



Reproducible Research

Niklaus Zemp
22.01.2020

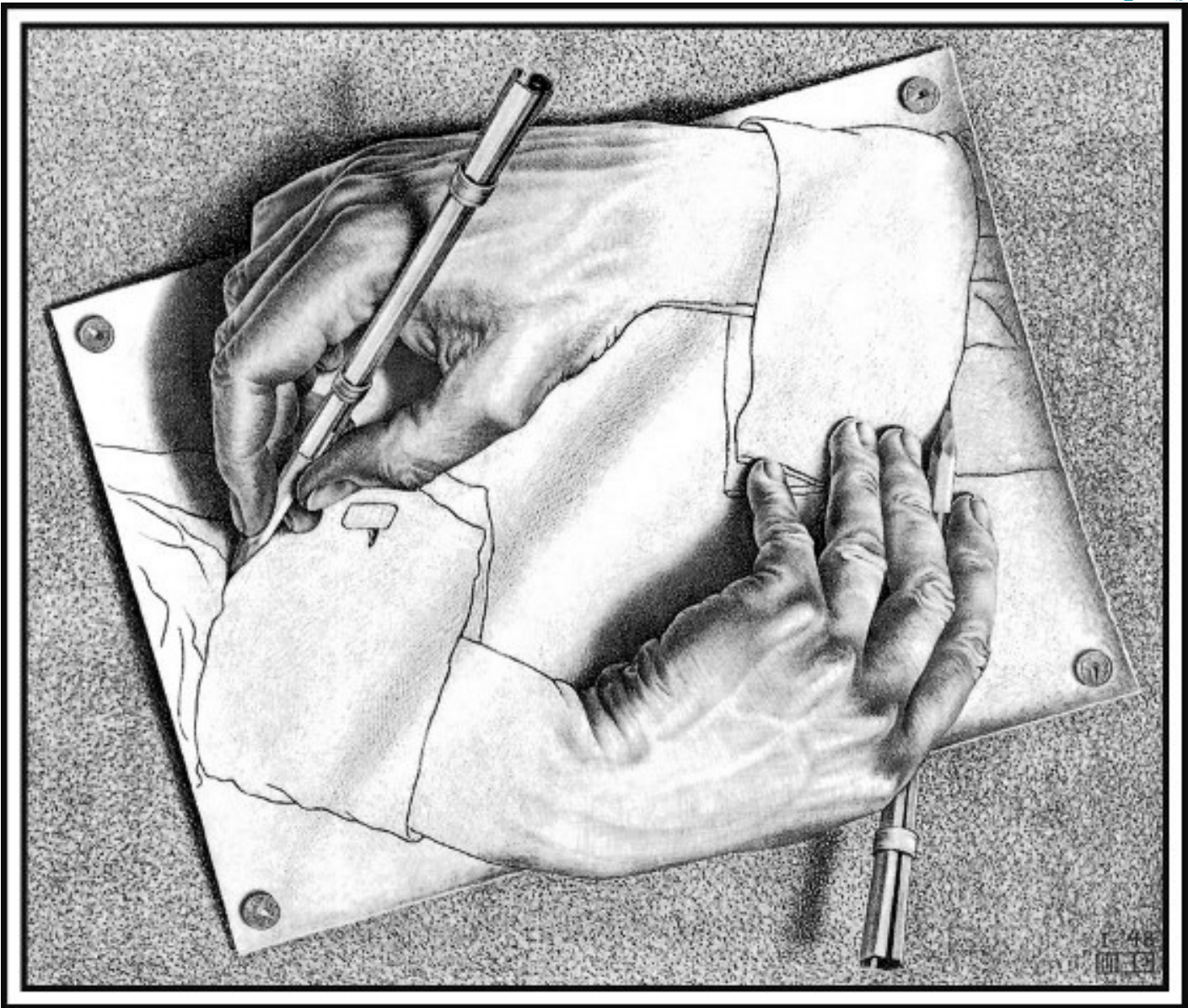
Genetic Diversity Centre (GDC)
Bioinformatics
ETH Zurich



Input (13.30-14.00)

Exercise (14.00-14.45)

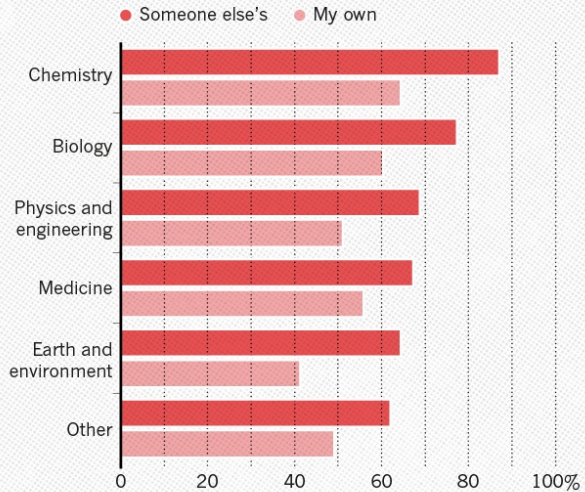
Discussion (14.45-15.00)



Scientific recipe-publications

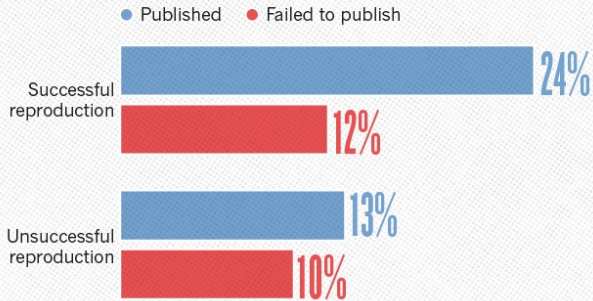
HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



HAVE YOU EVER TRIED TO PUBLISH A REPRODUCTION ATTEMPT?

Although only a small proportion of respondents tried to publish replication attempts, many had their papers accepted.



Number of respondents from each discipline:
Biology 703, Chemistry 106, Earth and environmental 95,
Medicine 203, Physics and engineering 236, Other 233 ©nature

IS THERE A REPRODUCIBILITY CRISIS?



©nature



Reality check on reproducibility

A survey of *Nature* readers revealed a high level of concern about the problem of irreproducible results. Researchers, funders and journals need to work together to make research more reliable.

25 May 2016



Recipe

Alles Leben strömt aus dir

Karoline Rudolphi
1754-1811

J.-H. Tobler
1777-1838

Etwas bewegt

S/A

1. Al - les Le - ben strömt aus dir, al - les
2. Das ich füh - le, was ich bin, das ich
3. Dei - ner Ge - gen - wart Ge - fühl, dei - ner

T/B

1. Le - ben strömt aus dir und durch - wallt in tau - send
2. füh - le, was ich bin, daß ich dich, du Gro - ßer,
3. Ge - gen - wart Ge - fühl sei mein En - gel, der mich

1. und durch
2. daß ich
3. sei mein

1. Bä - chen und durch - wallt in tau - send Bä - chen al - le
2. ken - ne, daß ich dich, du Gro - ßer, ken - ne, daß ich
3. lei - te, sei mein En - gel, der mich lei - te, daß mein

1. wallt in tau - send Bä - chen,
2. dich, du Gro - ßer, ken - ne
3. En - gel, der mich lei - te,

1

(C) Jürgen Knuth

Ingredients

- 3 fresh red chillies
- 2 onions
- 4 cloves of garlic
- 4 large plum tomatoes
- 1 bunch of fresh coriander
- 4 large free-range chicken legs, skin on
- olive oil
- 2 teaspoons garam masala
- 1 tablespoon crumbled dried curry leaves
- 1 tablespoon mustard seeds
- 2 tablespoons white wine vinegar
- fat-free natural yoghurt

Method

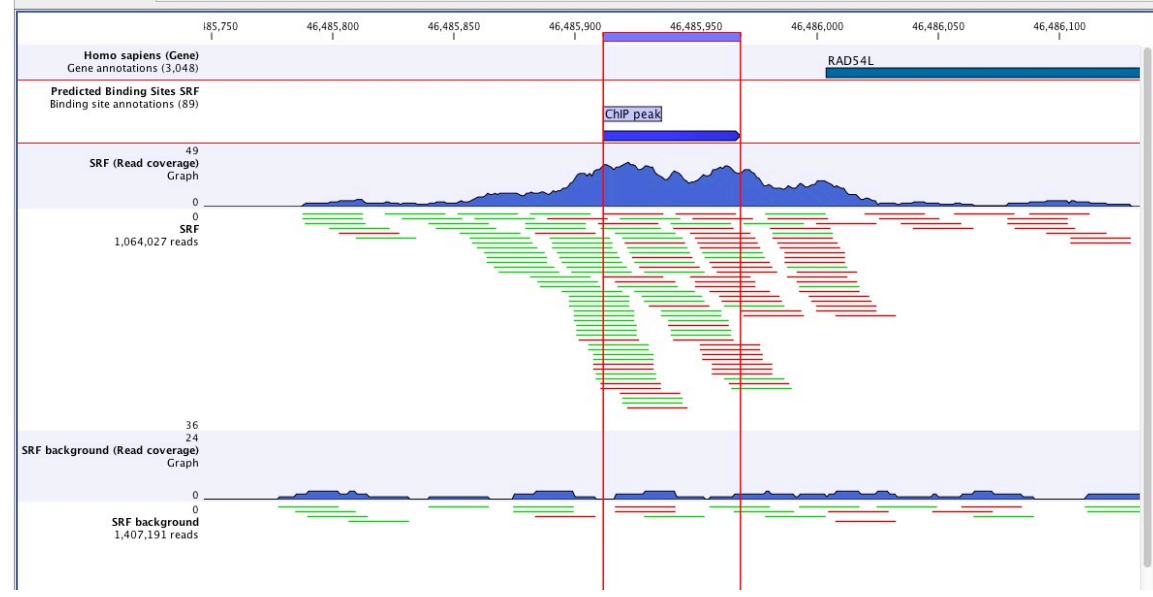
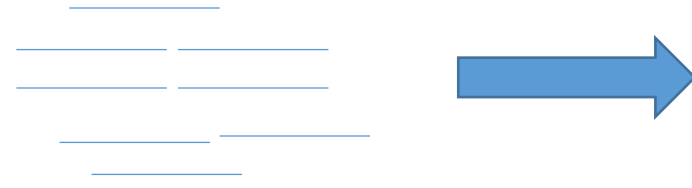
1. Halve the chillies (deseed if you like), peel and finely slice the onions, then peel and crush the garlic. Quarter the plum tomatoes, and pick the coriander leaves.
2. Rub the chicken legs all over with a drizzle of oil and the garam masala, then transfer to a large non-stick ovenproof pan.
3. Add another drizzle of oil and fry the chicken over a medium heat for 10 to 15 minutes, or until lovely and dark golden all over. Be brave and let it get really nice and dark – it will make such a difference to the end result if you get it right at this stage. Drain off any excess fat.
4. At this point, preheat the oven to 180°C/350°F/gas 4.
5. Next, add the curry leaves, mustard seeds, chillies, onion and garlic to the chicken. Cook, stirring often, for 5 minutes, then add the tomatoes and white wine vinegar.

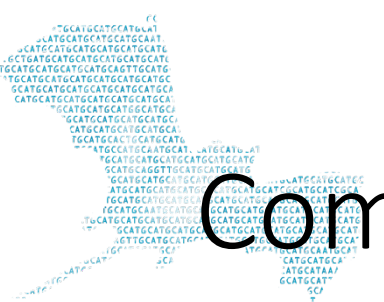
- Transfer the pan to the oven. Cook, uncovered, for 50 minutes, or until the chicken is cooked through and falling off the bone.
6. Pop the pan on the hob and reduce the liquid until sticky. Scatter with the coriander leaves and serve with the cooling yoghurt. Delicious with rice or couscous, and a crisp, refreshing



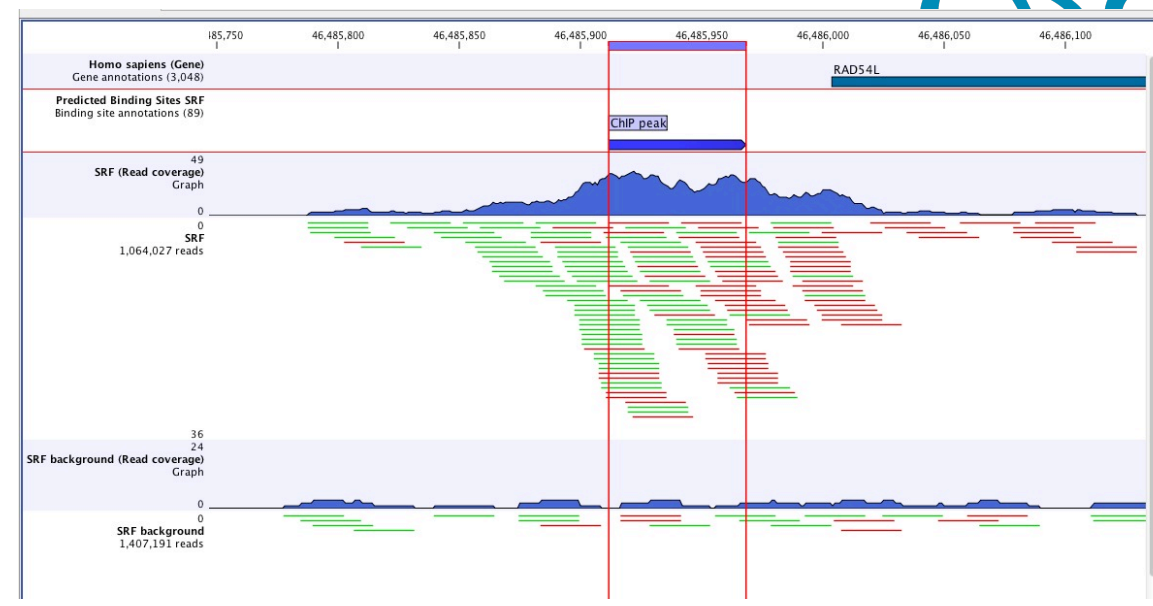


GUI tools





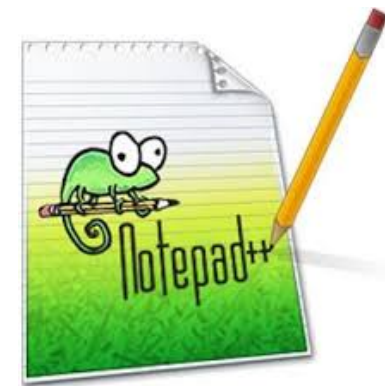
Command line Tools



```
bwa mem Ref Read_R1 Read_R2 > alignment.sam
```



Scientific recipe



Log-file

```
#####
Nik Zemp, niklaus.zemp@env.ethz.ch, GDC, ETH Zurich
#####
####ddRAD log file, Ivo Widmer, p432
#####

#####
####Download data
#####

module load eth_proxy
bsub -n1 -W 4:00 -R "rusage[mem=1000]" \
"wget -r --no-parent --reject="index.htm*" http://gc3fstorage.u
bsub -n1 -W 4:00 -R "rusage[mem=1000]" \
"wget -r --no-parent --reject="index.htm*" http://gc3fstorage.u
bsub -n1 -W 4:00 -R "rusage[mem=1000]" \
"wget -r --no-parent --reject="index.htm*" http://gc3fstorage.u
bsub -n1 -W 4:00 -R "rusage[mem=1000]" \
```

script

```
#!/bin/bash
#BSUB -J "processradtags"
#BSUB -R "rusage[mem=10000]"
#BSUB -n 1
#BSUB -W 24:00

module load gcc/4.9.2 gdc perl/5.18.4
export PATH=$PATH:/cluster/project/gdc/shared/tools/stacks-1.48
#module load gcc/4.8.2 gdc perl/5.18.4 stacks/1.40
source /cluster/apps/gdc/perl5/etc/bashrc

mkdir samples

process_radtags -i gzfastq -f /cluster/project/gdc/people/buckleyj/Hiseq
```



Scientific recipe

- (original) author or source
- your name
- date
- version of the tool
- version of the script
- reproducible code with comments (more comments than code)
- use style guides
- syntax coloring

```
#!/bin/bash
#BSUB -J "processradtags"
#BSUB -R "rusage[mem=10000]"
#BSUB -n 1
#BSUB -W 24:00

module load gcc/4.9.2 gdc perl/5.18.4
export PATH=$PATH:/cluster/project/gdc/s
#module load gcc/4.8.2 gdc perl/5.18.4 s
source /cluster/apps/gdc/perl5/etc/bashr

mkdir samples

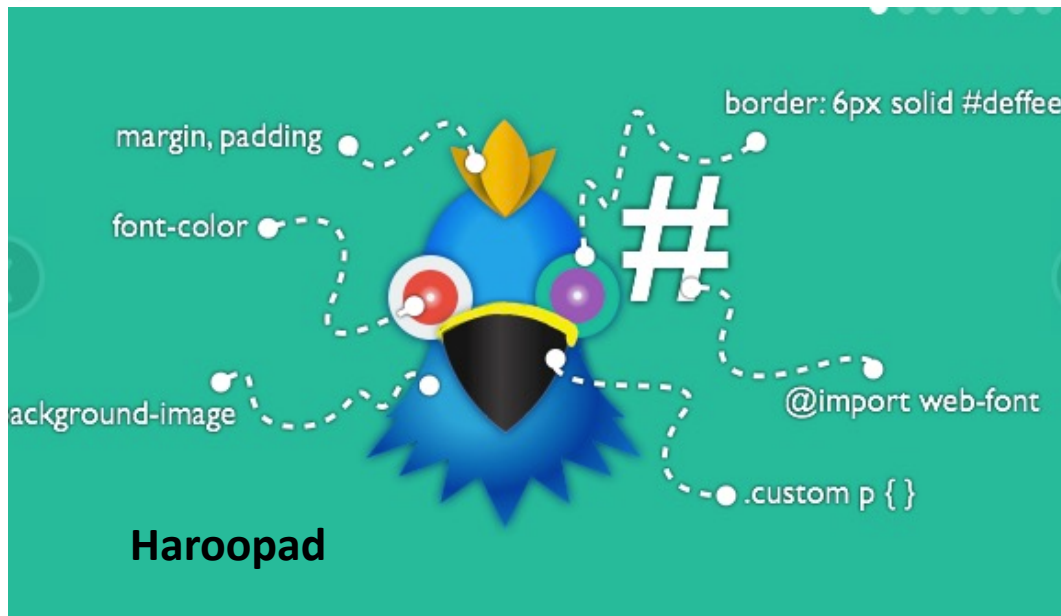
process_radtags -i gzfastq -f /cluster/p
```


Markdown

Slightly modified based on dDocent Version 2.6.1; overlapping paired-end reads; September 2018; Nik Zemp based on the [Tutorial](#)

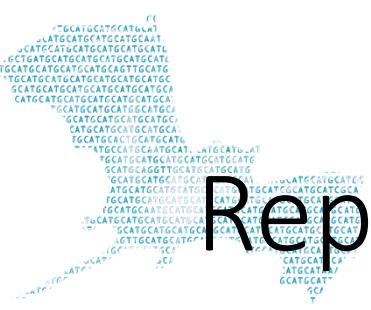
```
bsub -n 2 -W 4:00 -R "rusage[mem=5000]" -Is bash
```

```
module load gcc/4.8.2 gdc python/2.7.11 java/1.8.0_73 perl/5.18.4 freebayes/0.9.20 trimmomatic/0.35 bwa/0.7.12 cd-hit,  
export PATH="$PATH:/cluster/project/gdc/shared/tools/pear-0.9.6-bin-64"  
export PATH="/cluster/project/gdc/shared/tools/seqtk:$PATH"
```



[Cheat sheet](#)





Reproducible Research

scripts
logs



scripts
logs



scripts
logs



Violin

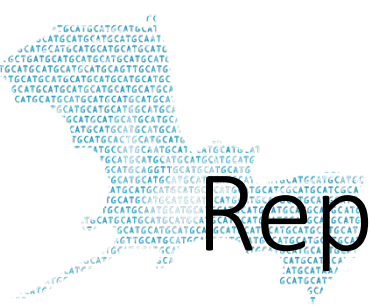


Density



Histogram



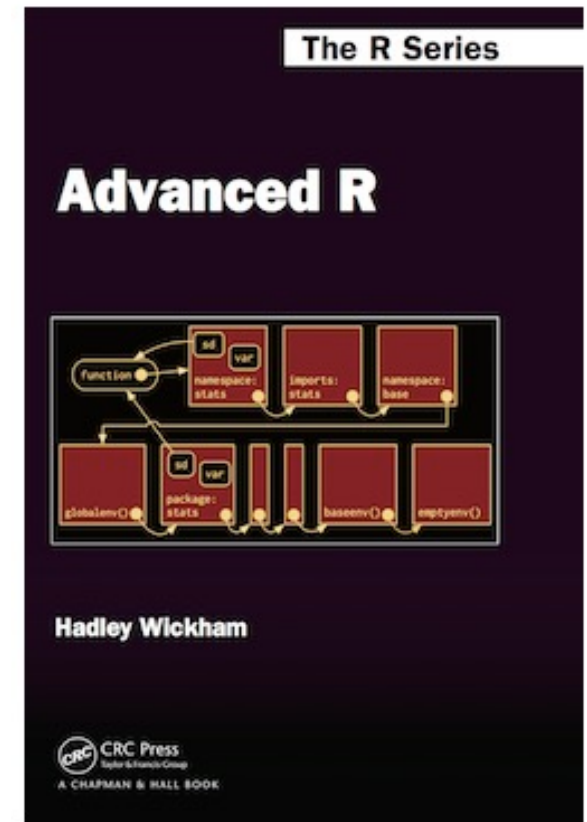
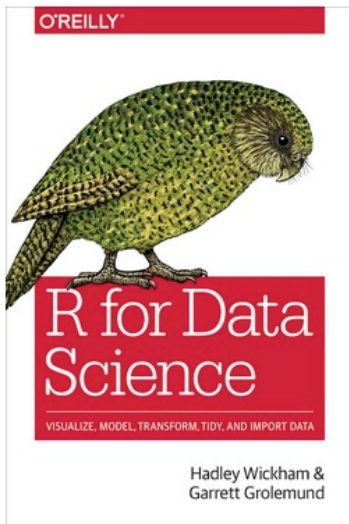


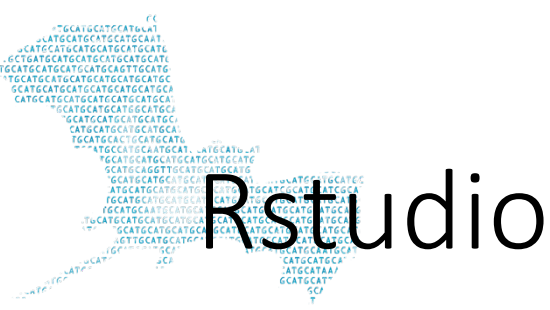
Reproducible Research



Resources

Many many tutorials, forum, YouTube videos posts and books available





Rstudio

The screenshot shows the RStudio interface with four main panes:

- Script:** Contains R code:


```
1 source("My_functions.R")-
2 plot_points(10000)|
```
- Environment:** Shows "Global Environment" and "Environment is empty".
- Files:** Shows a file explorer for the "Test" directory:

Name	Size	Modified
..		
.Rhistory	46 B	Nov 27, 2019, 7:22 PM
Test.Rproj	205 B	Nov 28, 2019, 9:18 AM
Untitled.html	696.6 KB	Nov 27, 2019, 6:58 PM
Untitled.Rmd	796 B	Nov 27, 2019, 6:58 PM
- Console:** Shows the R startup message:


```
~/Desktop/Test/ ➤
Plattform: x86_64-apple-darwin15.6.0 (64-bit)

R ist freie Software und kommt OHNE JEGLICHE GARANTIE.
Sie sind eingeladen, es unter bestimmten Bedingungen weiter zu verbreiten.
Tippen Sie 'license()' or 'licence()' für Details dazu.

R ist ein Gemeinschaftsprojekt mit vielen Beitragenden.
Tippen Sie 'contributors()' für mehr Information und 'citation()',
um zu erfahren, wie R oder R packages in Publikationen zitiert werden können.

Tippen Sie 'demo()' für einige Demos, 'help()' für on-line Hilfe, oder
'help.start()' für eine HTML Browserschnittstelle zur Hilfe.
Tippen Sie 'q()', um R zu verlassen.

> |
```

Script

Variables

Console

Plots and other information



Data management

R project

- source ("function1.R")
- source ("function2.R")
- table.csv
- table2.csv
- Data.RData
- analysis1.R
- Analysis.Rmd



Packages

Available Packages

Currently, the CRAN package repository features 13884 available packages.

[Table of available packages, sorted by date of publication](#)

[Table of available packages, sorted by name](#)



Install »

- Discover [1649 software packages](#) available in *Bioconductor* release 3.8.

Get started with *Bioconductor*

- [Install *Bioconductor*](#)
- [Get support](#)
- [Latest newsletter](#)
- [Follow us on twitter](#)
- [Install R](#)

Tidyverse



```
install.packages("tidyverse")
```




Set up an R script

MDA 2020, GDC Nik Zemp, January 2020, Version 0.1

Set working directory

```
setwd("~/Desktop/")
```

Remove all variables

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

set seed

```
set.seed(1000)
```

load packages

```
library(tidyverse)
```

Source functions

```
source("My_Functions.R")
```



Use R packages

```
library(tidyverse)
```

```
dplyr::filter()
```

```
stats::filter()
```

“ *@ijlyttle a package is a like a book, a library is like a library; you use library() to check a package out of the library #rsats*

— Hadley Wickham (@hadleywickham) [December 8, 2014](#)

```
datawox=jittermap(datawox, amount=1e-6);dataauto=jittermap(dataauto, amount=1e-6) autosome;data=jittermap(data, amount=1e-6);data2=data;dataauto$pheno$gender=as.numeric(dataauto$pheno$gender)
dataauto$pheno$ffd[68] <- NA;dataauto$pheno$ffd=(dataauto$pheno$ffd-
min(dataauto$pheno$ffd,na.rm=T))/(max(dataauto$pheno$ffd,na.rm=T)-min(dataauto$pheno$ffd,na.rm=T));
dataauto$pheno$ffd=asin(sqrt(dataauto$pheno$ffd));dataauto$pheno$ffd_height=log(dataauto$pheno$ffd_height)
dataauto$pheno$open_flowers_ffd5=sqrt(dataauto$pheno$open_flowersffd5);dataauto$pheno$open_flowers_rate=log(dataauto$pheno$open_flowers_rate)
```

```
dataauto=jittermap(dataauto, amount=1e-6) # w/ X chromosome, but treated as autosome
data=jittermap(data, amount=1e-6)
data2=data
dataauto$pheno$gender=as.numeric(dataauto$pheno$gender)
dataauto$pheno$ffd[68] <- NA
dataauto$pheno$ffd=(dataauto$pheno$ffd-min(dataauto$pheno$ffd, na.rm=T))/(max(dataauto$pheno$ffd, na.rm=T)-min(dataauto$pheno$ffd, na.rm=T))
dataauto$pheno$ffd=asin(sqrt(dataauto$pheno$ffd))
dataauto$pheno$ffd_height=log(dataauto$pheno$ffd_height)
dataauto$pheno$open_flowers_ffd5=sqrt(dataauto$pheno$open_flowers_ffd5)
dataauto$pheno$open_flowers_rate=log(dataauto$pheno$open_flowers_rate)
```

```
# w/o X chromosome
datawox <- jittermap(datawox, amount = 1e-6)

# w/ X chromosome, but treated as an autosome
dataauto <- jittermap(dataauto, amount = 1e-6)
data <- jittermap(data, amount = 1e-6)
data2 <- data

# convert f to 0, m to 1
dataauto$pheno$gender <- as.numeric(dataauto$pheno$gender) - 1 # convert f to 0, m to 1
dataauto$pheno$ffd[68] <- NA
dataauto$pheno$ffd <- (dataauto$pheno$ffd - min(dataauto$pheno$ffd, na.rm = T)) / (max(dataauto$pheno$ffd, na.rm = T) - min(dataauto$pheno$ffd, na.rm = T))
dataauto$pheno$ffd <- asin(sqrt(dataauto$pheno$ffd))
dataauto$pheno$ffd_height <- log(dataauto$pheno$ffd_height)

dataauto$pheno$open_flowers_ffd5 <- sqrt(dataauto$pheno$open_flowers_ffd5)
dataauto$pheno$open_flowers_rate <- log(dataauto$pheno$open_flowers_rate)
```



The tidyverse style guide

Hadley Wickham

<https://style.tidyverse.org/index.html>

Google R-style

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



styler

The goal of styler is to provide non-invasive pretty-printing of R source code while adhering to the [tidyverse](#) formatting rules. styler can be customized to format code according to other style guides too.

Installation

You can install the package from CRAN:

```
install.packages("styler")
```

Some examples

```
# Good
day_one
day_1

# Bad
DayOne
dayone
```

```
# Good
x <- 5

# Bad
x = 5
```

```
# Good
if (y < 0) {
  stop("Y is negative")
}
```

```
# Bad
if (y < 0) stop("Y is negative")
```

```
# Good
"Text"
'Text with "quotes"'
'<a href="http://style.tidyverse.org">A link</a>'

# Bad
'Text'
'Text with "double" and \'single\' quotes'
```





save specific objects to a file

```
save(iris,file="iris.RData")
```

load object

```
load("iris.RData")
```

save workspace

```
save.image(file='image.RData')
```

load workspace

```
load("image.RData")
```


Customized functions

```
Myfunction <- function(variables) {
  Functions
}
```

```
plot_points <- function(n_points) {
  dat <- rnorm(n_points, 100, 5)
  plot(dat)
}
```

My_functions.R

```
1 source("My_functions.R")
2 plot_points(10000)
```

Shiny Apps

Shiny from R Studio

[Back to Gallery](#) [Get Cod](#)

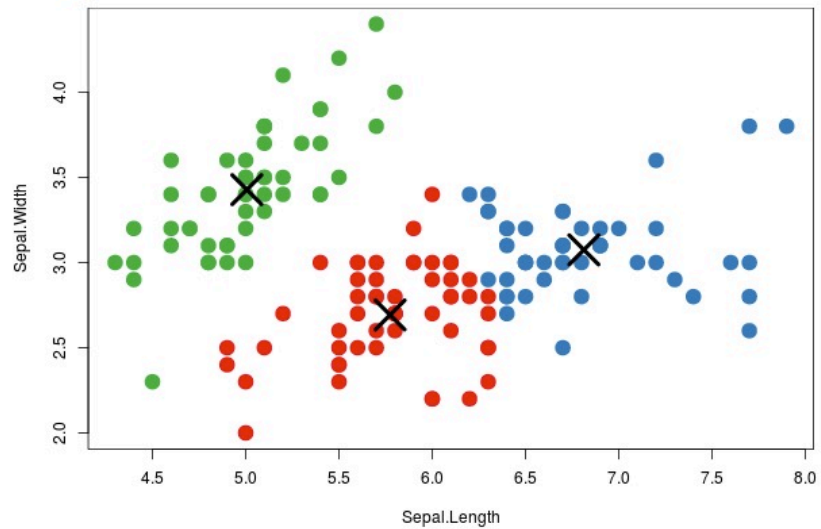


Iris k-means clustering

X Variable
Sepal.Length

Y Variable
Sepal.Width

Cluster count
3



server.R

ui.R

```
function(input, output, session) {  
  
  # Combine the selected variables into a new data frame  
  selectedData <- reactive({  
    iris[, c(input$xcol, input$ycol)]  
  })  
  
  clusters <- reactive({  
    kmeans(selectedData(), input$clusters)  
  })  
  
  output$plot1 <- renderPlot({  
    palette(c("#E41A1C", "#377EB8", "#4DAF4A", "#984EA3",  
             "#FF7F00", "#FFFF33", "#A65628", "#F781BF", "#999999"))  
  
    par(mar = c(5.1, 4.1, 0, 1))  
    plot(selectedData(),  
          col = clusters()$cluster,  
          pch = 20, cex = 3)  
    points(clusters()$centers, pch = 4, cex = 4, lwd = 4)  
  })  
}
```

<https://shiny.rstudio.com/gallery/kmeans-example.html>

Workflows on a HPC cluster

```

#!/usr/bin/env Rscript                                     Read.R
args <-  commandArgs(trailingOnly=TRUE)

##read table
samples <- read.table(args[1], header = F)
  
```

```

nik$ Rscript --vanilla Read.R table.csv|
  
```



- keep all scripts
- deposit raw data to public database
- use default settings or mention if not
- provide version information
- provide commands in supplement
- deposit scripts on github/gitlab
- Check-list

Comment

THE 'REAPPRAISED' CHECKLIST FOR EVALUATION OF PUBLICATION INTEGRITY

Not all items will be applicable to every publication, and other questions might be relevant for individual categories.

R — Research governance

- Are the locations where the research took place specified, and is this information plausible?
- Is a funding source reported?
- Has the study been registered?
- Are details such as dates and study methods in the publication consistent with those in the registration documents?

- 'P-hacking': biased or selective analyses that promote fragile results
- Other unacknowledged multiple statistical testing
- Is there outcome switching — that is, do the analysis and discussion focus on measures other than those specified in registered analysis plans?

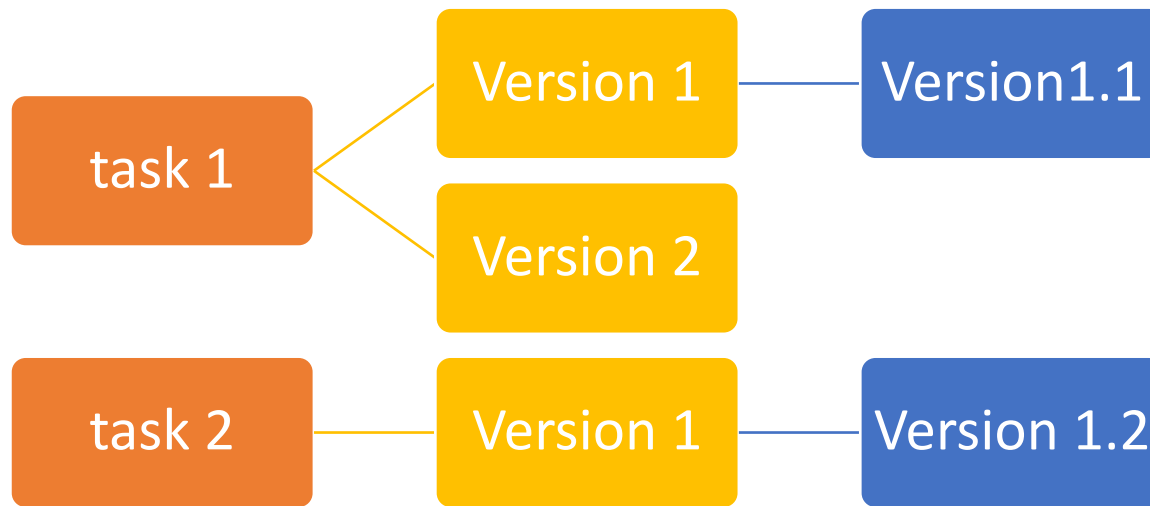
E — Ethics

- Is there evidence that the work has been approved by a specific, recognized committee?

I — Image manipulation

- Is there evidence of manipulation or duplication of images?

Exercise







Take home message

- Do reproducible research
- Rstudio has a powerful functions implemented

