4-43, A-36*
- lm.fit, *2-12, 2-13*
- lm.influence, *2-13*
- lm.wfit, *2-13*
- load, *2-37, 2-38, A-29*
- locator, *A-38*
- locfit, *4-6*
- loess, **2-34**
- loess.control, *2-35, 2-36*
- loglin, *1-29*
- lowess, *2-36*
- ls, *1-47, A-2, A-9, A-25*
- ls.str, *A-2*
- Marginalverteilung, **4-7**
- MASS, *2-41*
- match, *A-15*
- matrix, *4-17, A-13*
- max, *A-41*
- mean, *1-33, A-41*
- median, *A-41*
- merge, *A-15*

- methods, 2-40, A-9
 min, A-41
 missing, A-19
 mode, 2-39, A-2, A-3, A-7
 model.frame, 3-11, 3-14
 model.matrix, 2-12, 2-16, 2-41, A-36
 model.matrix.default, 2-12
 model.offset, 2-12
 Modell
 einfaches lineares, **2-8**
 lineares, 2-2
 Modellfunktion, **2-1**
 mtext, 1-44, 4-17, A-39
 mva, 4-27
 mvr, 4-43

 NA, 3-6
 na.exclude, 2-11
 na.fail, 2-11
 na.omit, 2-11
 names, A-7
 ncol, A-13
 nlm, A-36
 nls, A-36
 nrow, A-13

 objects, A-1, A-25
 offset, 2-12
 options, 2-11, A-9
 order, A-41
 ordered, A-11
 outer, 1-23, A-14

 package.manager, A-31
 package.skeleton, 1-48, 1-49, A-31
 pairs, 4-4, **4-7**, 4-7, 4-10, 4-17, 4-33, A-37
 panel.smooth, 4-17
 par, 4-17, A-9
 parallel, 4-4
 Parameter
 default, 1-3
 parse, 1-46, A-45
 persp, 4-1, 4-2, 4-4, A-37
 plot, 1-5
 plot, 1-12, 1-13, 2-30, 2-40, 4-4, 4-22, 4-35,
 A-5, A-37
 plotmath, 1-13
 pmatch, A-15
 points, 4-17, A-38
 polygon, A-38
 polymorph, 1-3, 1-47, 2-39, 2-40, A-5
 power.prop.test, 3-23
 power.t.test, 3-20
 ppoints, 3-7
 prcomp, 4-27, 4-40
 predict, 2-13, 2-25, 2-41

 predict.lm, 2-13, 2-25, 4-43
 predict.loess, 2-36
 print, 1-46, 1-47, 2-40, 4-2, 4-11, A-5
 print.anova (*anova*), 2-20
 print.lm (*lm*), 2-11
 probability plot, **1-38**
 prod, A-41
 Profiling, A-21
 projection pursuit, **4-12**
 prop.test, 3-12, 3-22
 psignrank, 3-16
 punif (*Uniform*), 1-4
 pwilcox, 3-16

 q, 1-3
 qq, 4-4
 qqline (*qqnorm*), 3-6, A-37
 qqmath, 4-4
 qqnorm, 3-5, **3-6**, 4-4, A-37
 qqplot, 3-5, 4-4, A-37
 qqplot (*qqnorm*), 3-6
 qr, A-33
 quantile, 1-34, A-41
 Quantilplot, **1-38**
 qunif (*Uniform*), 1-4

 range, 4-17, A-41
 rank, A-41
 rbind, A-11, A-13
 read.fwf, A-30
 read.table, A-30
 recover, A-21
 Regression
 lineare, 2-2
 Regressor, **2-1**
 regsubsets, 4-38
 rep, 1-8
 replicate, A-14
 require, A-31
 reshape, A-15
 residuals, 2-13, 2-21, 2-41
 Residuum, **2-8**
 Respons, **2-1**
 rev, A-41
 rm, A-25
 RngKind, A-43
 rnorm, 1-5
 Rprof, A-22
 Rprofmem, A-22
 rug, 1-15
 runif, 1-4
 runif (*Uniform*), 1-4

 sample, A-43
 sapply, A-14
 save, 1-48, 2-37, A-29

- scan, *A-30*
- Scatterplot-Matrix, **4-7**
- screen, *A-39*
- sd, *1-34*
- search, *1-47, A-1, A-9*
- searchpaths, *A-1*
- seq, *1-8, A-15*
- Serienplot, **1-5**
- set.seed, *A-43*
- Shift-Familie, **3-3**
- sink, *A-27*
- Skala
 - kategorial, *2-4*
 - ordinal, *2-4*
- Skalen-Shiftfamilie, **3-4**
- smoothing, **1-10**
- solve, *A-33*
- sort, *1-12, A-41*
- source, *1-47, 1-48, A-27*
- split, *A-14, A-16*
- split.screen, *A-39*
- splom, *4-4*
- stack, *2-38*
- Standardabweichung, **1-33**
- stdres, *2-41*
- Stichproben
 - wiederholte, **1-30**
- Stichprobenvarianz, **1-33**
- stochastisch kleiner, **3-3**
- storage.mode, *2-39, A-7*
- str, *A-7*
- Streuungszerlegung, **2-19**
- stripchart, *4-4*
- stripplot, *4-4, 4-23*
- structure, *A-5, A-7*
- studres, *2-41*
- substitute, *A-45*
- substring, *A-15*
- sum, *A-41*
- summary, *1-34, 2-21, A-5*
- summary.lm, *2-13*
- summaryRprof, *A-22*
- svd, *A-33*
- Sweave, *1-47*
- sys.parent, *A-2*
- system, *A-9*
- system.time, *A-22*

- t, *A-15*
- t.test, **3-11, 3-16**
- table, *1-18, A-14, A-15*
- tapply, *A-11, A-14*
- terms, *2-12, A-36*
- Test
 - χ^2 , *1-28*
 - exakt, *3-13*
 - Kolmogorov-Smirnov, *1-27*
 - Median-, *1-27*
 - Monte-Carlo, **1-23**
 - t, *3-11*
 - Wilcoxon, *3-13*
- title, *4-17, A-38*
- tolower, *A-15*
- toupper, *A-15*
- trace, *A-21*
- traceback, *A-22*
- trellis.par.set, *A-9*
- try, *A-22*
- ts.intersect, *2-13*
- typedef, *A-3*
- typeof, *2-39, A-2, A-7*

- unclass, *2-39*
- undebug, *A-21*
- Uniform, **1-4**
- unique, *A-15*
- unsplit, *A-16*
- untrace, *A-22*
- update, *A-36*
- update.packages, *1-47*
- UseMethod, *2-40*

- var, *1-33, A-41*
- Varianz, **1-33**
 - residuelle, **2-8**
- Varianzanalyse, **2-19**
- Variationskoeffizient, **3-23**
- vcov, *2-13, 2-41*

- wilcox.exact, *3-16*
- wilcox.test, *3-13, 3-14*
- wilcox_test, *3-13*
- Wilkinson-Rogers-Notation, **2-3**
- wireframe, *4-4*
- write, *A-27*
- write.table, *A-30*

- xyplot, *4-4*

- Zeit
 - see DateTimeClasses, *A-30*
- Zufallszahlen, *1-4*
 - Pseudo-, *1-7*
 - reproduzierbare, *A-43*