About t	his Tu	utorial
---------	--------	---------



About this Tutorial	Basics of R 0000000 000000 0000000 00000000 000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 0000000 000000000 000000	Statistics with R	Modelling	Graphics	
					Outline					
					About this	Tutorial				
	An Introduction	to the R Env	ironment		Basics of I Objects	R s and arithmetic				
	Pe	ter Dalgaard			Matrix calculus Important functions					
	Ce Copenha	nter for Statistics gen Business School			Workin Progra	g with data fram mming	es			
	MPAS L	ecture April 2010			Statistics v	with R				
					Modelling					
					Graphics					
About this Tutorial	Basics of R 0000000 000000 0000000 00000000 000000	Statistics with R	Modelling	1/70 Graphics	About this Tutorial	Basics of R 0000000 000000 000000 0000000 00000000	Statistics with R	Modelling	2/70 Graphics	
Practical	ities				Plan					
SheFooSci	ort introduction (app cus on things releva ript of demos on MF	prox. 90 min) nt to your project AS web page			 Eleme D N W R Basic Mode 	entary things abo pata types and sor latrix calculus Vorking with data s as a programmin statistics and te ling tools	out R ne important functi sets ig language sts	ons		

Elementary Graphics

About	this	Tutorial	
-------	------	----------	--

The R environment

Source dialect of the S language

be using Linux here).

Extensible with user functions

Statistics with R

Built around the programming language R, an Open

Command-line execution based on function calls

Workspace containing data and functions

R is Free Software, and runs on a variety of platforms (I'll

Various graphics devices (interactive and non-interactive)

Modelling

Graphics

About this Tutorial

Graphics

Objects and arithmetic

R is a vectorized language

- The basic data type in R is a vector
- Vectors often represent data (e.g. the age for each participant in a study), but also other things like regression coefficients, plot limits, cut points, etc.

Statistics with R

- Data types: Numeric (integer/double), character (strings), logical (TRUE/FALSE)
- Factor (really integer + level attribute) for categorical variables
- Lists (generic vectors)

7/70									8/70
About this Tutorial	Basics of R 0000000 000000 0000000 00000000 000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R ○○○○○○ ○○○○○○ ○○○○○○○ ○○○○○○○○ ○○○○○○	Statistics with R	Modelling	Graphics
Objects and arithmetic					Objects and arithmetic				
Basic operation	ons				Demo 1				

- Standard arithmetic is vectorized: x + y adds each element of x to the corresponding element of y
- Recycling: If operating on two vectors of different length, the shorter one is replicated (with warning if it is not an even multiple)
- ▶ c concatenate: c(7, 9, 13)
- seq sequences: seq(1, 9, 2), short form: 1:5 is the same as seq(1, 5, 1)
- rep replication rep(1:3, 3:1) (1 1 1 2 2 3)
- sum, mean, range, ...

```
x <- round(rnorm(10,mean=20,sd=5)) # simulate data
x
mean(x)
m <- mean(x)
x - m # notice recycling
sqrt(sum((x - m)^2)/9)
sd(x)</pre>
```

About this Tutorial	Basics of R 0000000 00000 000000 0000000 0000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 0000000 000000000 000000	Statistics with R	Modelling	Graphics	
Objects and arithmetic					Objects and arithmetic					
Smart inde	xing			Extended data types						
a[5] single element					The base form m	asic vector type nore complex da	s can be combine ata structures	d and extende	d to	
a[c(5,6,7)] several elements					Attributes extend a basic type with further information.					
a[-6] all except the 6th					E.g., a vector can have a names attribute, for more readable printing					

- Classes have two main functions:
 - Hide details
 - Allow function dispatch (functions that behave differently depending on the class.



Factors

 Factors are used to describe nominal variables (the term originates from *factorial designs*)

a [b>200] index by logical vector

▶ a ["name"] by name

- Internally, they are just integer codes plus a set of names for the *levels*
- They have class "factor" making them (a) print nicely and (b) behave consistently
- A factor can also be ordered (class "ordered"), signifying that there is a natural sort order on the levels
- In model specifications, factors play a fundamental role by indicating that a variable should be treated as a classification rather than as a quantitative variable.

Lists (generic vectors)

- A vector where the elements can have different types
- Functions often return (classed) lists
- Indexing:
 - ▶ lst\$A
 - Ist[[1]] first element
 - Ist[1] list containing the first element

About this Tutorial	Basics of R 000000● 000000 0000000 00000000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R ○○○○○○○ ○○○○○○○ ○○○○○○○○ ○○○○○○○○○	Statistics with R	Modelling	Graphics		
Objects and arithmetic					Matrix calculus						
Demo 2					Elementary	v matrix man	ipulations				
(indexing, f	actors, lists)			 Matrices are implemented as vectors with a dim attribute (of length 2) Constructor function: matrix (1:4,2,2) Indexing in the usual way M[i,j], with all the features of "smart indexing". M[,j] is <i>j</i>th column, etc. Special feature for matrices and arrays: <i>Matrix indexing</i>, M[A] where A has as many columns as M has dimensions. 							
About this Tutorial	Basics of R	Statistics with R	Modelling	15/70 Graphics	About this Tutorial	Basics of R	Statistics with R	Modelling	16/70 Graphics		
Matrix calculus					Matrix calculus	000000 0000000 00000000000000000000000					
Matrix algel	ora				Demo 3						
 R cont A %*% solve matrix 	ains a pretty ful B for matrix m (A, b) for solv inverse)	l set of primitives ultiplication ving linear equatio	for matrix calcu ons. (solve (A	ulus	<pre>Permutatic perm <- sa n <- lengt M <- matr: M[cbind(1 M perm</pre>	on matrix (<i>Mx</i> per ample(5) # w/o th(perm) ix(0,n,n) :n,perm)] <- 1	ermutes the eleme	ents of <i>x</i>)			

• t (A) for transpose of a matrix.

```
M
perm
M %*% 1:n
```

al Basics of R Statistics with R Modelling Graphics About this Tutorial	Basics of R Statistics with R Modelling ○○○○○○○ ○○○○○○○ ○○○○○○○○ ○○○○○○○○○
---	--

Matrix calculus

Matrix calculus

About this Tutori

Other matrix techniques

- diag has multiple functions: creation of diagonal matrices, extracting, and manipulating the diagonal of a matrix.
 Beware: diag (v) is ambiguous if v can have length 1.
- row(X), col(X) are convenient for generating some forms of matrices.
- upper.tri and lower.tri generate indexes for accessing the upper/lower triangle of a matrix.
- Matrices can be "glued together" using cbind and rbind

Row and column matrices

- R usually treats vectors as row or column matrices "as appropriate" (i.e., it guesses)
- E.g., you can left- or right-multiply a vector by a matrix, even though the latter formally requires transposition
- ► And even do y %*% x to get the inner product y'x
- If you want to be explicit about it, you can use rbind or cbind to create the appropriate single-row or single-column matrix.



Using drop() and drop=FALSE

- Default: If a dimension has length one, it is dropped from results. M[1,] is a vector, not 1 × n matrix.
- Often convenient, but source of obscure bugs
- Watch out for extreme cases
- Use M[1, drop=FALSE] to prevent this
- Conversely sometimes you get a matrix and want a vector, as in drop (M %*% 1:n)

Some Basic Functions

- Constructors of simple objects
- Single-column modifications
- Modifying and subsetting data frames

Graphics

About this Tutorial	
---------------------	--

Important functions

Constructors

Statistics with R

Modelling

Graphics

About this Tutorial



Statistics with R

Graphics

000000 000000 000000 0000000

Important functions

Demo 4

R deals with many kinds of objects besides data sets

 Need to have ways of constructing them from the command line

Basics of R

000000

- \blacktriangleright We have (briefly) seen the c and <code>list</code> functions
- Notice the naming forms c (boys=1.2, girls=1.1)
- Extracting and setting names with names (x)
- For matrices and arrays, use the (surprise) matrix and array functions. data.frame for data frames.
- It is also fairly common to construct a matrix from its columns using cbind

x <- c(boys = 1.2, girls = 1.1) x names(x) names(x) <- c("M", "F") x matrix(1:4,ncol=2) cbind(x=0:3,"exp(x)"=exp(0:3))</pre>



he factor Function

- This is typically used when read.table gets it wrong
- E.g. group codes read as numeric
- Or read as factors, but with levels in the wrong order (e.g. c("rare", "medium", "well-done") sorted alphabetically.)
- Notice the slightly confusing use of levels and labels arguments.
- levels are the value codes on input
- labels are the value codes on output (and become the levels of the resulting factor)

About this Tutorial	Basics of R	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R ○○○○○○○ ○○○○○○ ○○○○○○ ○○○○○○○	Statistics with R	Modelling	Graphics
Important functions					Important functions				
The cut Fur	nction				Demo 6				

- The cut function converts a numerical variable into k groups according to a set of break points
- The breaks must include **all** k + 1 interval endpoints
- The intervals are left-open, right-closed by default (right=FALSE changes that)
- The lowest endpoint is not included by default (include.lowest=TRUE to change it.)

```
library(ISwR)
age <- juul$age
age <- age[age >= 10 & age <= 16]
range(age)
agegr <- cut(age, seq(10,16,2), right=FALSE, include.lowest=TR)
length(age)
table(agegr)
agegr2 <- cut(age, seq(10,16,2), right=FALSE)
table(agegr2)</pre>
```



- Like data set in other packages
- Technically: Lists of vectors/factors of same length
- Indexed like matrices (Beware, though: Data frames are not matrices) or lists
- Generate from read operation or with data.frame
- Many sample data frames are available in packages

- Normally variables in data frames must be *qualified* by the name of the data frame
- I.e., you have to say airquality\$Month, in case there is another Month in another data frame
- However, if you attach (airquality), then you can access its contents without the qualifier

About this	Tutorial
------------	----------

Working with data frames

mean(oz)

detach()

airquality\$Month

mean(oz, na.rm=TRUE)

mean(Ozone, na.rm=TRUE)

attach(airquality)

Demo 7

Statistics with R

oz <- airquality[airquality\$Month==5,]\$Ozone</pre>

tapply(Ozone, Month, mean, na.rm=TRUE)

Basics of R

000000000

airquality[airquality\$Month==5,]

Modelling

About this Tutorial

Graphics

Basics of R

Graphics

Reading data

Working with data frames

- Simple data vectors can be read using scan()
- Data frames can be read from most reasonably structured text file formats (space separated columns, tab- and comma-delimited files) using read.table() or read.delim2().

Statistics with R

- > A nice expedient is to read from the clipboard
 (read.delim2("clipboard"))
- The foreign package can read files from Stata, SPSS, ...
- Finding the full pathname of a file is tricky. Either change the working directory use file.choose() to browse to it.

```
31/70
                                                                                                                                                                                                                                                         32/70
About this Tutorial
                                Basics of R
                                                           Statistics with R
                                                                                         Modelling
                                                                                                                  Graphics
                                                                                                                                 About this Tutorial
                                                                                                                                                                                            Statistics with R
                                                                                                                                                                                                                           Modelling
                                                                                                                                                                                                                                                   Graphics
                                                                                                                                                                 Basics of R
                                                                                                                                                                  0000000000
                                0000000000
Working with data frames
                                                                                                                                 Working with data frames
```

The workspace

- The global environment contains R objects created on the command line.
- There is an additional search path of loaded packages and attached data frames.
- When you request an object by name, R looks first in the global environment, and if it doesn't find it there, it continues along the search path.
- The search path is maintained by library(), attach(), and detach()
- Notice that objects in the global environment may mask objects in packages and attached data frames

Demo 8

search()
library(ISwR)
data(intake) # From ISwR
ls()
attach(intake)
search()
ls("intake") # show variables in attached data frame
post - pre
rm(intake) # remove data frame
detach() # remove from search path

About this Tutorial	Basics of R 000000 000000 000000 000000000000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 0000000 0000000000	Statistics with R	Modelling	Graphics
Working with data frames					Working with data frames				

A Common Mistake

attach(juul2)
sex <- factor(sex)
tapply(height, sex, mean, na.rm=TRUE)
detach()
attach(subset(juul2, age > 25))
sex <- factor(sex)
tapply(height, sex, mean, na.rm=TRUE)</pre>

You get an error saying that ${\tt height}$ and ${\tt sex}$ are of different length. What went wrong?

Second time around, sex was found in the global environment *before* the attached data frame.

Subsetting Data Frames

- The syntax for indexing data frames easily gets heavy: airquality[airquality\$Month == 5 & airquality\$Ozone > 50,]
- The subset function allows you to say subset (airquality, Month == 5 & Ozone > 50). I.e., it evaluates the second argument within the data frame.



Attaching or not

- Attached data frames can lead to confusion, especially if you modify their variables
- On the other hand you do probably want the simplified notation for convenience (and things like axis labels do depend on it)
- I tend to recommend attaching only after making all necessary transformation. Others are more radical and discourage attach entirely
- In either case, it is useful to know some functions that allow you to work with the simplified notation, without actually attaching
- subset is one example, the others are transform, with, and within

Demo 9

Notice that within and transform are very similar, but within allows a bit more flexibility.

About this Tutorial	Basics of R	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R	Statistics with R	Modelling
	0000000					0000000		
	000000					000000		
	0000000					0000000		
	000000000					000000000		

Programming

R programming

- R is a full programming language
 - Flow control
 - User-written functions

00000000

- You will soon be writing your own R functions
- For repeated ad-hoc tasks
- Or, because functions are required input for certain tasks (e.g., optimizers)
- User-written functions are not substantially different from system functions, making R very smoothly extensible.

Programming

Argument matching

- > A very simple function logit <- function(p)
 log(p/(1-p))</pre>
- Usage: logit(0.5)

00000000

- The return value of a function is the result of the last expression, unless there is an explicit return()
- ► Formal arguments (p)
- Actual arguments (0.5)
- Positional matching: plot (x, y)
- Keyword matching:t.test(x ~ g, mu=2, alternative="less")
- Partial matching: t.test(x ~ g, mu=2, alt="l")

				39/70					40/70
About this Tutorial	Basics of R 0000000 000000 0000000 00000000 000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 0000000 000000000 000000	Statistics with R	Modelling	Graphics
Programming					Programming				
Scoping					Flow control				

- Objects defined in a function will generally be local to that function and *usually* become unavailable when the function exits
- However, the full story is more complicated
- ...so we skip most of it.
- One basic point is that a function can "see" objects in the parent in which it was defined (*lexical scoping*).

- ▶ if/else
- switch()
- ► for loops
- repeat, while

Graphics

About this Tutorial	Basics of R 0000000 000000 000000000000000000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 000000000000000000000	Statistics with R	Modelling
Programming					Programming			
Conditiona	I Expression	S			The for lo	оор		
if (paire xok < else { yok < xok < }	ed) - yok <- comple - !is.na(y) - !is.na(x)	te.cases(x, y)			for (i in for (i in for (ns i for(pkg i	n 1:n) n names) n list()) n getOption("de	 efaultPackages")) {

- twopi <- if(clockwise) -2*pi else 2*pi
 - Notice that the condition is a scalar. It doesn't vectorize; only one branch is taken. (Compare the ifelse() function)
 - Conditions can be combined using & & and operators
 - An if expression does have a result, which can be used as in the 2nd example above

About	this	Tutorial

Programming

Not Using for Loops

- Many applications of for loops have the following structure
 - Allocate a list/vector to hold the results
 - Loop, saving the results of each iteration in turn
- (A common buglet is that people *extend* the vector on every iteration, which can become terribly inefficient)
- This structure can be abstracted into a single function call lapply(lst, fun)

which causes fun to be applied to each element of lst, and returns a list of the results. (Further arguments can be added and will be passed on to fun)

Inside the body of a for loop, the loop variable takes on the value of each element in turn

Notice that the loop is over a vector or a list

- Even when it is a numeric sequence, the entire vector is stored in memory. Fortunately, this is rarely a problem in practice
- You can skip to the next element with next or exit the loop completely with break,

			43/70)				44/70
Basics of R 0000000 000000 0000000 0000000000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 0000000 00000000 000000	Statistics with R	Modelling	Graphics
				Programming				

Further Apply-functions

- lapply list-apply
- sapply simplifying apply
- tapply tabulating apply
- apply, sweep along slices of tables
- replicate repeat expression

Graphics

About this Tutorial	Basics of R
	0000000
	000000

Statistics with R

Modelling

Graphics About this Tutorial

Basics of R 00000000 000000 0000000 Graphics

Simple Descriptives

- mean, median, sd, etc.
- > quantile(x,p) where p is a vector of proportions
- (actually, there is nine different types of quantiles)
- summary gives some key quantities or a variable, depending on its type. This also works on entire data frames

Tabulation

- For simple tables of discrete variables, use the table function, as in table (sex, tanner), or xtabs
- For tables of descriptives the first choice is tapply, for example tapply (age, tanner, mean, na.rm=TRUE

Statistics with R

Explanation: age is split according to groups and mean is called on each piece with an *extra argument*, evaluating mean(age, na.rm=TRUE) within each group.



- Some variations over tapply is given by the by and aggregate functions
- Multiway tables are often hard to read and use for presentation purposes. Look into the ftable ("flattened" tables) and Martyn Plummer's stat.table function in the Epi package.

- Continuous data by group: t.test, wilcox.test, oneway.test, kruskal.test
- Categorical data: prop.test, chisq.test, fisher.test
- Correlations: cor.test, with options for nonparametrics

About this Tutorial	Basics of R 0000000 000000 0000000 00000000 000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 00000000 000000 0000000 000000000 00000	Statistics with R	Modelling	Graphics
Demo 10					Demo 11				
library(ISw attach(inta t.test(pre, detach()	R) ke) post, paired	=TRUE)			<pre>caesar.shoe chisq.test(c fisher.test) x <- caesar. n <- margin. rbind(x,n) prop.trend.t</pre>	caesar.shoe) (caesar.shoe) shoe[1,] table(caesar. cest(x,n)	shoe,2)		
About this Tutorial	Basics of R 0000000 000000 0000000000 000000000 0000	Statistics with R	Modelling	52/70 Graphics	About this Tutorial	Basics of R 00000000 000000 00000000000 000000000	Statistics with R	Modelling	53/70 Graphics

Modeling Tools: Overview

- Model formulas
- Model objects and summaries
- Comparing models
- Evaluating model fit (plot methods)
- Generalized linear models

Model formulas

- Linear model, $y = X\beta + \varepsilon$
- In practice something like

 $y = eta_0 + eta_1 imes ext{height} + eta_2 imes extsf{1}_{(type=2)} + eta_3 imes extsf{1}_{(type=3)} + arepsilon$

Wilkinson-Rogers formulas:

y = height + type

(Interpretation depends on whether variables are categorical or continuous)

Abou	t this	Tutorial	
------	--------	----------	--

Statistics with R

Modelling

Graphics About this Tutorial



Statistics with R

Graphics

Model formulas in R

- R representation y ~ height + type where type is a factor
- Interactions a:b, a*b = a + b + a:b
- Algebra (a: (b + c) = a:b + a:c etc.)
- Notice special interpretation of operators

Fitting linear models

- Im generates a fitted model object
- Extract information from model object
- Fit other models based on model object



About this Tutorial	Basics of R 00000000 0000000 00000000 0000000000	Statistics with R	Modelling	Graphics	About this Tutorial	Basics of R 0000000 000000 0000000 000000000 000000	Statistics with R	Modelling	Graphics
Demo 12					R graphics				

The standard interface (i.e., not the lattice package)

- Customizing plots
- Graphics parameters
- Math on plots



Standard R graphics

Basic x-y plots

- Ink on paper model; once something is drawn it cannot be erased.
- Sensible default plots
- Arguments can override defaults
- Options to turn off various elements of plots (e.g. the axes)
- Functions to *add* elements.

- The plot function with one or two numeric arguments
- Scatterplot or line plot (or both) depending on type argument: "l" for lines, "p" for points (the default), "b" for both, plus quite a few more.
- Functions for adding to a plot: lines, points, segments, abline, text, mtext, axis
- Also: formula interface, plot (y~x)

About	this	Tutorial	
-------	------	----------	--

Statistics with R

Modelling

Graphics About this Tutorial

Specific plots

- Histograms hist(x)
- Density plots plot (density (x))
- Boxplots boxplot (x)
- Barplots barplot (x) (x can be a matrix)
- ▶ Pies pie()
- Matrix plots (multiple y columns) matplot()

Graphical parameters

- Arguments to plot et al. (67 possibilities!)
- The par function can be used to set most of them persistently. Most info is found via help(par)
- Look them up! Here are some of the more commonly used:
 - Point and line characteristics: pch, col, lty, lwd

Statistics with R

- Multiframe layout: mfrow, mfcol
- ► Axes: xlim, ylim, xaxt, yaxt, log

About this Tutorial Basics of R Statistics with R Modelling Graphics About this Tutorial Basics of R Statistics with R Modelling Graphics About this Tutorial Basics of R Statistics with R Modelling Graphics Observed Statistics are served by the statistic served by the s

- Sort of like TeX
- Works on unevaluated expressions (quote(alpha), expression(alpha))
- Special conventions: ^,[] sub/superscript, special names alpha, sum, int
- See help (plotmath)

Graphics

Modelling

Demo 14

70/70