

Useful open science tools to do research in psychiatry



Thomas Gargot
thomas_gargot@hotmail.com
EPA Congress,
Tuesday March 6th



Useful open science tools to do research in psychiatry



Randomizer.org



Studio[®]

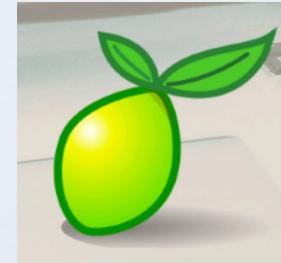
Analyse data



equator
network

Write an article

Collect data
Limesurvey



Share data
OSF and
GitHub



Make the knowledge
more visible

Goals of the workshop

- Understand the challenges of open science
- Learn the concepts of some tools
- Consider your own research/projects
- Rate cakes

Who I am

- Psychiatrist in Paris
- PhD Student on Computer Science :
use of IT in school



La pitié Salpêtrière



- Interest in IT and psychotherapy (CBT)
- Master in Cognitive science in ENS
- Former IT secretary of EFPT
- Psychotherapy WG chairman



Disclosure

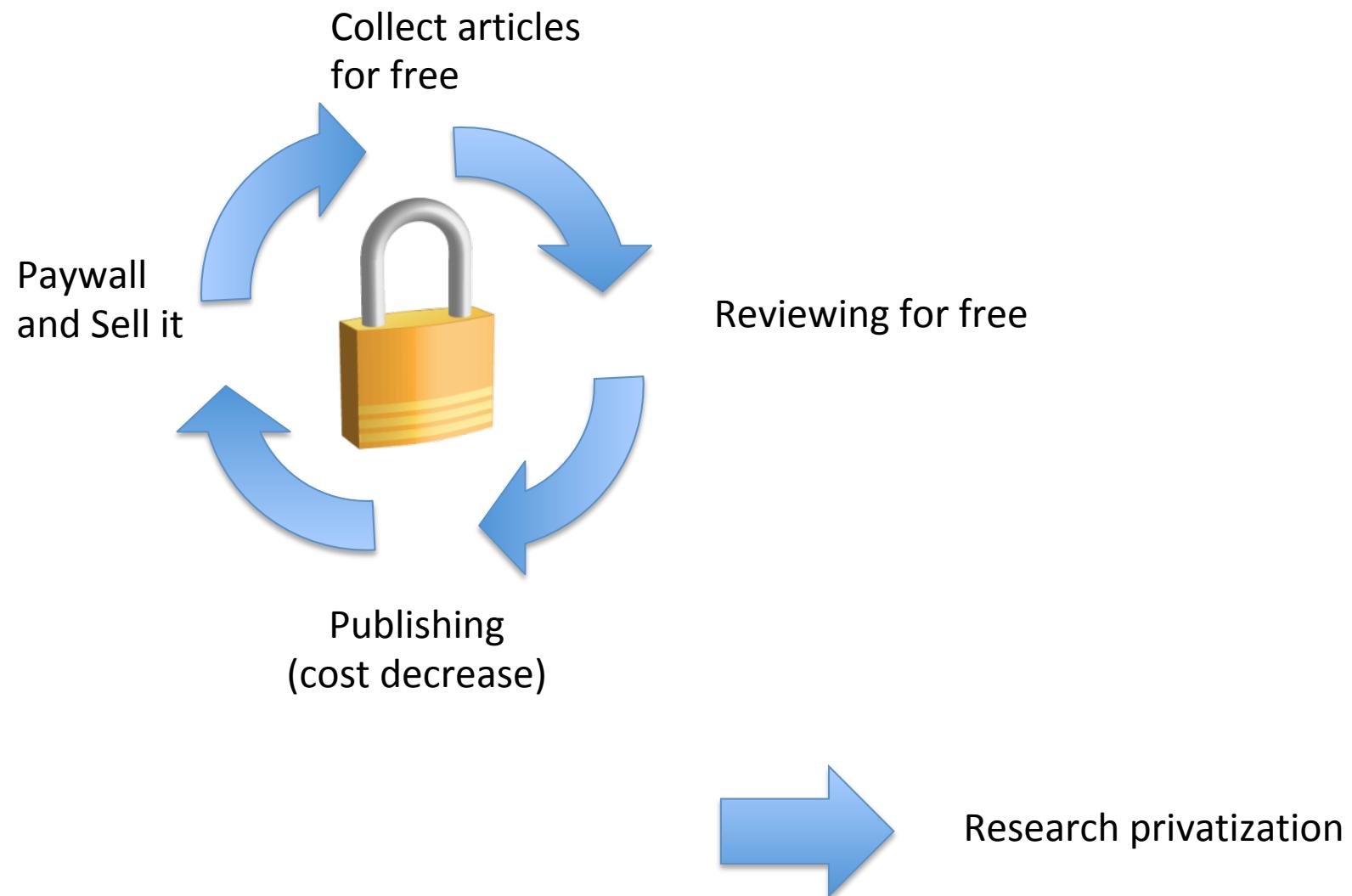
- None

Ethical issues

Problems

- **Cost:** How much will it cost to run this study?
- **Method:** What do they mean when they did this analysis?
- **Licence:** Does someone has the last version of this software ?
- **Memory:** Where is my data, how to re-use, how share it?
- **Paywall:** I can't get this article
- **Visibility:** Nobody cares about our work

Publication lucrative vicious circle





<https://www.youtube.com/watch?v=WnxqoP-c0ZE&list=PLKipY1cRnemK1Fn6PbVe5zkOfxLII02YY&index=6>

Goals

- Decrease costs
- Improve transparency
- Improve accessibility to softwares, articles, and re-use of data
- Improve visibility for public

→ Is it really possible ?!?

Why is there a need for open science?

nature
human behaviour

PERSPECTIVE

PUBLISHED: 10 JANUARY 2017 | VOLUME: 1 | ARTICLE NUMBER: 0021

OPEN

A manifesto for reproducible science

Marcus R. Munafò^{1,2*}, Brian A. Nosek^{3,4}, Dorothy V. M. Bishop⁵, Katherine S. Button⁶, Christopher D. Chambers⁷, Nathalie Percie du Sert⁸, Uri Simonsohn⁹, Eric-Jan Wagenmakers¹⁰, Jennifer J. Ware¹¹ and John P. A. Ioannidis^{12,13,14}

Data from many fields suggests reproducibility is lower than is desirable^{8–14}; one analysis estimates that 85% of biomedical research efforts are wasted¹⁴, while 90% of respondents to a recent survey in *Nature* agreed that there is a ‘reproducibility crisis’¹⁵.

<https://www.nature.com/articles/s41562-016-0021>

Information world

- Share information don't decrease the ownership of the sharer
- Rather, better accessibility improve possibilities of not anticipated use (hacking) and feedback
- Open softwares are more flexible and have fewer bugs than closed/proprietary one



Free / Open?



Open software: Free use, study, modification, duplication and sharing
≠ Freemium strategy (demonstration)
≠ centralisation and user data collection

An example of Open science pathway

Hands-on

Research protocol

- Make a data base of cookies tasting
- Rate them in an international context
- Select the best cakes in the world



Sampling process during EPA congress, Madrid, 2016

Useful open science tools to do research in psychiatry



Randomizer.org

Randomizer.org



RESEARCH RANDOMIZER

Randomizer.org



RESEARCH RANDOMIZER

RESULTS

PRINT

DOWNLOAD

CLOSE

1 Set of 10 Numbers

Range: From 1 to 2

Set #1

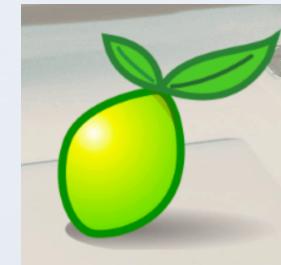
2, 2, 2, 1, 1, 2, 2, 1, 2, 1

Useful open science tools to do research in psychiatry



Randomizer.org

Collect data
Limesurvey



Limesurvey



LimeSurvey - the most popular
Free Open Source Software survey tool on the web.

Create professional online surveys
LimeService

"Truth fears no questions."



<http://thecakereports.limequery.com>

How does it look like ?

- 1
- 2
- 3
- 4
- 5
- No answer

How does it taste ?

- 1
- 2
- 3
- 4
- 5
- No answer

How is the texture ?

- 1
- 2
- 3
- 4
- 5
- No answer

How is the packaging ?

- 1
- 2
- 3
- 4
- 5
- No answer

Useful open science tools to do research in psychiatry



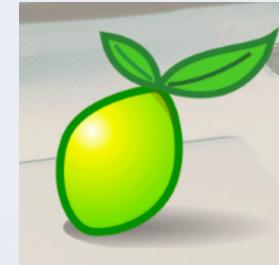
R Studio[®]

Analyse data



Randomizer.org

Collect data
Limesurvey





RStudio

File Edit Code View Project Workspace Plots Tools Help

Project: (None)

diamondPricing.R* formatPlot.R diamonds

Source on Save Go to file/function

1 Library(ggplot2)
2 source("plots/formatPlot.R")
3
4 View(diamonds)
5 summary(diamonds)
6
7 summary(diamonds\$price)
8 aveSize <- round(mean(diamonds\$carat), 4)
9 clarity <- levels(diamonds\$clarity)
10
11 p <- qplot(carat, price,
12 data=diamonds, color=clarity,
13 xlab="Carat", ylab="Price",
14 main="Diamond Pricing")
15

15:1 (Top Level) R Script

Console

```
      x          y          z
Min. : 0.000  Min. : 0.000  Min. : 0.000
1st Qu.: 4.710  1st Qu.: 4.720  1st Qu.: 2.910
Median : 5.700  Median : 5.710  Median : 3.530
Mean   : 5.731  Mean   : 5.735  Mean   : 3.539
3rd Qu.: 6.540  3rd Qu.: 6.540  3rd Qu.: 4.040
Max.  :10.740  Max.  :58.900  Max.  :31.800
> summary(diamonds$price)
  Min. 1st Qu. Median 3rd Qu.  Max.
  326    950   2401   3933   5324  18820
> aveSize <- round(mean(diamonds$carat), 4)
> clarity <- levels(diamonds$clarity)
> p <- qplot(carat, price,
+             data=diamonds, color=clarity,
+             xlab="Carat", ylab="Price",
+             main="Diamond Pricing")
>
> format.plot(p, size=24)
> |
```

Workspace History

Load Save Import Dataset Clear All

Data diamonds 53940 obs. of 10 variables

Values aveSize 0.7979 clarity character[8] p ggplot[8]

Functions format.plot(plot, size)

Plots Packages Help

Zoom Export Clear All

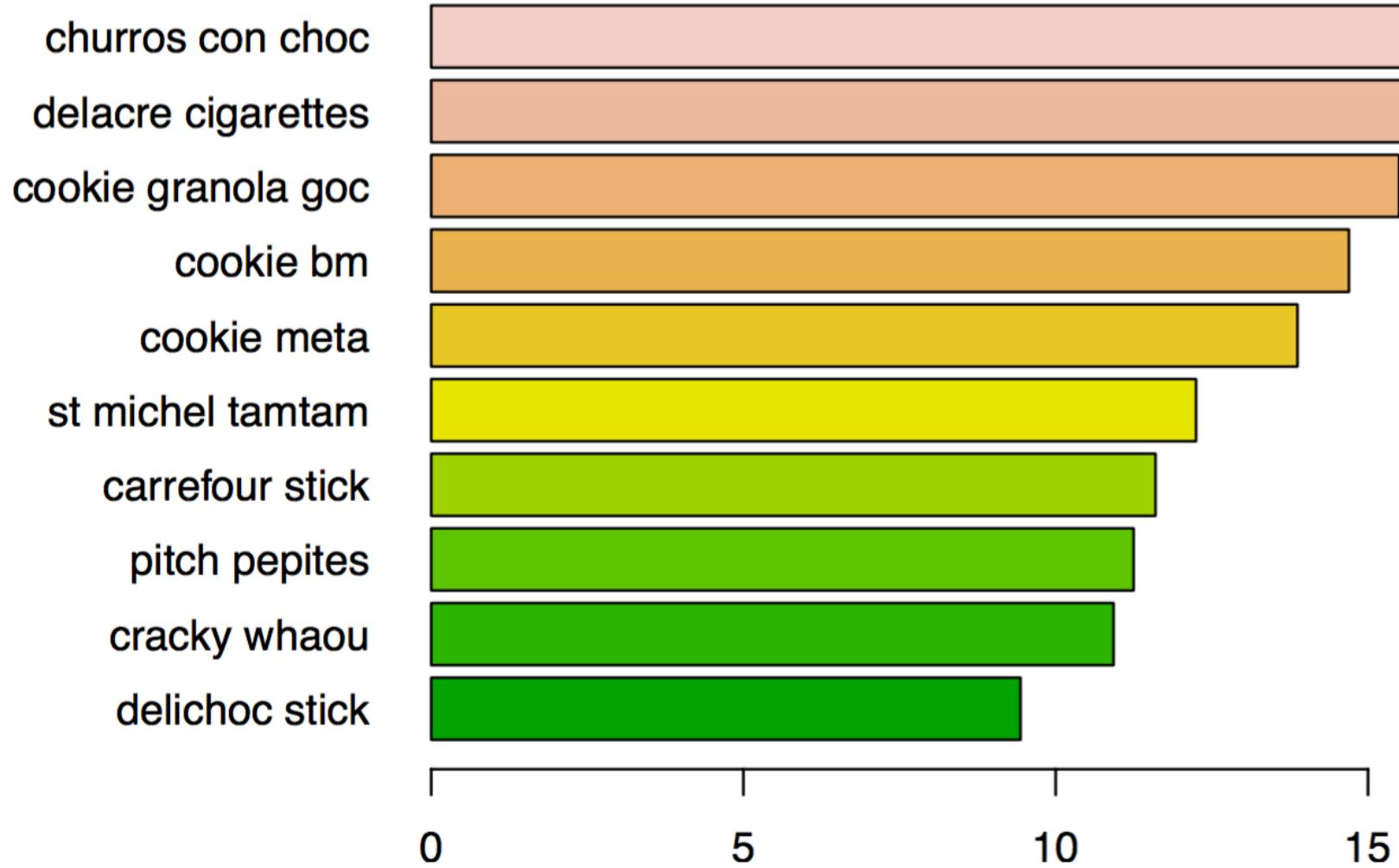
Diamond Pricing

A scatter plot titled "Diamond Pricing" showing the relationship between Carat (X-axis, ranging from 0.0 to 3.5) and Price (Y-axis, ranging from 0 to 15000). The data points are colored according to their Clarity level, as indicated by the legend on the right:

- I1 (Red)
- SI2 (Orange)
- SI1 (Yellow-green)
- VS2 (Green)
- VS1 (Blue)
- VVS2 (Light Blue)
- VVS1 (Purple)
- IF (Pink)

The plot shows a strong positive correlation between Carat and Price, with points clustered along the diagonal line where Price equals Carat. The distribution of points is more spread out at lower carat values and becomes more concentrated at higher carat values.

Total grades for each cake

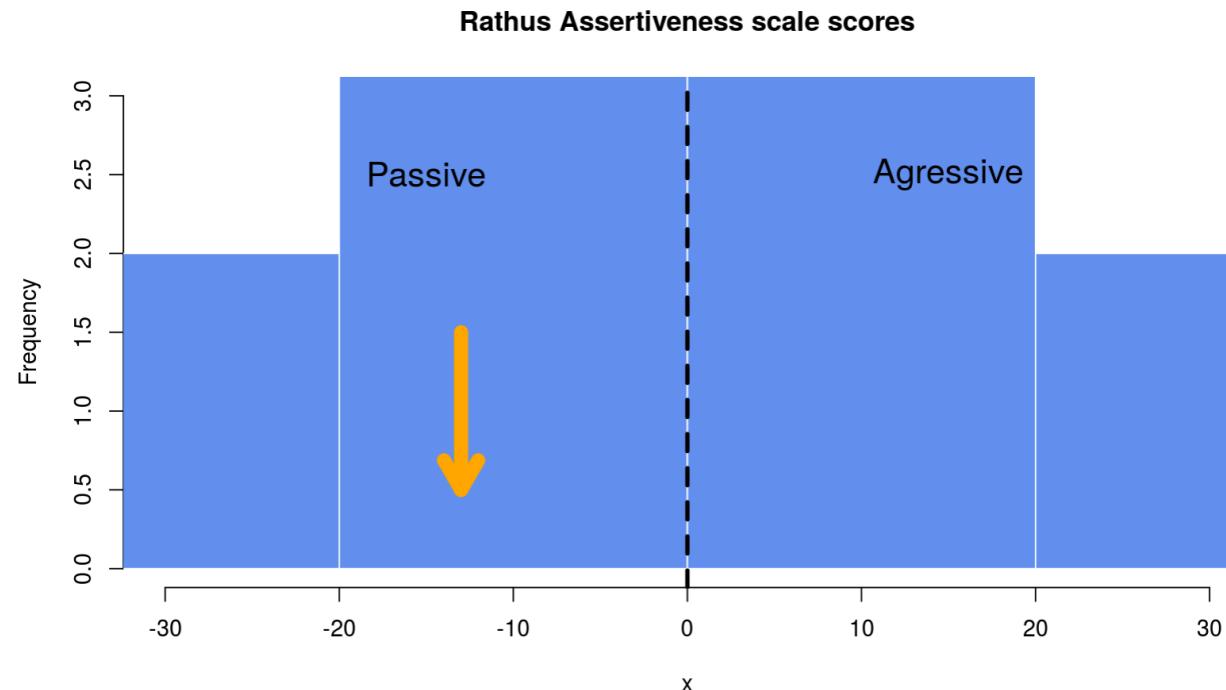


Connect Limesurvey with R ? : Shinyapps

Assertiveness Rathus Scale

What is your pseudo ?

Number of observations to view:



Your pseudo is thog1. Your score is -13. You are passive.

A score toward -90 is a sign than the subject has a lack of assertivity (passive style).

A score around 0 is found in normal assertive behaviour.

Python



- A little more fundamental and less user friendly, very used by computer scientists.
- Python has very large possibilities :
 - Collect data during a psychological experiment with pygame
 - Data analysis
 - Code a robot
 - Make a website



Useful open science tools to do research in psychiatry



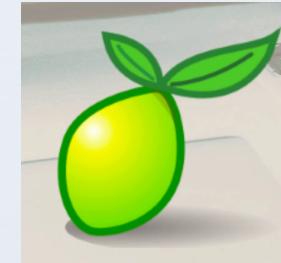
Studio[®]

Analyse data



Randomizer.org

Collect data
Limesurvey



Share data
OSF and
GitHub



Git Hub

- Share data

Screenshot of a GitHub repository page for "Ouphix / TCR".

The repository has the following statistics:

- 6 commits
- 1 branch
- 0 releases
- 2 contributors

Branch: master

Actions: New pull request, Create new file, Upload files, Find file, Clone or download

Latest commit: 174c9af on 5 Apr by AnneCha: yann vous sauve la vie

File	Type	Commit Message	Time
.RData	Rmarkdown v1	yann vous sauve la vie	7 months ago
.Rhistory	Rmarkdown v1	yann vous sauve la vie	7 months ago
README.md		yann vous sauve la vie	7 months ago
TheCakeReport.Rmd		yann vous sauve la vie	6 months ago
TheCakeReport.Rproj		yann vous sauve la vie	6 months ago
TheCakeReport.pdf		yann vous sauve la vie	6 months ago
carrefour stick.png		yann vous sauve la vie	6 months ago



Framagit, Gitlab

Useful open science tools to do research in psychiatry



Randomizer.org



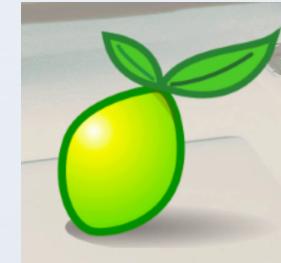
Studio[®]

Analyse data



Write an article

Collect data
Limesurvey



Share data
OSF and
GitHub



Equator Network



Enhancing the QUAlity and Transparency Of health Research



EQUATOR resources in
[Portuguese](#) | [Spanish](#)

[Home](#) [Library](#) [Toolkits](#) [Courses & events](#) [News](#) [Blog](#) [Librarian Network](#) [About us](#) [Contact](#)

Essential resources for writing and publishing health research



Library for health research reporting

The Library contains a comprehensive searchable database of reporting guidelines and also links to other resources relevant to research reporting.



[Search for reporting guidelines](#)



[Not sure which reporting guideline to use?](#)



[Reporting guidelines under development](#)



[Visit the library for](#)



Reporting guidelines for main study types

Randomised trials	CONSORT	Extensions	Other
Observational studies	STROBE	Extensions	Other
Systematic reviews	PRISMA	Extensions	Other
Case reports	CARE		Other
Qualitative research	SRQR	COREQ	Other
Diagnostic / prognostic studies	STARD	TRIPOD	Other
Quality improvement studies	SQUIRE		Other
Economic evaluations	CHEERS		Other
Animal pre-clinical studies	ARRIVE		Other
Study protocols	SPIRIT	PRISMA-P	Other

Publication School (2-day workshop)
The secrets of success in writing and publishing research articles



December 2016
Bond University, Gold Coast, QLD
AUSTRALIAN EQUATOR CENTRE

Do you want to get your health research published, and be praised for it?
Do you want your institution to be recognised for its excellent publication record?
Do you want to make a real difference with your research?
Then, this is the course for you!

A large number of published health research articles are reported badly. They provide insufficient, misleading or ambiguous information, and cannot be used to inform future research or improve healthcare for patients. It doesn't have to be this way.

Join us this summer for an intensive, practical course in the heart of the beautiful Gold Coast - and learn how to write a publishable research article in two days.

EQUATOR's flagship Publication School aims to develop essential writing skills to help you achieve success in planning, writing, publishing and communicating research through traditional journals and other channels.

Course tutors include:

Prof Paul Glasziou, As Prof Elaine Bellar (Bond University), and Dr David Moher (Canadian EQUATOR Centre)

For more details of the course and how to register, see www.rebelp.net.au or contact rebel@bond.edu.au

Collaborative writing

The figure shows the Overleaf interface. The left side displays the LaTeX source code with line numbers. The right side shows the generated PDF document.

LaTeX Source Code:

```
1 \documentclass[a4paper]{article}
2
3 %% Language and font encodings
4 \usepackage[english]{babel}
5 \usepackage[utf8x]{inputenc}
6 \usepackage[T1]{fontenc}
7
8 %% Sets page size and margins
9 \usepackage[a4paper,top=3cm,bottom=2cm,left=3cm,right=3cm,marginparwidth=1.75cm]{geometry}
{geometry}
10
11 %% Useful packages
12 \usepackage{amsmath}
13 \usepackage{graphicx}
14 \usepackage{colorinlistoftodos}{todonotes}
15 \usepackage[colorlinks=true, allcolors=blue]{hyperref}
16
17 \title{Your Paper}
18 \author{You}
19
20 \begin{document}
21 \maketitle
22
```

Generated PDF Document:

Your Paper
You
December 14, 2017

Abstract
Your abstract.

1 Introduction

Your introduction goes here! Some examples of commonly used commands and features are listed below, to help you get started. If you have a question, please use the help menu ("?") on the top bar to search for help or ask us a question.

2 Some examples to get started

2.1 How to add Comments

Comments can be added to your project by clicking on the comment icon in the toolbar above. To reply to a comment, simply click the reply button in the lower right corner of the comment, and you can close them when you're done.

2.2 How to include Figures

First you have to upload the image file from your computer using the upload link the project menu. Then use the `\includegraphics` command to include it in your document. Use the `\figure` environment and the `\caption` command to add a number and a caption to your figure. See the code for Figure 1 in this section for an example.

2.3 How to add Tables

Use the `\table` and `\tabular` commands for basic tables — see Table 1, for example.



Figure 1: This frog was uploaded via the project menu.



Welcome to StackEdit!

Hi! I'm your first Markdown file in **StackEdit**. If you want to learn about StackEdit, you can read me. If you want to play with Markdown, you can edit me. If you have finished with me, you can just create new files by opening the **file explorer** on left corner of the navigation bar.

Files

StackEdit stores your files in your browser, which means all your files are automatically saved locally and are accessible **offline!**

> **Note:**

- >
- > - StackEdit can be used offline thanks to the application cache.
- > - Your local files are not shared between different browsers or computers unless you use the [\[synchronization mechanism\]](#)

Welcome to StackEdit!

Hi! I'm your first Markdown file in **StackEdit**. If you want to learn about StackEdit, you can read me. If you want to play with Markdown, you can edit me. If you have finished with me, you can just create new files by opening the **file explorer** on left corner of the navigation bar.

Files

StackEdit stores your files in your browser, which means all your files are automatically saved locally and are accessible **offline!**

Note:

- StackEdit can be used offline thanks to the application cache.
- Your local files are not shared between different browsers or



- Select, gather, write and share bibliography

The screenshot shows the Zotero application window. On the left is a sidebar with a tree view of library collections and group libraries. The main area displays a bibliography table with three entries:

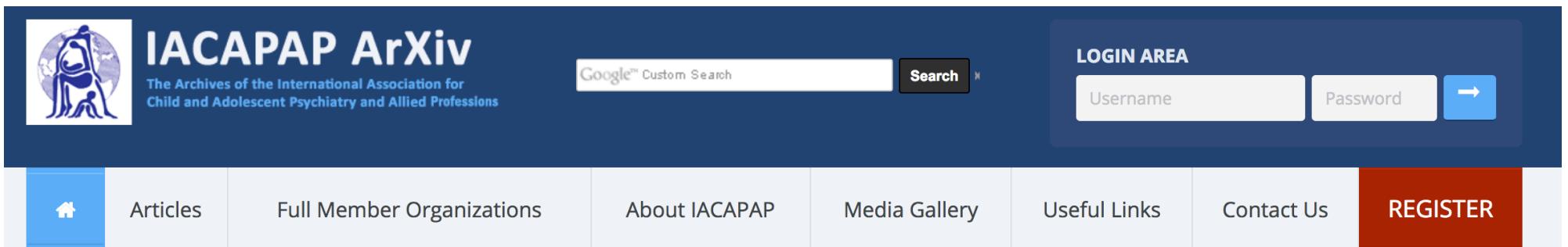
Titre	Créateur
Handwriting or Typewriting? Th...	Kiefer et al.
The influence of writing practice...	Longcamp et al.
Comparing Memory for Handwri...	Smoker et al.

On the right, the details for the first entry are expanded:

Type de document Article de revue
Titre Handwriting or Typewriting? The Influence of Pen- or Keyboard-Based Writing Training on Reading and Writing Performance in Preschool Children
Auteur Kiefer, Markus
Auteur Schuler, Stefanie
Auteur Mayer, Carmen
Auteur Trumpp, Natalie M.
Auteur Hille, Katrin
Auteur Sachse, Steffi
Résumé Digital writing devices associated with the use of computers, tablet PCs, or mobile phones are increasingly replacing writing by hand. It is, however, controversially discussed how writing modes influence reading and writing performance in children at the start of literacy. On the one hand, the easiness of typing on digital devices may accelerate reading and writing in young children, who have less developed sensory-motor skills. On the other hand, the meaningful coupling between action and perception during handwriting, which establishes sensory-motor memory traces, could facilitate written language acquisition. In order to decide between these theoretical alternatives, for the present study, we developed an intense training program for preschool children attending the German kindergarten with 16 training sessions. Using closely matched letter learning games, eight

New ways to publish ?

ArXiv : Preprints



The screenshot shows the IACAPAP ArXiv website. At the top left is the logo featuring a stylized figure holding a globe. To its right, the text "IACAPAP ArXiv" is displayed in large, bold, white letters, with "The Archives of the International Association for Child and Adolescent Psychiatry and Allied Professions" in smaller blue text below it. A search bar with a "Google Custom Search" placeholder and a "Search" button are positioned next to the logo. On the right side, there is a "LOGIN AREA" with "Username" and "Password" fields and a blue "→" button. Below the header is a navigation bar with tabs: "Articles" (highlighted in blue), "Full Member Organizations", "About IACAPAP", "Media Gallery", "Useful Links", "Contact Us", and a red "REGISTER" button.



About IACAPAP ArXiv

"IACAPAP ArXiv" is a facility that gives child and adolescent mental health professionals the opportunity to upload clinical or research documents in their own language (with an abstract in English). These documents are then freely available to all internet users. The quality of the articles is assured by formal approval by the national child and adolescent mental health organization to which the authors belong, or by IACAPAP if the national organization does not exist or does not participate to the ArXiv.



SUBMIT A MANUSCRIPT
Submit your paper to IACAPAP



FOLLOW IACAPAP
Follow Us on Social Media

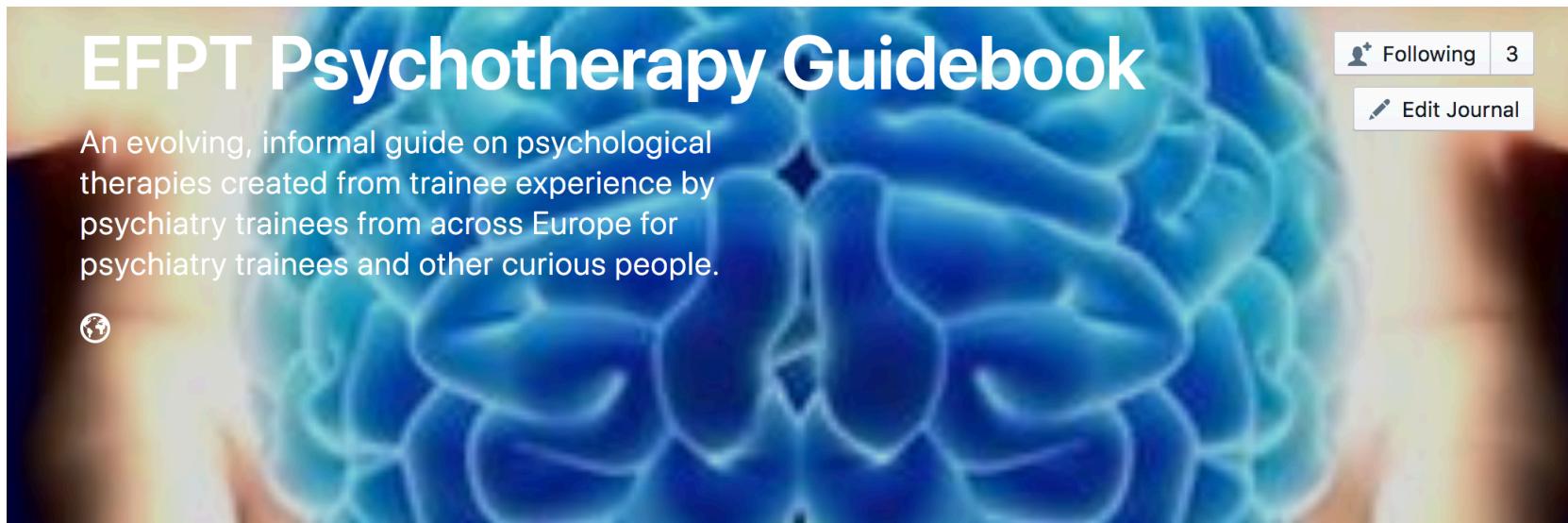


SUPPORT
Web support information pages

bioRxiv
beta



<https://www.pubpub.org/>



Featured

Submitted

Pages

People

Sort ▾



Rational Emotive Behaviour Therapy

EFPT Psychotherapy Guidebook chapter on Rational Emotive Behaviour Therapy

Featured on June 18, 2017

Pages +



Psychoeducation

EFPT Psychotherapy Guidebook chapter on Psychoeducation

Featured on June 18, 2017

• Authorea

Several way of practicing Open Access

- Not only Open Access Gold : journal where you need to pay to publish (new publisher's strategy)
- Use Open Repository as HAL in France, Zenodo, etc.
- Don't forget the licences (Creative Commons, etc)

LICENSES



TERMS



Attribution

Others can copy, distribute, display, perform and remix your work if they credit your name as requested by you



No Derivative Works

Others can only copy, distribute, display or perform verbatim copies of your work



Share Alike

Others can distribute your work only under a license identical to the one you have chosen for your work

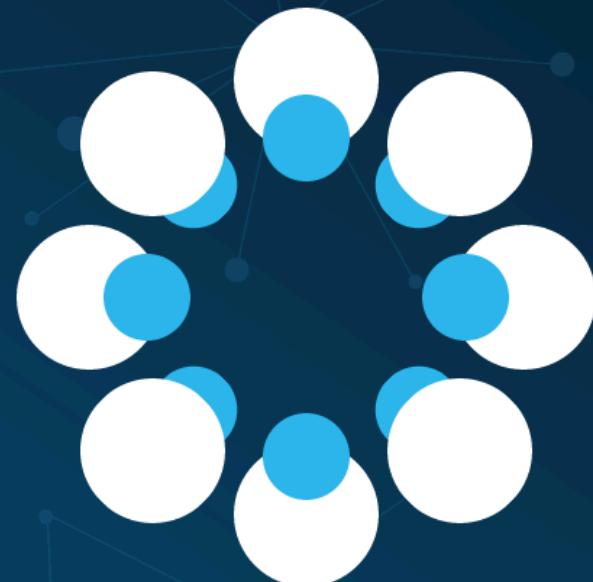


Non-Commercial

Others can copy, distribute, display, perform or remix your work but for non-commercial purposes only.

Open Science Framework

A scholarly commons to connect the entire research cycle



Share with general public

Useful open science tools to do research in psychiatry



Randomizer.org



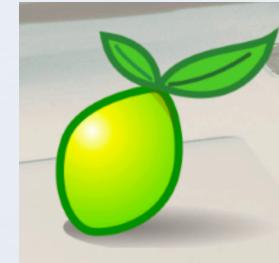
Studio[®]

Analyse data

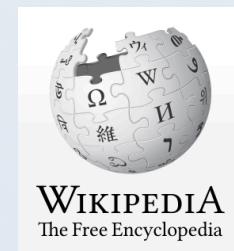
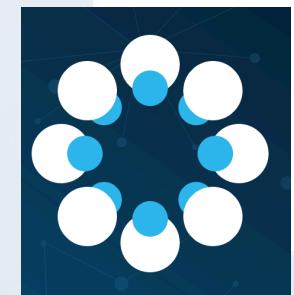


Write an article

Collect data
Limesurvey



Share data
OSF and
GitHub



**Make the knowledge
more visible**



WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia
Wikipedia store

Interaction
Help
About Wikipedia
Community portal
Recent changes
Contact page

Tools
What links here
Related changes
Upload file
Special pages
Permanent link
Page information

Article Talk

Read Edit View history

Not logged in Talk Contributions Create account Log in

Biscuit

From Wikipedia, the free encyclopedia

See also: [Biscuit \(bread\)](#) and [Cookie](#)

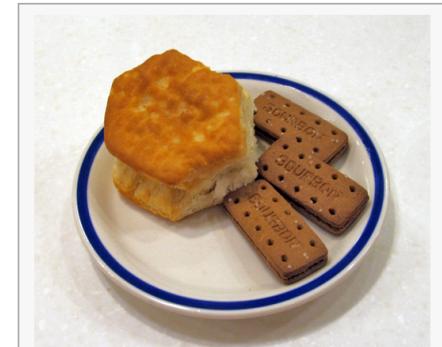
For other uses, see [Biscuit \(disambiguation\)](#).



This article may lack focus or may be about more than one topic. Please help improve this article, possibly by splitting the article and/or by introducing a [disambiguation page](#), or discuss this issue on the [talk page](#). (May 2016)

Biscuit is a term used for a diverse variety of [baked](#), commonly [flour-based](#) food products. The term is applied to two distinct products in [North America](#) and the [Commonwealth of Nations](#) and [Europe](#). The North American biscuit is typically a soft, leavened [quick bread](#), and is covered in the article [Biscuit \(bread\)](#). This article covers the other type of biscuit, which is typically hard, flat and unleavened.

Biscuit



American biscuit (left) and one variety of British biscuit (right) – the American biscuit is

Contents [hide]

- 1 Variations in meaning
- 2 Etymology
- 3 History
 - 3.1 Biscuits for travel
 - 3.2 Confectionery biscuits
- 4 Biscuits today
 - 4.1 Commonwealth of Nations and Europe

<https://en.wikipedia.org/>



WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia
Wikipedia store

Interaction
Help
About Wikipedia
Community portal
Recent changes
Contact page

Tools
What links here
Related changes
Upload file
Special pages
Permanent link
Page information

Article Talk

Read Edit View history

Search



Not logged in Talk Contributions Create account Log in

Biscuit

From Wikipedia, the free encyclopedia

See also: [Biscuit \(bread\)](#) and [Cookie](#)

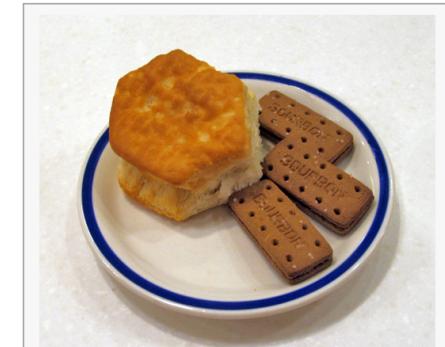
For other uses, see [Biscuit \(disambiguation\)](#).



This article may lack focus or may be about more than one topic. Please help improve this article, possibly by splitting the article and/or by introducing a [disambiguation page](#), or discuss this issue on the [talk page](#). (May 2016)

Biscuit is a term used for a diverse variety of [baked](#), commonly [flour-based](#) food products. The term is applied to two distinct products in [North America](#) and the [Commonwealth of Nations](#) and [Europe](#). The North American biscuit is typically a soft, leavened [quick bread](#), and is covered in the article [Biscuit \(bread\)](#). This article covers the other type of biscuit, which is typically hard, flat and unleavened.

Biscuit



American biscuit (left) and one variety of British biscuit (right) – the American biscuit is

Contents [hide]

- 1 Variations in meaning
- 2 Etymology
- 3 History
 - 3.1 Biscuits for travel
 - 3.2 Confectionery biscuits
- 4 Biscuits today
 - 4.1 Commonwealth of Nations and Europe

<https://en.wikipedia.org/>



WORDPRESS.ORG

Showcase Themes Plugins Mobile Support Get Involved About Blog Hosting

Search WordPress.org



Download WordPress

Meet WordPress

WordPress is open source software you can use to create a beautiful website, blog, or app.

The screenshot shows the WordPress admin dashboard with a blue header bar. The main content area is titled "Add a New Post". A sidebar on the left lists navigation items: Dashboard, Jetpack, Posts (which is selected and highlighted in blue), All Posts, Add New, Categories, Tags, Media, Pages, Comments, Events, and Feedback. A message box in the center states: "Wordfence could not get an API key from the Wordfence scanning servers when it activated. You can try to fix this by going to the Wordfence "options" page and hitting "Save Changes". This will cause Wordfence to retry fetching an API key for you. If you keep seeing this error it usually means your WordPress server can't connect to our scanning servers. You can try asking your WordPress host to allow your WordPress server to connect to noc1.wordfence.com." Below this message is a title input field with the placeholder "Enter title here". Underneath the title field is a toolbar with buttons for "Add Media" and "Add Contact Form", and tabs for "Visual" and "Text". Below the toolbar is a rich text editor toolbar with buttons for bold, italic, link, b-quote, del, ins, img, ul, ol, li, code, more, close tags, contact form, and proofread. To the right of the post editor is a "Publish" sidebar containing "Save Draft" and "Preview" buttons, and sections for "Status: Draft" (with an edit link), "Visibility: Public" (with an edit link), and "Publish immediately" (with an edit link). At the bottom of the sidebar is a "Publicize: Not Connected" section with a "Show" button. A "Publish" button is located at the very bottom right of the sidebar.

Badges

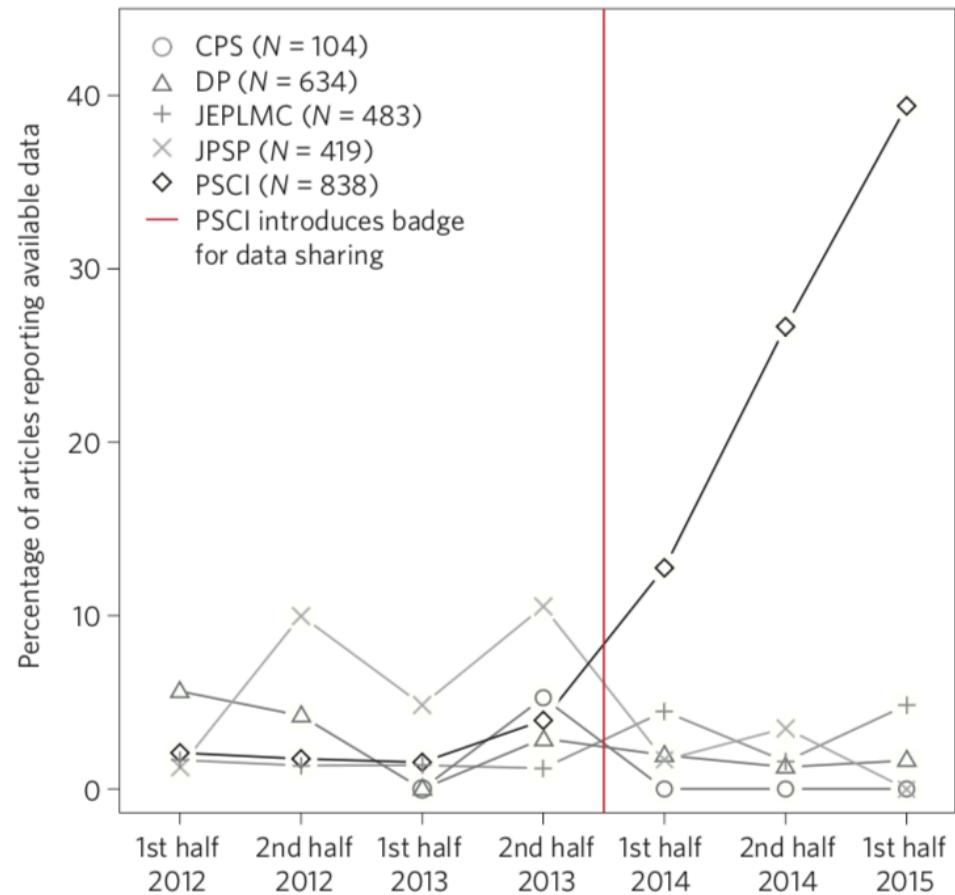


Figure 2 | The impact of introducing badges for data sharing. In January 2014, the journal *Psychological Science* (PSCI) introduced badges for articles with open data. Immediately afterwards, the proportion of articles with open data increased steeply, and by October 2015, 38% of articles in *Psychological Science* had open data. For comparison journals (*Clinical Psychological Science* (CPS), *Developmental Psychology* (DP), *Journal of Experimental Psychology: Learning, Memory and Cognition* (JEPLMC) and *Journal of Personality and Social Psychology* (JPSP)) the proportion of articles with open data remained uniformly low. Figure adapted from ref. 75, PLoS.

Useful open science tools to do research in psychiatry



Randomizer.org



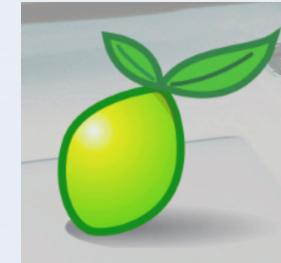
Studio[®]

Analyse data

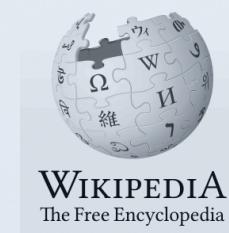


Write an article

Collect data :
Limesurvey



Share data
OSF and
GitHub



WIKIPEDIA
The Free Encyclopedia

Make the knowledge
more visible



Thank you for your attention



Hack your phd

<http://hackyourphd.org/>

thomas_gargot@hotmail.com

Base R Cheat Sheet

Getting Help

Accessing the help files

?mean

Get help of a particular function.

help.search('weighted mean')

Search the help files for a word or phrase.

help(package = 'dplyr')

Find help for a package.

More about an object

str(iris)

Get a summary of an object's structure.

class(iris)

Find the class an object belongs to.

Using Packages

install.packages('dplyr')

Download and install a package from CRAN.

library(dplyr)

Load the package into the session, making all its functions available to use.

dplyr::select

Use a particular function from a package.

data(iris)

Load a built-in dataset into the environment.

Working Directory

getwd()

Find the current working directory (where inputs are found and outputs are sent).

setwd('C://file/path')

Change the current working directory.

Use projects in RStudio to set the working directory to the folder you are working in.

Vectors

Creating Vectors

c(2, 4, 6)	2 4 6	Join elements into a vector
2:6	2 3 4 5 6	An integer sequence
seq(2, 3, by=0.5)	2.0 2.5 3.0	A complex sequence
rep(1:2, times=3)	1 2 1 2 1 2	Repeat a vector
rep(1:2, each=3)	1 1 1 2 2 2	Repeat elements of a vector

Vector Functions

sort(x)	rev(x)	Return x sorted.
table(x)	unique(x)	See counts of values.

Selecting Vector Elements

By Position

x[4] The fourth element.

x[-4] All but the fourth.

x[2:4] Elements two to four.

x[-(2:4)] All elements except two to four.

x[c(1, 5)] Elements one and five.

By Value

x[x == 10] Elements which are equal to 10.

x[x < 0] All elements less than zero.

x[x %in% c(1, 2, 5)] Elements in the set 1, 2, 5.

Named Vectors

x['apple'] Element with name 'apple'.

Programming

For Loop

```
for (variable in sequence){  
  Do something  
}
```

Example

```
for (i in 1:4){  
  j <- i + 10  
  print(j)  
}
```

While Loop

```
while (condition){  
  Do something  
}
```

Example

```
while (i < 5){  
  print(i)  
  i <- i + 1  
}
```

Functions

```
function_name <- function(var){  
  Do something  
  return(new_variable)  
}
```

Example

```
square <- function(x){  
  squared <- x*x  
  return(squared)  
}
```

Reading and Writing Data

Also see the `readr` package.

Input	Output	Description
df <- read.table('file.txt')	write.table(df, 'file.txt')	Read and write a delimited text file.
df <- read.csv('file.csv')	write.csv(df, 'file.csv')	Read and write a comma separated value file. This is a special case of read.table/write.table.
load('file.RData')	save(df, file = 'file.Rdata')	Read and write an R data file, a file type special for R.

Conditions

a == b	Are equal	a > b	Greater than	a >= b	Greater than or equal to	is.na(a)	Is missing
a != b	Not equal	a < b	Less than	a <= b	Less than or equal to	is.null(a)	Is null

Types

Converting between common data types in R. Can always go from a higher value in the table to a lower value.		
as.logical	TRUE, FALSE, TRUE	Boolean values (TRUE or FALSE).
as.numeric	1, 0, 1	Integers or floating point numbers.
as.character	'1', '0', '1'	Character strings. Generally preferred to factors.
as.factor	'1', '0', '1', levels: '1', '0'	Character strings with preset levels. Needed for some statistical models.

Maths Functions

log(x)	Natural log.	sum(x)	Sum.
exp(x)	Exponential.	mean(x)	Mean.
max(x)	Largest element.	median(x)	Median.
min(x)	Smallest element.	quantile(x)	Percentage quantiles.
round(x, n)	Round to n decimal places.	rank(x)	Rank of elements.
signif(x, n)	Round to n significant figures.	var(x)	The variance.
cor(x, y)	Correlation.	sd(x)	The standard deviation.

Variable Assignment

```
> a <- 'apple'
> a
[1] 'apple'
```

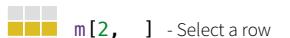
The Environment

ls()	List all variables in the environment.
rm(x)	Remove x from the environment.
rm(list = ls())	Remove all variables from the environment.

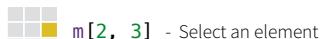
You can use the environment panel in RStudio to browse variables in your environment.

Matrices

`m <- matrix(x, nrow = 3, ncol = 3)`
Create a matrix from x.



`m[, 1]` - Select a column



`t(m)`

Transpose

`m %*% n`

Matrix Multiplication

`solve(m, n)`

Find x in: $m^* x = n$

Lists

`l <- list(x = 1:5, y = c('a', 'b'))`

A list is a collection of elements which can be of different types.

`l[[2]]`

Second element of l.

`l[1]`

New list with only the first element.

`l$x`

Element named x.

`l['y']`

New list with only element named y.

Also see the `dplyr` package.

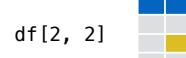
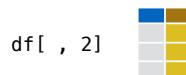
Data Frames

`df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))`

A special case of a list where all elements are the same length.

x	y
1	a
2	b
3	c

Matrix subsetting



List subsetting



Understanding a data frame

`View(df)` See the full data frame.

`head(df)` See the first 6 rows.

`nrow(df)` Number of rows.

`ncol(df)` Number of columns.

`dim(df)` Number of columns and rows.

`cbind` - Bind columns.

`rbind` - Bind rows.

Strings

Also see the `stringr` package.

`paste(x, y, sep = ' ')`

Join multiple vectors together.

`paste(x, collapse = ' ')`

Join elements of a vector together.

`grep(pattern, x)`

Find regular expression matches in x.

`gsub(pattern, replace, x)`

Replace matches in x with a string.

`toupper(x)`

Convert to uppercase.

`tolower(x)`

Convert to lowercase.

`nchar(x)`

Number of characters in a string.

Factors

`factor(x)`

Turn a vector into a factor. Can set the levels of the factor and the order.

`cut(x, breaks = 4)`

Turn a numeric vector into a factor by 'cutting' into sections.

Statistics

`lm(y ~ x, data=df)`

Linear model.

`glm(y ~ x, data=df)`

Generalised linear model.

`summary`

Get more detailed information out a model.

`t.test(x, y)`

Perform a t-test for difference between means.

`pairwise.t.test`

Perform a t-test for paired data.

`aov`

Analysis of variance.

Distributions

	Random Variates	Density Function	Cumulative Distribution	Quantile
Normal	<code>rnorm</code>	<code>dnorm</code>	<code>pnorm</code>	<code>qnorm</code>
Poisson	<code>rpois</code>	<code>dpois</code>	<code>ppois</code>	<code>qpois</code>
Binomial	<code>rbinom</code>	<code>dbinom</code>	<code>pbinom</code>	<code>qbinom</code>
Uniform	<code>runif</code>	<code>dunif</code>	<code>punif</code>	<code>qunif</code>

Plotting

Also see the `ggplot2` package.

`plot(x)`

Values of x in order.

`plot(x, y)`

Values of x against y.

`hist(x)`

Histogram of x.

Dates

See the `lubridate` package.

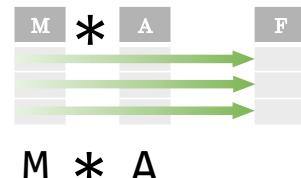
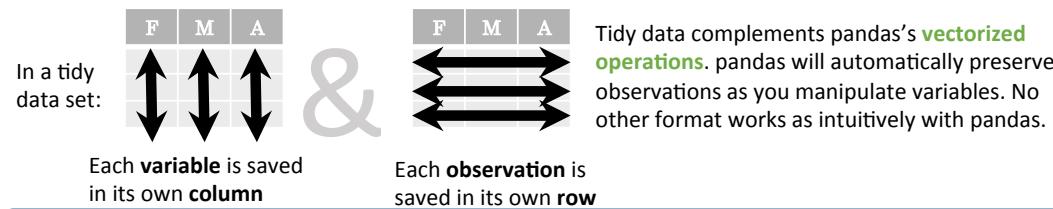
Data Wrangling

with pandas

Cheat Sheet

<http://pandas.pydata.org>

Tidy Data – A foundation for wrangling in pandas



Syntax – Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.
```

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

	a	b	c
n	v		
d	1	4	7
e	2	5	11
	2	6	9
			12

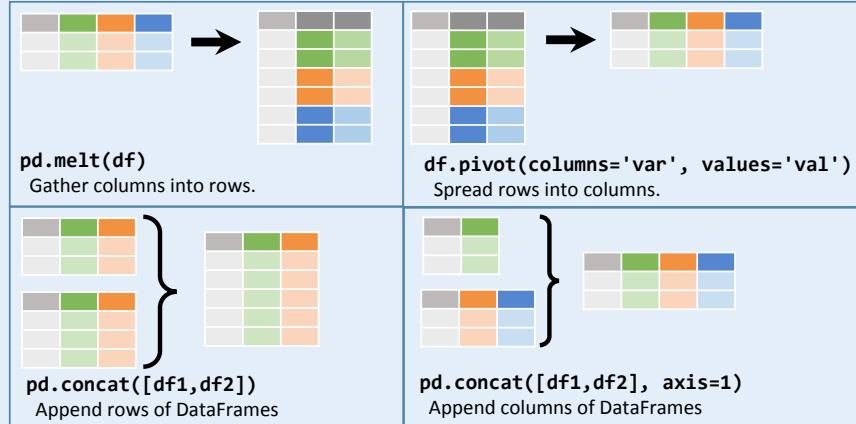
```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=['n','v']))
Create DataFrame with a MultiIndex
```

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={
          'variable' : 'var',
          'value' : 'val'})
      .query('val >= 200')
)
```

Reshaping Data – Change the layout of a data set



```
df.sort_values('mpg')
Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)
Order rows by values of a column (high to low).

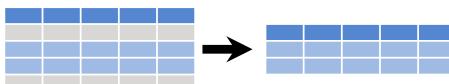
df.rename(columns = {'y':'year'})
Rename the columns of a DataFrame

df.sort_index()
Sort the index of a DataFrame

df.reset_index()
Reset index of DataFrame to row numbers, moving index to columns.

df.drop(['Length','Height'], axis=1)
Drop columns from DataFrame
```

Subset Observations (Rows)



```
df[df.Length > 7]
Extract rows that meet logical criteria.

df.drop_duplicates()
Remove duplicate rows (only considers columns).

df.head(n)
Select first n rows.

df.tail(n)
Select last n rows.

df.sample(frac=0.5)
Randomly select fraction of rows.

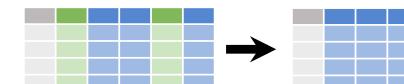
df.sample(n=10)
Randomly select n rows.

df.iloc[10:20]
Select rows by position.

df.nlargest(n, 'value')
Select and order top n entries.

df.nsmallest(n, 'value')
Select and order bottom n entries.
```

Subset Variables (Columns)



```
df[['width', 'length', 'species']]
Select multiple columns with specific names.

df['width'] or df.width
Select single column with specific name.

df.filter(regex='regex')
Select columns whose name matches regular expression regex.
```

regex (Regular Expressions) Examples

'.'	Matches strings containing a period '.'
'Length\$'	Matches strings ending with word 'Length'
'^Sepal'	Matches strings beginning with the word 'Sepal'
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^(?!Species\$).*'	Matches strings except the string 'Species'

```
df.loc[:, 'x2':'x4']
Select all columns between x2 and x4 (inclusive).

df.iloc[:,1,2,5]
Select columns in positions 1, 2 and 5 (first column is 0).

df.loc[df['a'] > 10, ['a','c']]
Select rows meeting logical condition, and only the specific columns .
```

Logic in Python (and pandas)

<	Less than	!=	Not equal to
>	Greater than	df.column.isin(values)	Group membership
==	Equals	pd.isnull(obj)	Is NaN
<=	Less than or equals	pd.notnull(obj)	Is not NaN
>=	Greater than or equals	&, , ~, ^, df.any(), df.all()	Logical and, or, not, xor, any, all

Summarize Data

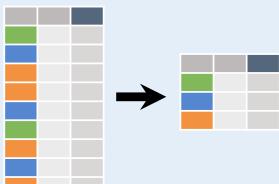
```
df['w'].value_counts()
Count number of rows with each unique value of variable
len(df)
# of rows in DataFrame.
df['w'].nunique()
# of distinct values in a column.
df.describe()
Basic descriptive statistics for each column (or GroupBy)
```



pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

sum()	Sum values of each object.	min()	Minimum value in each object.
count()	Count non-NA/null values of each object.	max()	Maximum value in each object.
median()	Median value of each object.	mean()	Mean value of each object.
quantile([0.25, 0.75])	Quantiles of each object.	var()	Variance of each object.
std()	Standard deviation of each object.	apply(function)	Apply function to each object.

Group Data



df.groupby(by="col")
Return a GroupBy object, grouped by values in column named "col".

df.groupby(level="ind")
Return a GroupBy object, grouped by values in index level named "ind".

All of the summary functions listed above can be applied to a group. Additional GroupBy functions:

size()
Size of each group.

agg(function)
Aggregate group using function.

Windows

```
df.expanding()
Return an Expanding object allowing summary functions to be applied cumulatively.
df.rolling(n)
Return a Rolling object allowing summary functions to be applied to windows of length n.
```

Handling Missing Data

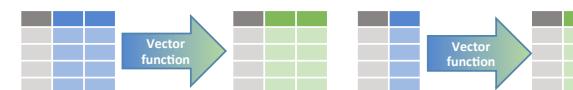
```
df.dropna()
Drop rows with any column having NA/null data.
df.fillna(value)
Replace all NA/null data with value.
```

Make New Columns



```
df.assign(Area=lambda df: df.Length*df.Height)
Compute and append one or more new columns.
df['Volume'] = df.Length*df.Height*df.Depth
Add single column.
```

pd.qcut(df.col, n, labels=False)
Bin column into n buckets.



pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

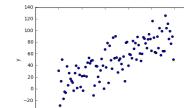
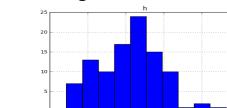
max(axis=1)	Element-wise max.	min(axis=1)	Element-wise min.
clip(lower=-10,upper=10)	Trim values at input thresholds	abs()	Absolute value.

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

shift(1)	Copy with values shifted by 1.	shift(-1)	Copy with values lagged by 1.
rank(method='dense')	Ranks with no gaps.	cumsum()	Cumulative sum.
rank(method='min')	Ranks. Ties get min rank.	cummax()	Cumulative max.
rank(pct=True)	Ranks rescaled to interval [0, 1].	cummin()	Cumulative min.
rank(method='first')	Ranks. Ties go to first value.	cumprod()	Cumulative product.

Plotting

```
df.plot.hist()
Histogram for each column
df.plot.scatter(x='w',y='h')
Scatter chart using pairs of points
```



Combine Data Sets

adf	bdf
x1 x2	x1 x3
A 1	A T
B 2	B F
C 3	D T

Standard Joins

```
pd.merge(adf, bdf,
        how='left', on='x1')
Join matching rows from bdf to adf.
```

```
pd.merge(adf, bdf,
        how='right', on='x1')
Join matching rows from adf to bdf.
```

```
pd.merge(adf, bdf,
        how='inner', on='x1')
Join data. Retain only rows in both sets.
```

```
pd.merge(adf, bdf,
        how='outer', on='x1')
Join data. Retain all values, all rows.
```

Filtering Joins

```
adf[adf.x1.isin(bdf.x1)]
All rows in adf that have a match in bdf.
```

```
adf[~adf.x1.isin(bdf.x1)]
All rows in adf that do not have a match in bdf.
```

ydf	zdf
x1 x2	x1 x2
A 1	B 2
B 2	C 3
C 3	D 4

Set-like Operations

```
pd.merge(ydf, zdf)
Rows that appear in both ydf and zdf (Intersection).
```

```
pd.merge(ydf, zdf, how='outer')
Rows that appear in either or both ydf and zdf (Union).
```

```
pd.merge(ydf, zdf, how='outer',
        indicator=True)
.query('_merge == "left_only"')
.drop(['_merge'], axis=1)
Rows that appear in ydf but not zdf (Setdiff).
```