# R and R-Studio

## Bjarki Þór Elvarsson and Einar Hjörleifsson

Marine Research Institute

# Overview

1. Introduction to R

2. RStudio

3. A bit of R–programming

# Overview

1. Introduction to R

2. RStudio

3. A bit of R–programming

# Why R?

R has become the lingua france of statistical analysis and data wrangling

- Its free! If you are a teacher or a student, the benefits are obvious

- It runs on a variety of platforms including Windows, Unix and MacOS

- It provides an unparalleled platform for programming new statistical methods in an easy and straightforward manner

- It offers powerful tools for data exploration and presentation

# Why not Excel

On the accuracy of statistical procedures in Microsoft Excel 97

B.D. McCullough ✉, Berry Wilson

⊞ Show more

Get rights and content

Abstract

The reliability of statistical procedures in Excel are assessed in three areas: estimation (both linear and nonlinear); random number generation; and statistical distributions (e.g., for calculating $p$-values). Excel's performance in all three areas is found to be inadequate. Persons desiring to conduct statistical analyses of data are advised not to use Excel.

# Why not Excel



Computational Statistics & Data Analysis

Volume 31, Issue 1, 28 July 1999, Pages 27–37

Computational Statistics & Data Analysis

Volume 40, Issue 4, 28 October 2002, Pages 713–721

## On the accuracy of statistical procedures in Microsoft Excel 2000 and Excel XP

B D. McCullough[a], ✉ , Berry Wilson[b], ✉

⊞ Show more

Get rights and content

### Abstract

The problems that rendered Excel 97 unfit for use as a statistical package have not been fixed in either Excel 2000 or Excel 2002 (also called "Excel XP"). Microsoft attempted to fix errors in the standard normal random number generator and the inverse normal function, and

# Why not Excel

# Why not Excel



**Computational Statistics & Data Analysis**

Volume 49, Issue 4, 15 June 2005, Pages 1244–1252

**Computational Statistics & Data Analysis**

Volume 52, Issue 10, 15 June 2008, Pages 4570–4578

## On the accuracy of statistical procedures in Microsoft Excel 2007

B.D. McCullough[a], David A. Heiser[b]

⊞ **Show more**

doi:10.1016/j.csda.2008.03.004                    Get rights and content

### Abstract

Excel 2007, like its predecessors, fails a standard set of intermediate-level accuracy tests in three areas: statistical distributions, random number generation, and estimation. Additional errors in specific Excel procedures are discussed. Microsoft's continuing inability to correctly fix errors is discussed. No statistical procedure in Excel should be used until
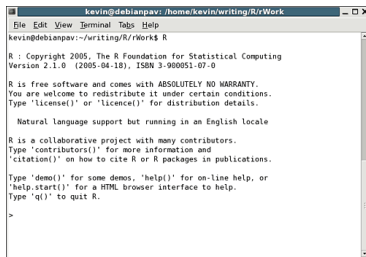
# What is R?

# What is R (cont)

- R is commandline driven:
  - its biggest appeal as one can reuse commands
  - its biggest hurdle in widespread use

- R is open-source:
  - Other statistical software packages can be extremely expensive
  - Large user base with almost all statistical methods implemented

# R - interfaces

# R - interfaces

# R - interfaces

# R - interfaces

# R - interfaces

# R - interfaces

# Overview

# RStudio

- RStudio allows the user to run R in a more user-friendly environment. It is open- source (i.e. free) and available at www.rstudio.com

- built to help you write R code, run R code, and analyze data with R

- text editor, version control, project handling, markdown support, keyboard shortcuts, debugging tools, and much more

# Using R-studio

## R-console

One can type commands directly into the console to perform calculations:

```
1+1 ## this should be 2

## [1] 2

mean(c(1,2,3,4,5)) ## mean of all numbers from 1 to 5

## [1] 3
```

and assign the results with the "<-" sign:

```
x <- 1+1 ## save the result into a variable named x
x        ## print out the value of x to screen

## [1] 2
```

# Using R as a calculator

R can be used as calculator:

```
1+1       # 1 plus 1 equals 2;)
2*2       # 2 times 2 equals 4
3^2.5     # 3 to the power of 2.5
exp(5)    # e to the power of 5
log(5)    # natural logarithm (ln) of 5
log10(5)  # Base 10 log of 5
1/2       # 1 over 2 equals 0.5
```

## Data types

R has a number of data types to handle the various inputs from the user:

- Numbers: integers, reals and complex

- Strings: letters, words, files etc..

- Logicals: TRUE or FALSE

- Factors: integer numbers that correspond to a fixed (limited) set of values

- NA's: Not Available, used when data is missing

# Data types

```r
x <- 1.2          ## number
y <- 'A'          ## string
z <- TRUE         ## logical
w <- factor('A')  ## factor
```

## Data structures

One of R's many strengths are multiple data structures:

- **vectors** (1d)

- matricies (2d)

- arrays (nd)

- **dataframes** (2d)

- lists (1d, sort of)

- ...

# Vectors

Vectors can contain:

- Numbers
- Strings
- Logicals
- NA's

```
x <- c(1,2,3)          ## vector of numbers
y <- c('a','b','c')    ## vector of strings
z <- c(TRUE,FALSE)     ## vector of locigals
class(x)               ## returns the type of vector
u <- NA                ## NA
```

## Useful commands for vectors

```
seq(1,10)        ## makes a vector from 1 to 10
1:10             ## same
rep(2,3)         ## makes a vector that contains 2, 3 times
length(y)        ## returns the length of a vecor named y
sort(y)          ## sorts a vector into ascending or
                 ## descending order
cut(y)           ## divides the range of y into intervals and
                 ## makes a factor variable
as.character(y)  ## changes y to a character vector
cbind(x1,x2)     ## binds vectors by columns
rbind(x1,x2)     ## binds vectors by rows
```

# Useful commands for vectors (II)

```
mean()        # mean
median()      # median
quantile()    # quantiles
summary()     # depends on what we feed it with...
sd()          # standard deviation
var()         # variance
range()       # range
min()         # smallest value
max()         # largest value
```

## Logical statements

```
a == b    ## a is equal to b
a != b    ## a is not equal to b
a > b     ## a is greater than b
a >= b    ## a is greater or equal to b
a < b     ## a is less than b
a <= b    ## a is less or equal to b
a & b     ## a and b
a | b     ## a or b
!a        ## not a
is.na(a)  ## is a equal to NA (missing)
```

# Selecting parts of a vector

```
x[1]           ## get the first element of a vector
x[c(1,2)]      ## get the first and second element
x[-1]          ## get every element except the first one
x[-c(1,2)]     ## get every element except 1 and 2
x[x==3]        ## all elements that are equal to 3
x[!is.na(x)]   ## all elements that are not NA
x[x!=3&x<5]    ## all elements that not equal to 3 and <5
```

## Data-frames

- Data frames are collections of vectors (columns) of the same length, similar to a table in a database

- The columns do not need to be of the same type

- One can access the columns:

  ```
  data$Col.name        ## get column named Col.name
  data[,1]             ## get column number 1
  ```
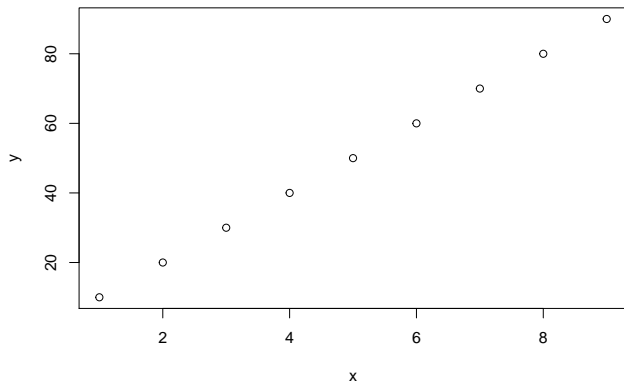
- and the entries in the column:

  ```
  data$Col.name[1]     ## get entry number 1 from Col.name
  data[2,1]            ## get entry number 2 from column 1
  ```
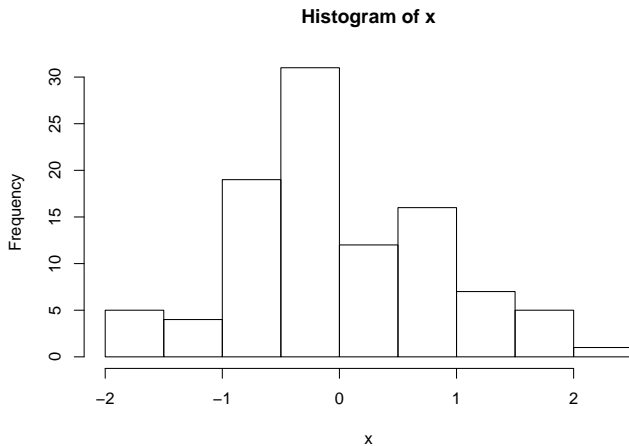
# Base R scatter plot

```
x <- 1:9
y <- 10*(1:9)
plot(x,y)
```

# Base R histogram

```
x <- rnorm(100)
hist(x)
```



**Histogram of x**

# Using R-studio

# Working in R

- In R we use *commands* to create new *objects*. When we do that we have a choice to do two things:
    - The command spits out the results to the screen and then these results forgotten
    - The command saves the output so we can reuse it.

- R can't save the object without giving it a name

```
name <- command
```

- Note that, like in real life, there is not undo button

- One should rather "remember" the commands used to create the object

## Script editor

- Is opened when selecting File->New file->R script from the menu panel

- A R script is basically a series of R commands that can be run in the console

- It is recommended to store all commands that "work" in a script for later reconstruction of a particular analysis

- It is possible to run a single line in a script by pressing [ctrl]+Enter or the ⇥ Run button in the top right corner

# Using R-studio

# Session info



- **Environment**: Contains a list of all declared variables. If you have a dataset you can double click the variable to view the data.

- **History**: Contains a list of all previous commands entered into the console.

# Getting data into R

The environment tab allows you also to import data into R:

# Getting data into R

The environment tab allows you also to import data into R:

# Getting data into R

The environment tab allows you also to import data into R:

# Getting data into R

The environment tab allows you also to import data into R:

# Using R-studio

## Other items–files

**Files**: Shows files on the computer. The user can change the working directory (where R reads data-files and scripts from) by selecting "more->Set as working directory"

# Other items–plots

**Plots**: Graphical output from R. The
user can export these the figures to
file (as jpeg, png or pdf) or to
clipboard

# Other items–help

**Help**: Gives a browsable interface to R's in-line help pages. The user can search for a topic or a specific function

# On help

A help page for a function usually has the following categories:

- **Description**: short description of the function
- **Usage**: how the function is used/called from R
- **Arguments**: What are the function inputs
- **Details**: The nitty gritty on how the function should work
- **Examples**: example uses
- **See also**: other useful functions that are similar

# Help commands

One can also get help from the console:

```
help.start()        # General help page
help(nameOffuntion) # Gives the help page needed
?nameOffunction     # shorthand
help.search('topic') # searches for specific topic
```

# Packages

- One of R's greatest strengths comes from specialized add-ons, called *packages* that the user can download and install

- Users download packages from CRAN's website (Comprehensive R Archive Network)

- Currently at 7811 packages (January 20th 2016)

- All packages comes with both an in-line and pdf manual

# Package growth

# MRI's R packages

- geo – nautical charts (also available on CRAN)

- ora – data base connectivity

- fjolst – all samples from commercial and survey vessels

- Logbooks – logbooks from the commercial fishery

- Rgadget – specialised package to interact with Gadget

## Using packages

- One can use the install.packages command to download and install packages

- To use the package in an R session one needs to load it into memory using the library command

- In R-studio one can also select, load and install packages in the "other items" part

```
install.packages('ggplot2') ## install package named ggplot2
library(ggplot2)            ## load it into memory
```

# Projects

- Rstudio allows us to make things a little bit easier to isolate the task we are working with a any given point by defining specific **projects**

- Projects save the state between sessions. This includes:
    - Working directories
    - Open files/scripts
    - Workspaces
    - Color scheme

- Projects are created by selecting File->New project

# Creating projects

# Creating projects

# Creating projects

# Turning R off

When shutting down R and Rstudio it reminds you to save your work:



In general you should **not** save your workspace unless the analysis takes a long time to reproduce.

# Don't save your workspace

## Excercise 1

- Create a new directory anywhere on your computer and names it "Classwork"
- Open R-studio and create a new project called "classEx1" and associate it with the newly created directory
- Install lubridate, stringr,geo, ggplot2, tidyr, dplyr and gridExtra if you haven't done so already
- Create a new script named 'Ex1.R'
- Look at the help page for 'sum', 'mean' and 'sd'
- Read in the minke whale data using the "import dataset" menu
- Write a script that calculates the sum, mean and standard deviation of the stomach.weight of the minke whales

# Overview

## if sentences

- It is often so that we would like some parts of our code to run only if some condition is fulfilled

- We use `if` sentences for that

- When setting up the conditions the operators come in handy

```
if(condition){
  ## this runs if the condition(s) are fulfilled
}
```

## if sentences - example

```
x <- 10
## check if x is greater than 0
if(x>0){
    print('x is positive')
}
```

## if and else sentences

- Sometimes we would like to check for more than one condition

- We use a combination of if and else sentences for that

```
if(condition_1){
  ## this runs if condition 1 is fulfilled,
  ## then we skip the else-sentences below
} else if(condition_2){
  ## this runs if condition 2 is fulfilled,
  ## then we skip the else-sentence below
} else{
  ## this runs of neither condition 1 or 2 are fulfilled
}
```

# if - else sentences - example

```r
x <- 10

## check if x > 10
if(x>10){
    print('x is greater than 10')
} else if(x>=0){
    print('x is less than 10 but greater than or equal to 0')
} else {
  print('x is negative')
}
```

# ifelse sentences

- It is often so that we only have two conditions

- We use ifelse sentences for these cases

```
ifelse(condition,
       'this runs if the condition is fulfilled',
       'this runs if the condition is not fulfilled')
```

## ifelse sentences - example

Replace NA's from data

```
x <- c(1,2,3,NA,4,5)
## find all missing entries and replace them
x <- ifelse(is.na(x),0,x)
```

## for-loop

- We use a for-loop if we need some part of our code to run more than one time and we know how often it sould run

- We use an index/counter that counts how often we enter the loop

- We try to avoid for-loops because they are rather slow

- There are several commands in R that help us to avoid for-loops - see the end of the slides

```
for(index in vector){
  ## This code runs as often as the length of the vector
}
```

## while-loop

- While-loops are similar to for-loops

- We use a `while`-loop if we need some part of our code to run more than one time but we dont know how often it sould run

- The loop runs while the condition stated is fulfilled

```
while(condition){
  ## This code runs while the condition(s) are fulfilled
}
```

## Loops, examples

```
## find the number of entries in the data
n <- nrow(minke)

## calculate the mean length old fashion way
ML <- 0
for(i in 1:n){
  ML <- ML + dat$length[i]
}
ML <- ML/n
```

## User–defined functions in R

- Even though there are a great number of built in functions/commands/methods in R we sometimes need to write our own functions

- We use function for that

- We can write a name of a built in function in the console to to get the code behind it

- We can take the code and change it as we please

## function

- A function has a name, takes arguments/settings and returns an object

- Object within the function are local within the function

- Objects the function returns can be of any type

```
nameOfFunction <- function(argument1, argument2, ... ){
  ## The code of the function
  return(object)
}
```

Note that you can name your function whatever you want

# function

```
add1 <- function(x){
  x <- x + 1
  return(x)
}

add1(10)

## [1] 11

x

## Error in eval(expr, envir, enclos):  object 'x' not found
```

## Useful string operations

The 'stringr' package adds a number of string operations:

```
str_c()          ## glues strings together
str_length()     ## measure the length of a string
str_sub()        ## select parts of the string
str_str() <-     ## assign parts of the string with new values
str_dup()        ## duplicates string
str_trim()       ## removes trailing white space
str_pad()        ## adds whitespace
```

# Find and replace

Stringr also does find and replace:

```
str_detect()
str_subset()
str_locate()
str_locate_all()
str_extract()
str_extract_all()
str_match()
str_match_all()
str_split()
```

## Class excercise 2

- Read in the minke whale data

- Define a function that calculates the condition factor of a whale using the formula:
$$C = 100 * (W/L^3)$$

- Use the function to calculate the condition factor for all whales that were weighed

- Investigate the distribution of the condition factor

- Assign classification to all whales with body condition less than 0.8 'Malnourished', for 0.8 and above assign "OK" (hint use ifelse)

- Impute the weight for those whales that were not weighed