

Introduction to R

Marlene Müller

Slides version: April 13, 2012

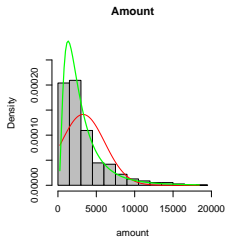
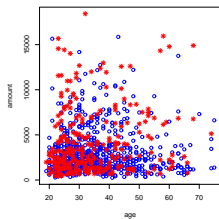


Table of contents

What is this R?

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics

“Advanced” mathematics

Basics in programming

References

What is this R?

What is this R?
How do I start?
Working with R under Unix/Linux
Working with R under Windows

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics


"Advanced" mathematics

Basics in programming

References

What is this R?

Programming Language S = developed at Bell Labs for statistics, simulation, graphics (Becker and Chambers, 1984)

- S-PLUS: commercial implementation
- R: implementation under GPL (GNU General Public License), open source
 - + interpreted program code, object orientation
 - + easily extensible by self-written routines, packages, DLLs
 - + many types of graphics (mainly static)
 - + standardized, simple-to-used data format ( `data.frame`)
 - + well developed format for fitting (regression) models
 - + active developers team, helpful mailing list
 - + increasing number of books on R on the market
 - (up to now) no "standard" GUI
 - available routines/packages sometimes difficult to find

How do I start?

What is this R?
How do I start?
Working with R under Unix/Linux
Working with R under Windows

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics

"Advanced" mathematics

Basics in programming

References

How do I start?

R is command-line oriented, so start simply by typing expressions like:

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

```
> 1+2*3^4  
[1] 163
```

```
> x <- 1; y <- 2  
> x+y  
[1] 3
```

in the following all R Code is printed in a way that can be directly copied:

```
x <- seq(-pi,pi,by=0.1)  
plot(x,sin(x),type="l",col="red",main="Sinuskurve")
```

Working with R under Unix/Linux

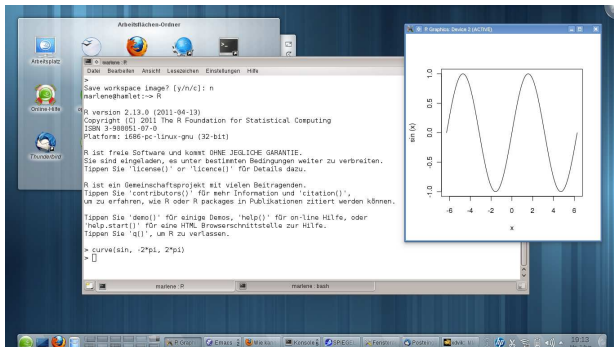


Figure: R in a Unix/Linux shell

R in a Windows desktop

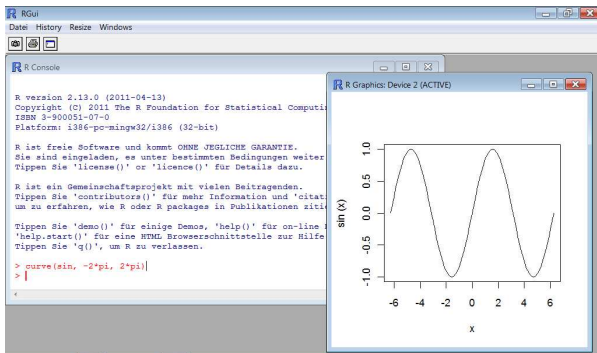


Figure: R in a Windows desktop

How to get help?

What is this R?

How to get help?

WWW

Mailing lists

Book selection (German & English)

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics

"Advanced" mathematics

Basics in programming

References

How to get help?

local help pages:

- ▶ help for a function:

```
help(<function name>) or ?<function name>
```

- ▶ help for a package:

```
library(help=<Package>)
```

usually, the texts in the local help pages correspond to those in the package documentation

WWW

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

R home page, there are in particular FAQs as well as a Google site search, and additionally:

- ▶ manuals (<http://cran.r-project.org/manuals.html>)
introduction, language definition, “Writing R Extensions” (DLLs, packages), introduction written in different languages (German, French, etc.)
- ▶ CRAN (<http://cran.r-project.org>)
Comprehensive R Archive Network (→ R software for download)
- ▶ mailing lists (<http://www.r-project.org/mail.html>)
- ▶ books list (<http://www.r-project.org/doc/bib/R-books.html>)
- ▶ conference announcements, related projects, . . .

Mailing lists

- ▶ R-help
main list for R user questions, take care to read <http://www.r-project.org/posting-guide.html> before!
→ also available as a (usenet-) news group gmane.comp.lang.r.general auf <http://news.gmane.org>
- ▶ R-announce, R-packages, R-devel
announcements, package announcements, developers list (→ more for R specialists)
- ▶ R-sig-* (special interests groups)
e.g. R-sig-finance = Special Interest Group for 'R in Finance'

for subscribing and archives see <http://www.r-project.org/mail.html> or <http://news.gmane.org/index.php?prefix=gmane.comp.lang.r>

helpful for search is <http://www.rseek.org>

Book selection (German & English)

data analysis:

- ▶ Wollschläger (2010): Grundlagen der Datenanalyse mit R
- ▶ Dalgaard (2002): Introductory Statistics with R
- ▶ Murrell (2005): R Graphics
- ▶ Venables and Ripley (2002): Modern Applied Statistics with S (R complements: <http://www.stats.ox.ac.uk/pub/MASS4>)

programming:

- ▶ Ligges (2009): Programmieren mit R (see also: <http://www.statistik.uni-dortmund.de/~ligges/PmitR/>)
- ▶ Gentleman (2008): R Programming for Bioinformatics
- ▶ Venables and Ripley (2000): S Programming (see also: <http://www.stats.ox.ac.uk/pub/MASS3/Sprog>)

more books:

→ <http://www.r-project.org/doc/bib/R-books.html>

Some calculations to start with

What is this R?

How to get help?

Some calculations to start with

- Data types

- Vectors, matrices, arrays, ...

- Operations (elementwise and/or vector-/matrixwise)

Data & files

Wonderful world of graphics

Some statistics

"Advanced" mathematics

Basics in programming

References

Some calculations to start with

```
demo()
demo(graphics)    ## nice graphics ;- )
demo(persp)       ## nice 3D graphics ;- )
demo(image)       ## more nice graphics ;- )

x <- 1
x <- 0 -> y
x <- y <- z <- NA    ## missing
x <- 0/0             ## not a number (NaN)
x <- NULL           ## no value

x <- rnorm(100)     ## vector of 100 N(0,1) random variables

hist(x, col="orange")          ## histogram
r <- hist(x, col="orange", freq=FALSE) ## same histogram?
g <- seq(-5,5,length=100)
ylim <- range(c(r$density,max(dnorm(g))))

hist(x, col="orange", freq=FALSE, ylim=ylim) ## same histogram?
lines(g, dnorm(g))                ## with N(0,1) pdf
```

Useful tools

```
ls()                ## list all R objects

x <- 1:3
x                  ## show object (vector: x)

print(x)           ## show object (vector: x), also within
                  ## R scripts and functions

fun <- function(x){ sin(x) }
fun                ## show object (function: fun)

median             ## show object (internal function: median)

rm(x)              ## delete object x

save.image()       ## save workspace (.RData, .Rhistory)
load(".RData")     ## load workspace (.RData, .Rhistory)

date()             ## date and time

q()                ## quit R
```


Data types

numeric:

```
x <- 1
y <- pi      ## predefined pi = 3.1415926535898
```

character:

```
x <- "a"
y <- "my text"
```

logical:

```
x <- TRUE
y <- 1 > 2
```

```
> y
[1] FALSE
```

more complex data types can be constructed by combining these three simple types into vectors, matrices, arrays and lists

Vectors

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

```
x <- 1:3
```

```
y <- c(1,1,1)
```

```
y <- rep(2,10)
```

```
z <- as.character(1:3)
```

```
z <- c("a", "b", "c")
```

```
length(z)
```

```
names(x) <- z
```

```
x[2:3]
```

```
x["b"]
```

all elements of a vector are of the same type (numeric, character, logical)!

Matrices

```
x <- 1:20  
x <- matrix(x, 5,4)      ## matrix(x, nrow=5,ncol=4)
```

```
x[2,3]  
x[c(1,5),2:4]  
x[,2:4]
```

```
dim(x)  
nrow(x); ncol(x)
```

```
length(x)  
as.vector(x)
```

```
dimnames(x) <- list(paste("row",1:nrow(x), sep=""),  
                    c("a","b","c","d"))
```

```
x[,"b"]  
x[,c("a","b")]
```

all elements of a matrix are of the same type (numeric, character, logical)!

Generating Vectors and Matrices

vectors from vectors:

```
x <- c(2,6,3)
```

```
y <- 1:3
```

```
c(x,y)          ## concatenate two vectors
```

```
c(x,1:5,y,6)    ## concatenate vectors and scalars
```

matrices from vectors:

```
x <- c(2,6,3)
```

```
y <- 1:3
```

```
cbind(x,y)      ## vertical concatenation
```

```
rbind(x,y)      ## horizontal concatenation
```

```
cbind(x,y,rep(0,3)) ## vertical concatenation
```

Arrays

```
x <- 1:60  
x <- array(x, c(5,4,3))
```

```
x[2,3,1]  
x[1,2:4,3]  
x[, ,1]
```

```
dim(x)  
nrow(x)  
ncol(x)
```

```
length(x)  
as.vector(x)
```

```
dimnames(x) <- list(paste("row",1:nrow(x), sep=""),  
                   c("a","b","c","d"),c("x","y","z"))
```

all elements of an array are of the same type (numeric, character, logical)!

Lists

```
x <- list(One=11:15, Two=c("a", "b", "c"), Three=(1:4)>0)
y <- list(x=x, Four=1:3)
```

```
x$One
y$x$One
```

```
y$Four
y[[2]]
```

```
length(x)
length(y)
```

```
y$Five <- names(x)
```

lists may contain objects of different type, these objects can be called with `$<name>` by name or with `[[<number>]]` by their number

Data frames

```
x <- data.frame(N=11:14, C=c("a","b","c","d"), L=(1:4)>0)
```

```
dim(x)
```

```
nrow(x); ncol(x)
```

```
length(x)
```

```
as.vector(x)
```

```
names(x)
```

```
x[2,3]
```

```
x[,2:3]
```

```
x[,2]
```

```
x["C"]
```

```
x$C
```

data frames can be seen as lists, with all columns having the same length
→ Excel tables, saved as .csv, are typically read into R as a (data.frame)

Operations (elementwise and/or vector-/matrixwise)

```
x <- matrix( 1:20, 5, 4)    ## 5x4 matrix

x+1; x-1; x*1; x/1        ## elementwise operations
sin(x); exp(x)            ## elementwise function calls

y <- 1:5
x * y                      ## elementwise multiplication

z <- 1:4
x %*% z                    ## matrix multiplication

min(x)                     ## minimum of all elements of x
apply(x,1,min)             ## row minima
apply(x,2,min)             ## column minima

y <- c(TRUE, TRUE, FALSE, FALSE)
y & TRUE                   ## elementwise logical "AND"
y | FALSE                  ## elementwise logical "OR"
!y                          ## elementwise logical "NOT"

y && TRUE                  ## here only the first result
y || FALSE                 ## holds! ("AND" or "OR")
```


Data & files

What is this R?

How to get help?

Some calculations to start with

Data & files

- Reading and saving CSV files

- R script files

- Random numbers and probability distributions

Wonderful world of graphics

Some statistics

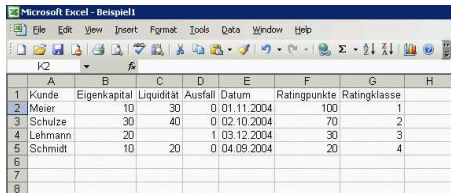
“Advanced” mathematics

Basics in programming

References

Data & files

example file in Excel:



The screenshot shows a Microsoft Excel spreadsheet with the following data:

	A	B	C	D	E	F	G	H
1	Kunde	Eigenkapital	Liquidität	Ausfall	Datum	Ratingpunkte	Ratingklasse	
2	Meier	10	30	0	01.11.2004	100	1	
3	Schulze	30	40	0	02.10.2004	70	2	
4	Lehmann	20		1	03.12.2004	30	3	
5	Schmidt	10	20	0	04.09.2004	20	4	
6								
7								
8								

→ save under Excel as CSV: Example1.csv

```
Kunde;Eigenkapital;Liquidität;Ausfall;Datum;Ratingpunkte;Ratingklasse  
Meier;10;30;0;01.11.2004;100;1  
Schulze;30;40;0;02.10.2004;70;2  
Lehmann;20;;1;03.12.2004;30;3  
Schmidt;10;20;0;04.09.2004;20;4
```

Reading and saving CSV files

reading the file Example1.csv:

```
x <- read.csv("Example1.csv", sep=";")
```

```
dim(x)
```

```
names(x)
```

```
x
```

result in R:

	Kunde	Eigenkapital	Liquidität	Ausfall	Datum	Ratingpunkte	Ratingklasse
1	Meier	10	30	0	01.11.2004	100	1
2	Schulze	30	40	0	02.10.2004	70	2
3	Lehmann	20	NA	1	03.12.2004	30	3
4	Schmidt	10	20	0	04.09.2004	20	4

saving the data to Example2.csv:

```
write.table(x, file="Example2.csv", sep=";", row.names=FALSE, quote=FALSE)
```

More functions for data input and output

functions for reading data:

- ▶ `read.table` (ASCII data)
- ▶ `scan` (scans any text file, further postprocessing necessary)

functions to convert data:

- ▶ `as.numeric`, `as.character`, `as.factor`, `as.Date`

other possibilities to communicate to Excel:

- ▶ `RODBC` (accessing data from databases)
- ▶ R-Excel-interface via DCOM server
(<http://cran.at.r-project.org/contrib/extra/dcom>)

R script files

run a script with R code:

```
> source("MyProgram.R")
```

saving R output to file:

```
sink("MyOutput.txt") ## from now all output goes to file  
sink()                ## and now to the screen again
```


Random numbers and probability distributions

examples for the normal distribution:


<code>rnorm(n, mean=0, sd=1)</code>	pseudo-random numbers
<code>dnorm(x, mean=0, sd=1)</code>	density (pdf)
<code>pnorm(x, mean=0, sd=1)</code>	cumulative distribution function (cdf)
<code>qnorm(p, mean=0, sd=1)</code>	quantiles

in the same manner:


uniform distribution

 `{r|d|p|q}unif`

t distribution

 `{r|d|p|q}t`


lognormal distribution

 `{r|d|p|q}lnorm`


gamma distribution

 `{r|d|p|q}gamma`


χ^2 distribution

 `{r|d|p|q}chisq`


beta distribution

 `{r|d|p|q}beta`

binomial distribution

 `{r|d|p|q}binom`


poisson distribution

 `{r|d|p|q}pois`


exponential distribution

 `{r|d|p|q}exp`

F distribution

 `{r|d|p|q}f`

...

→ it is possible to fix the seed by  `set.seed`

Example: Normal- vs. t distribution

```
x <- rnorm(100)
mean(x)
sd(x)
```

```
plot(rnorm(10000), rnorm(10000))
```

```
x <- seq(-5,5,by=0.1)
plot(x, dnorm(x), type="l", col="black", lwd=2)
lines(x, dt(x, df=1), col="blue")
lines(x, dt(x, df=5), col="orange")
lines(x, dt(x, df=20), col="red")
```

```
qnorm(0.95)
qnorm(0.975)
```

Multivariate normal distribution

```
library(help=mvtnorm)
library(mvtnorm)

mu <- c(0,0)      ## means
sigma <- c(1,1)   ## standard deviations
rho <- 0.5        ## correlation

S <- matrix(NA, 2,2)
diag(S) <- sigma^2
S[1,2] <- S[2,1] <- rho*prod(sigma)

x <- rmvnorm(n=10000, mean=mu, sigma=S)
plot(x)

x <- seq(-5*sigma[1]+mu[1], 5*sigma[1]+mu[1], length = 50)
y <- seq(-5*sigma[2]+mu[2], 5*sigma[2]+mu[2], length = 50)
f <- function(x,y) { dmvnorm(cbind(x,y), mean=mu, sigma=S) }
z <- outer(x, y, f)
persp(x, y, z, theta = 10, phi = 20, expand = 0.5,
      col = "lightblue", shade = 0.75)
```


Wonderful world of graphics

What is this R?

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

- Barplots

- Boxplots

- Histograms

- Scatterplots and curves

- Three-dimensional graphics

- Arranging plots

- Save plots to files

Some statistics

"Advanced" mathematics

Basics in programming

References

Wonderful world of graphics

credit scoring data:

http://www.stat.uni-muenchen.de/service/datenarchiv/kredit/kredit_e.html

```
file <- read.csv("D:\\kredit.asc", sep=" ")
y <- 1-file$ kredit          ## default set to 1

prev <- (file$moral >2)+0      ## previous loans were OK
employ <- (file$beszeit >1)+0  ## employed (>=1 year)
dura <- (file$laufzeit)       ## duration
d9.12 <- ((file$laufzeit >9)&(file$laufzeit <=12))+0 ## 9-12 months
d12.18 <- ((file$laufzeit >12)&(file$laufzeit <=18))+0 ## 12-18 months
d18.24 <- ((file$laufzeit >18)&(file$laufzeit <=24))+0 ## 18-24 months
d24 <- (file$laufzeit >24)+0   ## > 24 months
amount <- file$hoehe          ## amount of loan
age <- file$alter             ## age of applicant
savings <- (file$sparkont > 4)+0 ## savings >= 1000 DM
phone <- (file$telef==1)+0    ## applicant has telephone
foreign <- (file$gastarb==1)+0 ## non-german citizen
purpose <- ((file$verw==1)|(file$verw==2))+0 ## loan is for a car
house <- (file$verm==4)+0     ## house owner
```

Barplots

→ graphical representation of the frequency distribution for discrete variables

```
table(dura)          ## frequency table

barplot(table(dura), col="cyan", main="Duration of Loan")
                    ## absolute frequencies

n <- length(dura)
barplot(table(dura)/n, col="cyan", main="Duration of Loan")
                    ## relative frequencies

par(mfrow=c(1,3))   ## graphical display with 1 row, 3 columns

barplot(table(dura),    col="cyan",    main="Duration of Loan")
barplot(table(savings), col="orange",  main="Savings >1000 DM")
barplot(table(house),   col="magenta",  main="House Owner")

par(mfrow=c(1,1))   ## reset display to single plot
```

Example: Barcharts

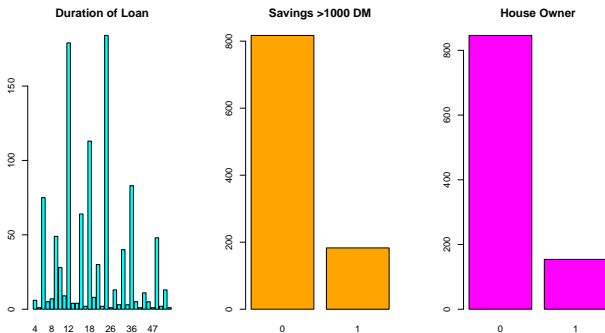


Figure: Examples for bar plots: duration of loan (left), savings (center) and house-owner indicator (right)

Boxplots

→ graphical representation of outliers, minima/maxima, 25%-, 50%-, and 75%-quantiles

```
boxplot(age)
```

```
boxplot(age, horizontal=TRUE)
```

```
boxplot(age, col="gray",horizontal=TRUE)
```

```
boxplot(age ~ y, col=c("gray","red"),  
        horizontal=TRUE, main="Age vs. Y")
```

```
boxplot(amount ~ y, col=c("gray","red"),  
        horizontal=TRUE, main="Amount vs. Y")
```

Examples: Boxplots

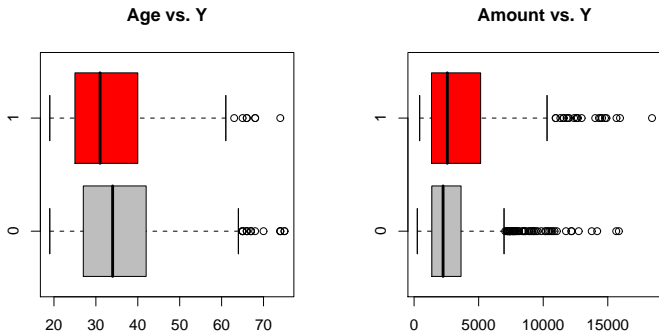


Figure: Age of credit applicant (left) and amount of loan (right) vs. default indicator (1 = default, 0 = non-default)

Histograms

→ graphical representation of the distribution (probability density function) of continuous variables

```
hist(age)
hist(age, freq=FALSE)
hist(age, freq=FALSE, col="gray")

hist(amount, freq=FALSE, col="gray", main="Amount")
xx <- seq(min(amount),max(amount), length=100)
lines(xx, dnorm(xx, mean(amount), sd(amount)), col="red")
lines(xx, dlnorm(xx, mean(log(amount)), sd(log(amount))), col="green", lwd=2)

## smaller intervals and better vertical scale
b <- seq(0,20000,by=1500)           ## new intervals
h <- hist(amount, freq=FALSE, breaks=b, plot=FALSE)
                                   ## histogram without display
xx <- seq(min(amount),max(amount), length=100)
d1 <- dnorm(xx, mean(amount), sd(amount))  ## normal pdf
d2 <- dlnorm(xx, mean(log(amount)), sd(log(amount)))
                                   ## lognormal pdf

ylim <- range( c(h$density, d1, d2) )

hist(amount, freq=FALSE, breaks=b, col="gray", main="Amount", ylim=ylim)
lines(xx, d1, col="red")
lines(xx, d2, col="green", lwd=2)
```

Example: Histogram

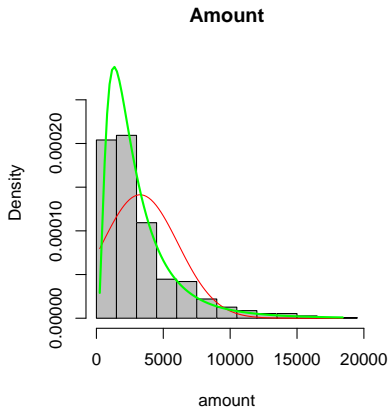


Figure: Probability distribution of the amount of the loan, histogram in comparison with normal and lognormal pdfs

Scatterplots

```
plot(age, amount)
```

```
color <- 1*(y==1) + 2*(y==0)
plot(age, amount, col=color)
```

```
color <- rep("", length(age))
color[y==1] <- "red"
color[y==0] <- "blue"
plot(age, amount, col=color)
```

```
plot(1:20,1:20,col=1:20, pch=1:20)
text(1:20,1:20,labels=as.character(1:20), pos=4)
```

```
symbol <- 8*(y==1) + 1*(y==0)
plot(age, amount, col=color, pch=symbol)
```

Scatterplots and curves

```
x <- seq(-pi,pi,length=100)
plot(x, sin(x), type="l")
lines(x, cos(x), col="red")

logit <- glm(y ~ age, family=binomial(link = "logit"))

plot(age, logit$fitted.values)

plot(age, logit$fitted.values, type="l")
      ## not this way ...

o <- order(age)
plot(age[o], logit$fitted.values[o], type="l")
      ## ... but that way! (sort data first)

plot(age[o], logit$fitted.values[o], type="l", lwd=2, ylim=c(0,1))
title("PDs")
points(age, y, col="red", pch=3, cex=0.5)
```

Example: scatterplots and curves

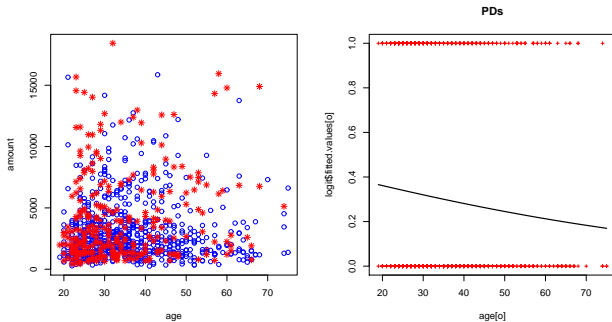


Figure: Scatterplot of age vs. amount (left), logit probabilities of default (right)

Surfaces, point clouds, contours

```
## bivariate normal pdf
library(mvtnorm)
x <- y <- seq(-5, 5, length = 50)
f <- function(x,y) { dmvnorm(cbind(x,y)) }
z <- outer(x, y, f)
persp(x, y, z, theta=10, phi=20, expand=0.5, col="lightblue")
persp(x, y, z, theta=10, phi=20, expand=0.5, col="lightblue",
      shade = 0.75)

## contours of the bivariate normal pdf
x <- y <- seq(-5, 5, length = 150)
z <- outer(x, y, f)
contour(x, y, z, nlevels=20)
contour(x, y, z, nlevels=20, col=rainbow(20))
contour(x, y, z, nlevels=20, col=rainbow(20), labels="")

## 3-dimensional normal data
library(scatterplot3d)
x <- matrix(rnorm(15000),ncol=3)
scatterplot3d(x)
scatterplot3d(x, angle=20)
```

Surfaces, point clouds, contours

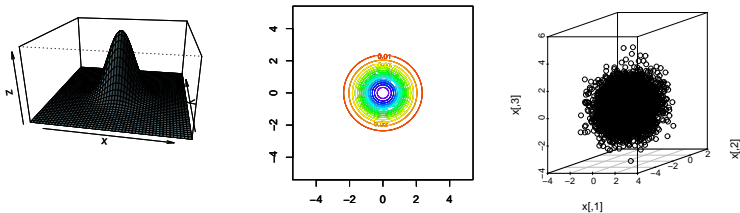


Figure: Bivariate normal pdf: 3D plot of the pdf (left), contour curves (center); 3D scatterplot (right)

Arranging plots

directly within the plot routines: → obtain help by `R?par`

- ▶ set colors with `R?col=...` (generate colors by → `R?rainbow`, `R?rgb`, `R?col2rgb`)
- ▶ set symbol style with `R?pch=...`, symbol size with `R?cex=...`
- ▶ set plot title with `R?main=...`, axes labels with `R?xlab=...`, `R?ylab=...`
- ▶ set plot drawing limits with `R?xlim=...`, `R?ylim=...`

after drawing a plot:

- ▶ add curves and points with `R?lines(...)` bzw. `R?points(...)`
- ▶ add labels (text) with `R?text(...)`
- ▶ add title with `R?title(...)`
- ▶ add legend with `R?legend(...)`

Save plots to files

- ▶ PostScript:

```
x <- matrix(rnorm(5000),ncol=2)
plot(x)
postscript("MyPlot.eps", width=5, height=4, horizontal=FALSE)
plot(x)
dev.off()
```

- ▶ other formats are e.g. `Rpdf`, `Rpictex`, `Rxfig`, `Rpng`, `Rjpeg`
→ see `R?Devices`

Some statistics

What is this R?

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics

Summary statistics

Tables

Linear Regression

Generalized linear model (GLM)

Other models for regression and times series analysis

Test for normality

Comparing distributions

Selected tests

"Advanced" mathematics

Basics in programming

References

Some statistics

credit scoring data:

http://www.stat.uni-muenchen.de/service/datenarchiv/kredit/kredit_e.html

```
file <- read.csv("D:\\\\kredit.asc", sep=" ")
y <- 1-file$kredit          ## default set to 1

prev <- (file$moral >2)+0      ## previous loans were OK
employ <- (file$beszeit >1)+0  ## employed (>=1 year)
dura <- (file$laufzeit)       ## duration
d9.12 <- ((file$laufzeit >9)&(file$laufzeit <=12))+0 ## 9-12 months
d12.18 <- ((file$laufzeit >12)&(file$laufzeit <=18))+0 ## 12-18 months
d18.24 <- ((file$laufzeit >18)&(file$laufzeit <=24))+0 ## 18-24 months
d24 <- (file$laufzeit >24)+0   ## > 24 months
amount <- file$hoehe          ## amount of loan
age <- file$alter             ## age of applicant
savings <- (file$sparkont > 4)+0 ## savings >= 1000 DM
phone <- (file$telef==1)+0    ## applicant has telephone
foreign <- (file$gastarb==1)+0 ## non-german citizen
purpose <- ((file$verw==1)|(file$verw==2))+0 ## loan is for a car
house <- (file$verm==4)+0     ## house owner
```

Summary statistics

```
kredit <- data.frame(y,age,amount,dura,prev,savings,house)
```

```
summary(kredit)
```

```
mean(kredit$age)
```

```
sd(kredit$age)
```

```
var(kredit$age)
```

```
cov(kredit[,1:3])
```

```
cor(kredit[,1:3])
```

```
median(kredit$age)
```

```
quantile(kredit$age,c(0.1,0.5,0.9))
```

```
library(help=e1071)
```

```
library(e1071)
```

```
skewness(kredit$age)
```

```
kurtosis(kredit$age)
```

```
skewness(rnorm(1000))
```

```
kurtosis(rnorm(1000))
```

Tables

```
length(kredit$age)
length(unique(kredit$age))
```

```
table(kredit$age)
table(kredit$dura)
table(kredit$savings)
```

```
table(kredit$y, kredit$savings)
table(kredit$y, kredit$savings)/nrow(kredit)
```

```
table(kredit$y, kredit$savings, kredit$house)
```

```
unique(kredit[,c("y", "savings", "house")])
```

Linear Regression

```
plot(kredit$age, kredit$dura)
```

```
lm <- lm( dura ~ age, data=kredit)
summary(lm)                ## dependence on age
abline(lm, col="red", lwd=2)
```

```
lm2 <- lm( dura ~ age + amount, data=kredit)
summary(lm2)               ## dependence on age+amount
```

```
lm3 <- lm( dura ~ amount, data=kredit)
summary(lm3)              ## dependence on amount
plot(kredit$amount, kredit$dura)
abline(lm3, col="red", lwd=2)
```

```
lm4 <- lm( dura ~ amount + I(amount^2), data=kredit)
summary(lm4)              ## dependence on amount (also squared)
o <- order(kredit$amount)
lines(kredit$amount[o], lm4$fitted.values[o], col="blue", lwd=2)
```

Linear Regression (cont'd)

→ duration of loan is clearly a function of amount:

```
> summary(lm4)
```

Call:

```
lm(formula = dura ~ amount + I(amount^2), data = kredit)
```

Residuals:

Min	1Q	Median	3Q	Max
-34.6115	-5.5761	-0.9547	5.0850	42.1110

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	8.410e+00	6.516e-01	12.906	< 2e-16 ***
amount	4.855e-03	2.961e-04	16.393	< 2e-16 ***
I(amount^2)	-1.815e-07	2.309e-08	-7.863	9.7e-15 ***

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

Residual standard error: 9.144 on 997 degrees of freedom

Multiple R-Squared: 0.4262, Adjusted R-squared: 0.425

F-statistic: 370.3 on 2 and 997 DF, p-value: < 2.2e-16

Linear Regression (cont'd)

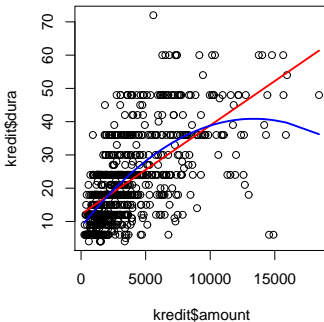
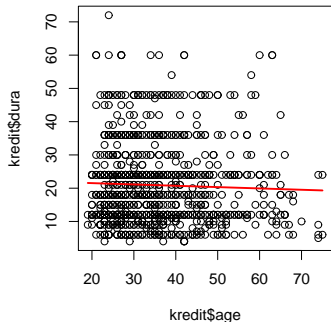


Figure: Dependence of duration on age (left) and amount (right)

Generalized linear model (GLM)

→ estimation of default probabilities with a logit model

```
logit <- glm(y ~ age + amount + dura + prev + savings + house,  
            family=binomial(link = "logit"))  
summary(logit)
```

```
logit2 <- glm(y ~ age + amount + I(amount^2) + dura + prev +  
             savings + house, family=binomial(link = "logit"))  
summary(logit2)
```

Examl: logit model

```
> summary(logit2)
```

Call:

```
glm(formula = y ~ age + amount + I(amount^2) + dura + prev +  
     savings + house, family = binomial(link = "logit"))
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.1244	-0.8495	-0.6196	1.0935	2.2584

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-4.637e-01	3.035e-01	-1.528	0.12652
age	-1.748e-02	7.159e-03	-2.442	0.01460 *
amount	-2.070e-04	9.348e-05	-2.214	0.02679 *
I(amount^2)	1.870e-08	6.941e-09	2.694	0.00707 **
dura	3.992e-02	8.106e-03	4.925	8.46e-07 ***
prev	-7.589e-01	1.619e-01	-4.688	2.76e-06 ***
savings	-9.897e-01	2.232e-01	-4.435	9.22e-06 ***
house	6.277e-01	2.073e-01	3.027	0.00247 **

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

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1221.7 on 999 degrees of freedom
Residual deviance: 1102.1 on 992 degrees of freedom
AIC: 1118.1

Number of Fisher Scoring iterations: 4

Other models for regression and times series analysis

model	routines / packages
linear	<code>Rlm</code> , <code>Ranova</code>
GLM	<code>Rglm</code> , <code>Rmgcv</code> / <code>Rgam</code> (additive nonparametric)
nonlinear	<code>Rnls</code>
nonparametric	<code>Rlocpoly</code> , <code>Rlocfit</code>
mixed models	<code>Rlmm</code> , <code>Rnlme</code> , <code>RglmmML</code> , <code>RglmmPQL</code>
times series	<code>Rar</code> , <code>Rarma</code> , <code>Rarima</code> , <code>Rarima0</code> , <code>Rgarch</code>
classification and regression trees	<code>Rtree</code> , <code>Rrpart</code> , <code>Rparty</code>

packages: `RMASS`, `Rstats`, `RKernSmooth`, `Rtseries`

Test for normality

```
library(KernSmooth)
f <- bkde(kredit$age)
plot(f, type="l", xlim=range(f$x), ylim=range(f$y))
                                                    ## normal distribution?
title("Distribution of Age")

t <- shapiro.test(kredit$age)
t
t$p.value

library(tseries)
t <- jarque.bera.test(kredit$age)
t
t$p.value
```










Comparing distributions




```
library(KernSmooth)
f0 <- bkde(kredit$age[y==0])
f1 <- bkde(kredit$age[y==1])
plot(f0, type="l", col="blue", xlim=range(c(f0$x,f1$x)),
      ylim=range(c(f0$y,f1$y)))
lines(f1, col="red")
                                     ## equal distributions?
title("Age vs. Default")

t <- ks.test(kredit$age[y==1],kredit$age[y==0])
t
t$p.value

t <- wilcox.test(kredit$age[y==1],kredit$age[y==0])
t
t$p.value
```

Selected tests

test	routine
comparing means (t-Tests)	 <code>t.test</code>
comparing variances (F-Tests)	 <code>var.test</code>
binomial tests	 <code>prop.test</code> ,  <code>binom.test</code>
correlation	 <code>cor.test</code>
rank tests	 <code>wilcox.test</code>
regression	 <code>anova</code>
unit roots (mean reversion)	 <code>adf.test</code> ,  <code>kpss.test</code>

packages:  `stats`,  `tseries`,  `exactRankTests`

“Advanced” mathematics

What is this R?

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics

“Advanced” mathematics
 Optimizing functions
 Interpolation
 Numerical integration

Basics in programming

References

Optimizing functions

example: linear model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$$

generate artificial data:

```
n <- 100; b <- c(-1,3)
x <- matrix(rnorm(n*length(b)),ncol=length(b))
                                     ## regressors
e <- rnorm(n)/4

y <- 1 + x %*% b + e                 ## linear model

l <- lm( y~x ); summary(l)          ## built-in lm
```

→ to optimize:

$$QS = \sum_i (y_i - \mathbf{x}_i^\top \beta)^2$$

(agreed, this is not a very useful example for iterative optimization ;-)

Optimizing functions (cont'd)

optimization (without intercept):

```
QS <- function(b, x, y){ sum( (y - x %*% b)^2 ) }  
                                ## objective function  
  
b0 <- c(0,0)  
opt <- optim(b0, QS, method="BFGS", x=x, y=y)  
                                ## optimization  
  
opt  
sum( (x %*% opt$par - mean(y))^2 )/sum( (y-mean(y))^2 ) ## R^2
```

→ coefficient of determination R^2 might be outside $[0, 1]$

optimization (with intercept):

```
b1 <- c(0,0,0)  
x1 <- cbind(rep(1,n),x)  
opt1 <- optim(b1, QS, method="BFGS", x=x1, y=y)  
                                ## optimization  
  
opt1  
sum( (x1 %*% opt1$par - mean(y))^2 )/sum( (y-mean(y))^2 ) ## R^2
```

Optimizing functions (cont'd)

optimization with gradient:

$$QS = \sum_i (y_i - x_i^T \beta)^2 = (y - X\beta)^T (y - X\beta), \quad \frac{\partial QS}{\partial \beta} = -2X^T y + 2X^T X\beta$$

```
D.QS <- function(b, x, y){ -2* t(x) %*% y + 2* t(x) %*% x %*% b }  
## gradient
```

```
opt2 <- optim(b1, QS, D.QS, method="BFGS", x=x1, y=y)  
opt2  
sum( (x1 %*% opt2$par - mean(y))^2 )/sum( (y-mean(y))^2 ) ## R^2
```

optimization with box constraints (e.g. $\beta_j \geq 0$):

```
b2 <- c(0,0,0)  
x1 <- cbind(rep(1,n),x)  
opt3 <- optim(b1, QS, D.QS, method="BFGS", lower=0, x=x1, y=y)  
opt3  
sum( (x1 %*% opt3$par - mean(y))^2 )/sum( (y-mean(y))^2 ) ## R^2
```


Optimizing functions (cont'd)

optimization with linear constraints (z.B. $\beta_0 \geq 0, \beta_1 + \beta_2 \leq 2$):

$$u\beta - c = \begin{pmatrix} 0 & -1 & -1 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} - \begin{pmatrix} -2 \\ 0 \end{pmatrix} \geq 0$$

```
u <- cbind( c(0,1), c(-1,0), c(-1,0) )
c <- c(-2,0)
```

```
applyDefaults <- function(fn, ...) { ## for transferring further
  function(x) fn(x, ...)           ## parameters to QS and D.QS
}
```

```
b4 <- rep(0.5,3)
opt4 <- constrOptim(b4, applyDefaults(QS, x=x1, y=y),
                   applyDefaults(D.QS, x=x1, y=y), ui=u, ci=c)
opt4
sum( (x1 %*% opt4$par - mean(y))^2 )/sum( (y-mean(y))^2 ) ## R^2
```

Interpolation

→ `R` `approx` for linear, `R` `spline` and `R` `interpSpline` for spline approximation


```
x <- seq(-5,5,by=1)
y <- sin(x)
```

```
xx <- seq(-5,5,by=0.1)
y.approx <- approx(x,y, xout=xx)$y
yy <- sin(xx)
```

```
plot(xx,yy, type="l", col="green")
lines(xx,y.approx, lwd=2)
```

```
library(splines)
sp <- interpSpline(x,y)
lines(predict(sp,xx), col="red")
```

Numerical integration

→  `integrate` for 1-dimensional integration

→  `adapt` for multidimensional integration

```
pnorm(0)
```

```
it <- integrate(dnorm, -Inf,0)
```

```
it
```

```
attributes(it)    ## result is object of class "integrate"
```

```
it$value
```

```
pmvnorm(c(0,0))
```

```
pmvnorm(c(0,0))[[1]]
```

```
library(adapt)
```

```
it <- adapt(2, c(-Inf,-Inf), c(0,0), functn=dmvnorm)
```

```
attributes(it)    ## result is object of class "integration"
```

```
it$value
```

Basics in programming

What is this R?

How to get help?

Some calculations to start with

Data & files

Wonderful world of graphics

Some statistics

"Advanced" mathematics

Basics in programming

- Functions in R

- Conditional instructions, loops

- "Set theory"

- Packages

- DLLs

- Tips & tricks

References

Functions in R

```
myfun <- function(x, a){  
  r <- a*sin(x)  
  return(r)  
}  
myfun(pi/2,2)
```

```
myfun1 <- function(x, a){ a*sin(x) } ## same as myfun  
myfun1(pi/2,2)
```

```
myfun2 <- function(x, a=1){      ## optional parameter with  
  a*sin(x)                       ## default value=1  
}  
myfun2(pi/2,2)  
myfun2(pi/2)
```

```
myfun3 <- function(x, a=NULL){  ## optional parameter  
                                ## without default value  
  if (!is.null(a)){ a*sin(x) }else{ cos(x) }  
}  
myfun3(pi/2,2)  
myfun3(pi/2)
```

Functions in R (cont'd)

```
myfun4 <- function(x, a=1){
  r1 <- a*sin(x); r2 <- a*cos(x)
  return(list(r1=r1,r2=r2))  ## result is a list
}
myfun4(pi/2)

myfun5 <- function(x, a=1, b=2){
  r1 <- a*sin(x); r2 <- b*cos(x)
  return(list(r1=r1,r2=r2))
}
myfun5(pi/2)                ## a=1, b=2 (defaults)
myfun5(pi/2,1,2)           ## a=1, b=2 (explicitely given)

myfun5(pi/2,2)             ## a=2, b=2 (only a explicitely given)
myfun5(pi/2,a=2)          ## a=2, b=2 (only a explicitely given)

myfun5(pi/2,b=3)          ## a=1, b=3 (only b explicitely given)
```

→ input parameters may be omitted (if reasonable); multiple output parameters are in fact elements of a list

Conditional instructions, loops

- ▶ **R**if & **R**else

```
x<- 1; if (x==2){ print("x=2") }  
x<- 1; if (x==2){ print("x=2") }else{ print("x!=2") }
```

- ▶ **R**for & **R**repeat

```
for (i in 1:4){ print(i) }  
for (i in letters[1:4]){ print(i) }  
i <- 0; while(i<4){ i <- i+1; print(i)}  
i <- 0; repeat{ i <- i+1; print(i); if (i==4) break }
```

- ▶ other: **R**ifelse, **R**switch

“Set theory”

```
A <- 1:3      ## vector as a "set"  
B <- 2:6      ## vector as a "set"  
A %in% B  
B %in% A  
  
C <- c("A", "B")  
D <- LETTERS[2:6]  
C %in% D  
D %in% C
```


Packages

- ▶ packages comprise (one or) more functions, are loaded with `library(<Package-Name>)`; available functions in a package can be queried with `library(help=<package-name>)`
- ▶ to create self-written packages, there exist two helpful functions:
 - `package.skeleton(<package-name>)`
generates the appropriate directory structure of the packages with templates for the necessary files
 - `prompt(<Funktion>)`
generates a template for the help text for a function
- ▶ packages may be installed with the according menu item under Windows or with `Rinstall.packages`
- ▶ packages covering a specific topic can be found under **task views** (→ <http://cran.at.r-project.org/web/views/>), to install a task view use e.g.:

```
install.packages("ctv")
library("ctv")
update.views("Econometrics")
```

DLLs

→ manual “Writing R Extensions”

(<http://cran.at.r-project.org/doc/manuals/R-exts.pdf>)

example of a simple C function under Unix/Linux:

```
#include <stdlib.h>
#include <math.h>

/* Compile shared library: gcc -shared -O2 -o mydll.so mydll.c */

int mysum(double *dim, double *x, double *y, double *z)
{
    long i, n;
    n=dim[0];

    for (i=0; i<n; i++)    /* loop over obs */
    {
        z[i] = x[i] + y[i];
    }
    printf ("mysum in C\n");
    return 0;
}
```

DLLs: call from R

```
dyn.load("mydll.so")           ## load the DLL
is.loaded("mysum")              ## is "mysum" available?

d <- 3
x <- 1:3
y <- 4:6
z <- rep(0,3)

r <- .C("mysum", dim=d, x=x, y=y, z=z )  ## that doesn't work!
r$z

d <- as.double(3); x <- as.double(1:3)
y <- as.double(4:6); z <- rep(0.0,3)
r <- .C("mysum", dim=d, x=x, y=y, z=z )  ## this works!
r$z
z                                         ## z is still =0

r <- .C("mysum", dim=d, x=x, y=y, z=z, DUP=FALSE)
                                         ## another way (without copying)
r$z
z                                         ## z contains the result

dyn.unload("mydll.so")                ## unload the DLL
```

Tips & tricks

editors:

- ▶ Windows: TinnR (<http://sciviews.org/Tinn-R/>)
- ▶ Windows/Unix/Linux: ESS = “Emacs speaks statistics” from <http://ess.r-project.org/> and add a line to `.emacs`:

```
(load "<path to ESS>/ess-5.1.24/lisp/ess-site")
```
- ▶ syntax highlighting for Windows is also available in WinEdt (<http://cran.r-project.org/web/packages/RWinEdt/index.html>)

alternative R GUIs:

- ▶ R Commander (<http://socserv.mcmaster.ca/jfox/Misc/Rcmdr/>) is a menu-assisted environment for data analysis, can be installed by:

```
install.packages("Rcmdr")
```
- ▶ R Studio (<http://rstudio.org/>) is an alternative GUI, available for different OS

Tips & tricks (cont'd)

programming:

- ▶ `eval` and `Rparse` evaluate strings as expressions, e.g.:

```
eval(parse(text=paste("x.", as.character(1:2), " <- 0", sep=" ")))  
print(x.1)
```

- ▶ there are two methods for OOP in R: S3- and S4-classes; for obtaining information about the components of a S3 class (former approach) one uses `Rclass` and `Rattributes` while for a S4 class (newer approach) `RgetClass`, `Rslot`, `RslotNames` are useful
- ▶ methods can be class-dependent, e.g. `methods(print)` gives all functions belonging to the print function

Tips & tricks (cont'd)

diverse tips:

- ▶ rounding and formatting of numbers works with `Rround`, `Rfloor`, `Rceiling`, `Rsignif`, `RformatC`
- ▶ strings (character vectors) can be edited with `Rpaste`, `Rsubstr`, `Rnchar`, `Rstrsplit`, `Rtoupper`, `Rtolower`, `Rsub`
- ▶ time dates can be generated with `Ras.POSIXlt` and `Rstrptime`, e.g.
`as.POSIXlt(strptime("20050101", "%Y%m%d"))+(0:364)*86400`
creates all days of the year 2005;
`d <- as.POSIXlt(strptime("20050926", "%Y%m%d"))`; `d$wday`
shows the weekday of Sep 26, 2005
- ▶ `Rsystem` executes an OS command, e.g. under Linux
`system("cal 09 2005")`
- ▶ `Rxtable` (package: `Rxtable`) and `Rlatex` (package: `RHmisc`) can save R object into LaTeX code

References

- Becker, R. A. and Chambers, J. M. (1984). *S. An Interactive Environment for Data Analysis and Graphics*. Wadsworth and Brooks/Cole, Monterey.
- Becker, R. A., Chambers, J. M., and Wilks, A. R. (1988). *The New S Language*. Chapman & Hall, London.
- Chambers, J. M. and Hastie, T. J. (1992). *Statistical Models in S*. Chapman & Hall, London.
- Dalgaard, P. (2002). *Introductory Statistics with R*. Springer. ISBN 0-387-95475-9.
- Gentleman, R. (2008). *R Programming for Bioinformatics*. Computer Science & Data Analysis. Chapman & Hall/CRC, Boca Raton, FL. ISBN 978-1-420-06367-7.
- Ligges, U. (2009). *Programmieren mit R*. Springer-Verlag, Heidelberg, 3rd edition. ISBN 978-3-540-79997-9, in German.
- Murrell, P. (2005). *R Graphics*. Chapman & Hall/CRC, Boca Raton, FL. ISBN 1-584-88486-X.
- Venables, W. N. and Ripley, B. D. (2000). *S Programming*. Springer. ISBN 0-387-98966-8.
- Venables, W. N. and Ripley, B. D. (2002). *Modern Applied Statistics with S. Fourth Edition*. Springer. ISBN 0-387-95457-0.
- Wollschläger, D. (2010). *Grundlagen der Datenanalyse mit R: Eine anwendungsorientierte Einführung*. Statistik und ihre Anwendungen. Springer.