

Analiza danych

prof. UAM dr hab. Tomasz Górecki

tomasz.gorecki@amu.edu.pl
<http://drizzt.home.amu.edu.pl>

Zakład Statystyki Matematycznej i Analizy Danych
Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza w Poznaniu



-  Biecek, P. (2016). *Odkrywać! Ujawniać! Objaśniać!* Wydawnictwo UW.
-  Biecek, P. (2017). *Przewodnik po pakiecie R*. GiS.
-  Deisenroth, M.P., Faisal, A.A., Ongm C.S. (2022). *Matematyka w uczeniu maszynowym*. Helion.
-  Gągolewski, M. (2016). *Programowanie w języku R*. PWN.
-  Ghatak, A. (2019). *Deep Learning with R*. Springer.
-  Gillespie, C., Lovelace, R. (2018). *Wydajne programowanie w R*. O'Reilly Media.
-  Górecki, T. (2011). *Podstawy statystyki z przykładami w R*. BTC.
-  Hastie, T., Tibshirani, R., Friedman, J. (2009). *The Elements of Statistical Learning*. Springer.

-  James, G. Witten, D., Hastie, T., Tibshirani, R. (2017). *An Introduction to Statistical Learning with Applications in R*. Springer.
-  Krzyśko, M., Wołyński, W., Górecki, T., Skorzybut, M. (2008). *Systemy uczące się*. WNT.
-  Lander, J.P. (2018). *Język R dla każdego*. APN Promise.
-  Nwanganga, F., Chapple, M. (2022). *Praktyczne uczenie maszynowe w języku R*. APN Promise.
-  Muraszkiwicz, M., Nowak, R. (2022). *Sztuczna inteligencja dla inżynierów. Metody ogólne*. Oficyna Wydawnicza PW.
-  Muraszkiwicz, M., Nowak, R. (2023). *Sztuczna inteligencja dla inżynierów. Istotne obszary i zastosowania*. Oficyna Wydawnicza PW.
-  Szeliga, M. (2017). *Data science i uczenie maszynowe*. PWN.



Wickham, H., Golemund, G. (2018). *Język R*. Helion.



Winke, C.O. (2020). *Podstawy wizualizacji danych*. Helion.

Data Scientist



Uses Statistics and Advanced Machine Learning algorithms to predict and provide key answer keys to business questions

Skills: Math, Programming, Statistics



Tech: SQL, Python, R, Cloud

Data Engineer



Develops systems for collecting, storing, and analyzing data at a large scale

Skills: Programming, BigData & Cloud



Tech: SQL, Python, Cloud, Distributed Computing

Data Analyst



Examines information that help an organization to provide teams to develop insights and business strategies

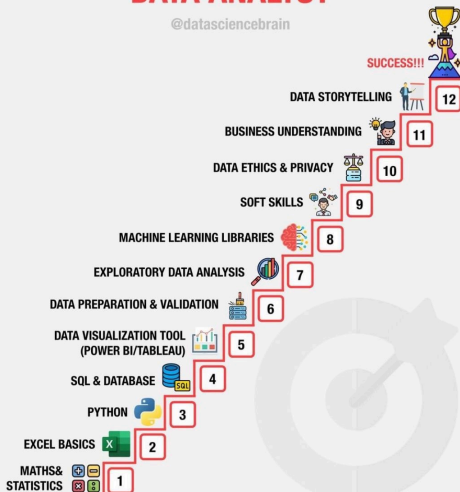
Skills: Communication, Business Knowledge



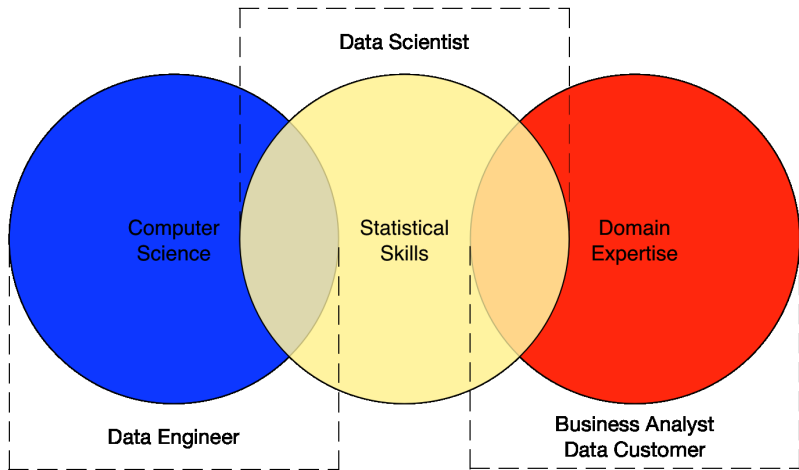
Tech: SQL, Excel, Tableau

SELF GUIDE TO BECOME A DATA ANALYST

@datasciencebrain



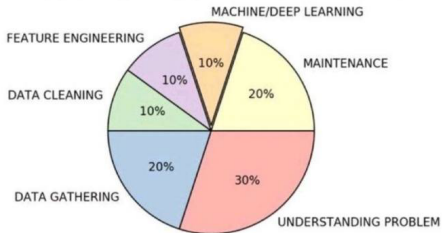
DATA SCIENCE BRAIN™



DATA SCIENTIST JOB - EXPECTATION



DATA SCIENTIST JOB - REALITY



Data Mining

It is the process of automatically discovering useful information in large data repositories.

Machine Learning

It is a set of techniques, which help in dealing with vast data in the most intelligent fashion (by developing algorithms or set of logical rules) to derive actionable insights (delivering search for users in this case)

- ...what we want is a machine that can learn from experience – **Alan Turing (1947)**.
- Machine learning is the study of computer algorithms that allow computer programs to automatically improve through experience – **Tom Mitchell (1997)**.
- Artificial intelligence is the science and engineering of making computers behave in ways that, until recently, we thought required human intelligence – **Andrew Moore (2008)**.



Alan Turing (1912-1954)

ANATOMY OF A DATA SCIENTIST

SALARY 💰

Average salary of data scientists is **\$120,000/year**

EDUCATION 🎓

- **88%** of all data scientists have at least a Master's degree
- **46%** of data scientists have a PhD

BENEFITS 🏠

- Harvard Business Review called data science the **"Sexiest Job of the 21st Century"**
- One of the fastest growing careers in the United States
- **94%** of data science graduates have found jobs since 2011

SKILLS 🧠

- Programming languages (R, Python, SQL, Hive, etc.)
- Statistics
- Multivariable calculus and linear algebra
- Machine learning
- Software engineering
- Wrangle, visualize, and communicate data to management

RESPONSIBILITIES 💡

- Conduct research
- Extract, clean, and analyze data from varied sources
- Solve problems
- Build automation tools
- Communicate findings to management

CAREER POSSIBILITIES 🌐

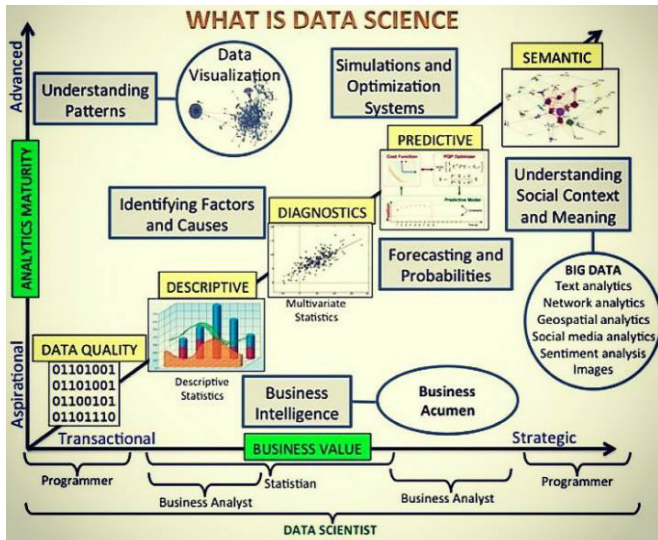
- The majority of data scientists work in the **technology industry**.
- Other options include marketing, consulting, healthcare and pharmaceuticals, finance, government, gaming, and many more.

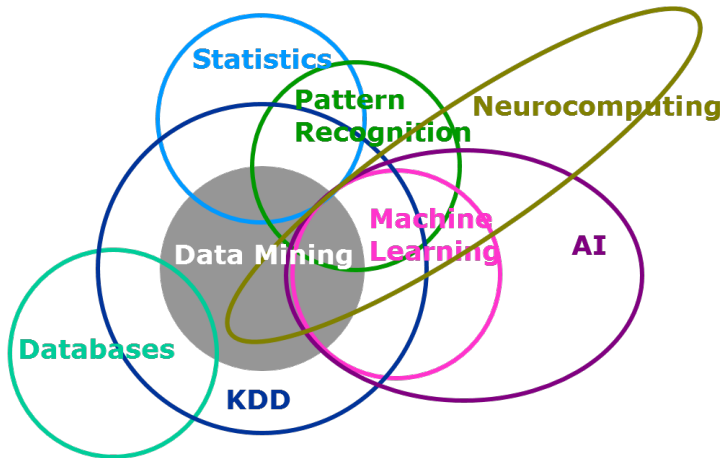
Artificial Intelligence

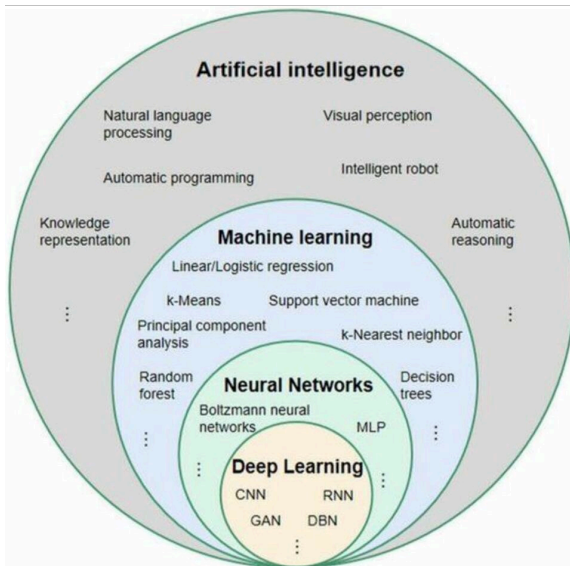
Machine Learning

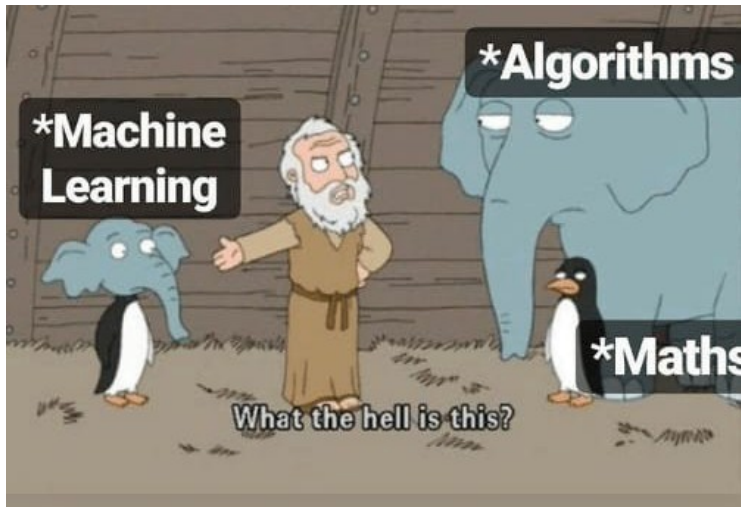
Deep Learning

Tools such as TensorFlow

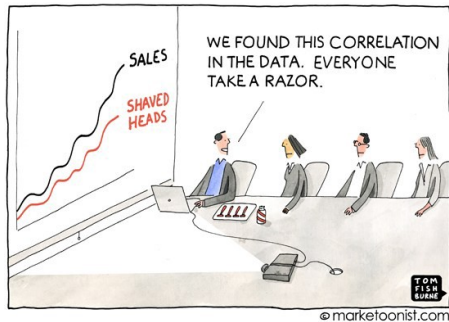




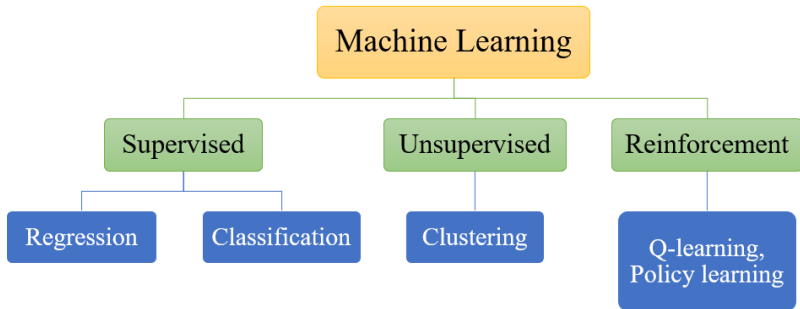




- Any figure that looks interesting or different is usually wrong.
- Any statistics that appears interesting is almost certainly a mistake (double check all findings).
- The more unusual or interesting the data, the more likely they are to have been the result of an error of one kind or another.

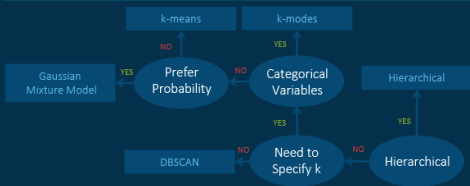


Ehrenberg, A.S.C. (1975). The Teaching of Statistics: Corrections and Comments. *Journal of the Royal Statistical Society. Series A* 138(4):543–545.



Machine Learning Algorithms Cheat Sheet

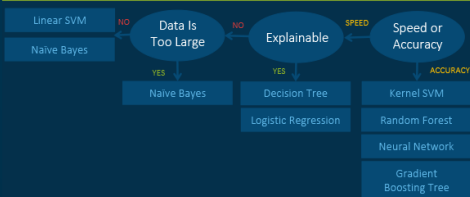
Unsupervised Learning: Clustering



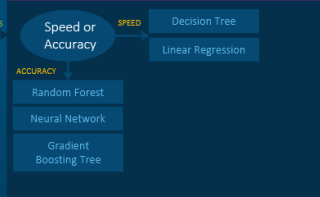
Unsupervised Learning: Dimension Reduction

START

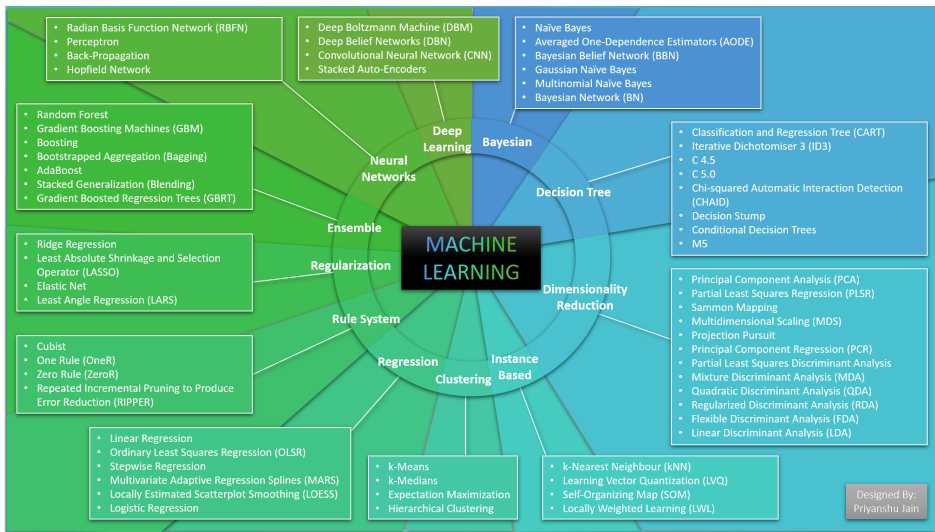
Supervised Learning: Classification



Supervised Learning: Regression



Podział algorytmów ML

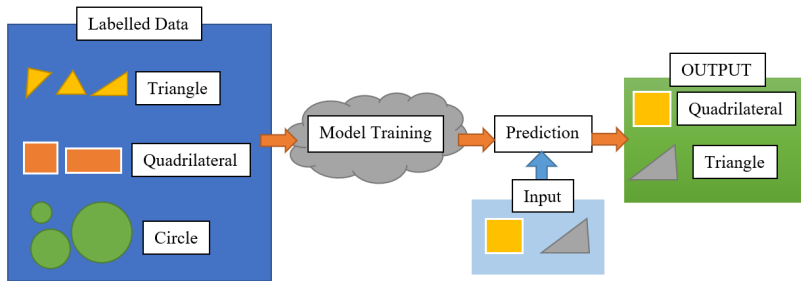


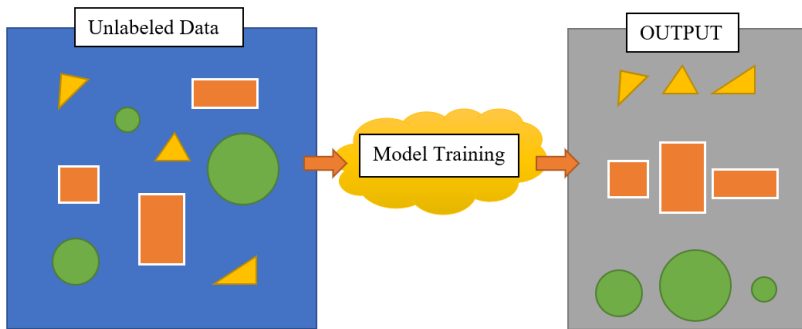
SUPERVISED VS UNSUPERVISED

In supervised learning, for every observation we have some feature values and some target vector or tensor. We use both to train a model that takes in some x values and outputs a predicted value.

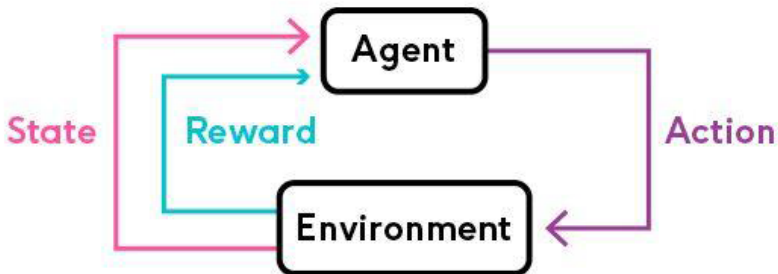
In unsupervised learning, we only have features and do not have a target. This makes the estimation problem more difficult. When possible, use supervised learning.

Chris Albon

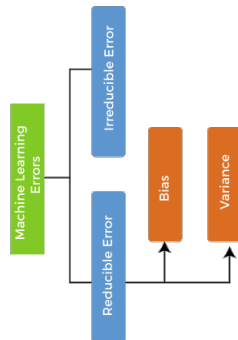




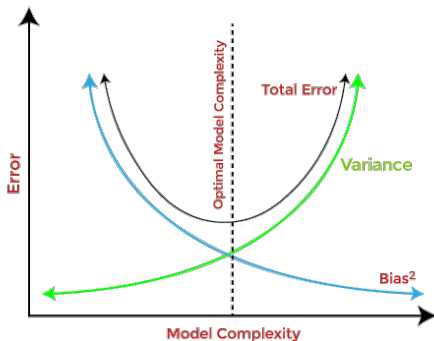
Uczenie przez wzmacnianie (*ang. reinforcement learning*) – w przeciwieństwie do uczenia nadzorowanego i nienadzorowanego nie przygotowuje się zestawu danych uczących. Zamiast tego system uczy się metodą prób i błędów. Z tego powodu sekwencja właściwych decyzji spowoduje wzmocnienie procesu, ponieważ stanowi on najlepsze rozwiązanie danego problemu.



- 1 **Błąd nieredukowalny** (*ang. irreducible error*) – błąd związany z danymi, który nie może być usunięty. Bywa również nazywany **błędem bayesowskim** (*ang. Bayes error*) – najmniejszy, możliwy do osiągnięcia poziom błędu. W wielu problemach za taki błąd możemy uznawać tzw. **błąd ludzki** (*ang. human-level error*), czyli taki jaki popełniają ludzie.
- 2 **Obciążenie** (*ang. bias*) – ogólne przesunięcie błędu względem najmniejszej możliwej wartości.
- 3 **Wariancja** (*ang. variance*) – zmienność wyników ze względu na nieco inne próbki.



Z obciążeniem mamy do czynienia w sytuacji gdy istnieje duża różnica pomiędzy błędem ludzkim (oczekiwanym), a uzyskanym przez model. Jeśli natomiast jest nieznaczna różnica pomiędzy tymi błędami, a duża różnica pomiędzy błędem na zbiorze uczącym i walidacyjnym, to mamy do czynienia z dużą wariancją.



Bias - Variance Tradeoff

$$\text{Error}(x) = \underbrace{\left(\underbrace{E[\hat{f}(x)]}_{\text{predicted}} - \underbrace{f(x)}_{\text{true}} \right)^2}_{\text{Bias}^2} + \underbrace{E\left[\underbrace{\hat{f}(x)}_{\text{predicted}} - \underbrace{E[\hat{f}(x)]}_{\text{average predicted value}} \right]^2}_{\text{Variance}} + \underbrace{\sigma_e^2}_{\text{irreducible error}}$$

Bias²

How much predicted values differ from true values.

Variance

How predictions made on the same value vary on different realizations of the model

BY CHRIS ALBON

Jak sobie radzić:

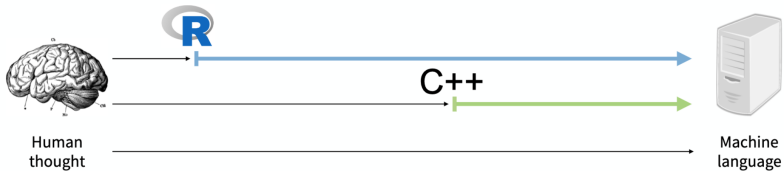
Obciążenie:

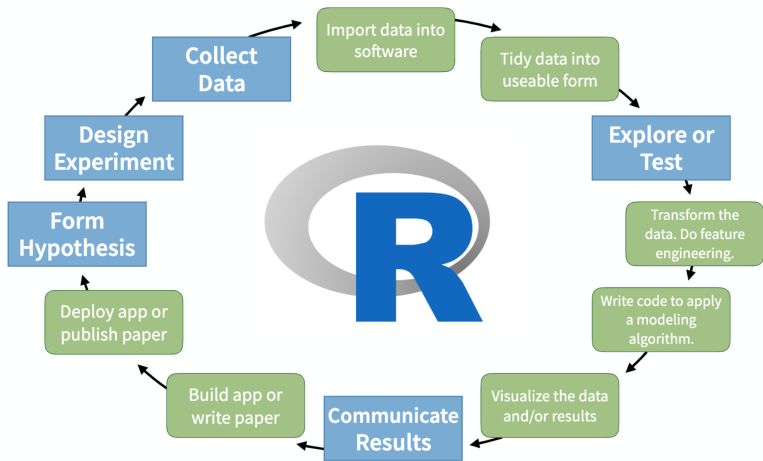
- Wyucz bardziej skomplikowany model.
- Ucz dłużej.
- Użyj lepszego algorytmu optymalizacji.
- Użyj innego modelu.
- Znajdź lepsze hiperparametry.

Wariancja:

- Użyj więcej danych.
- Użyj regularyzacji modelu.
- Zastosuj augmentację.
- Użyj innego modelu.
- Znajdź lepsze hiperparametry.

R - A computer language for scientists





WHAT IS R?

R is a scripting language for statistical data manipulation and analysis. It was inspired by, and is mostly compatible with, the statistical language S developed by AT&T. The name S, obviously standing for statistics, was an allusion to another programming language developed at AT&T with a one-letter name, C. S later was sold to a small firm, which added a GUI interface and named the result S-Plus.

WHY USE R FOR YOUR STATISTICAL WORK?

1 A public-domain implementation of the widely-regarded statistical language; R/S is the de facto standard among professional statisticians



2 comparable, and often superior, in power to commercial products in most senses



3 available for Windows, Macs, Linux



4 in addition to enabling statistical operations, it's a general programming language, so that you can automate your analyses and create new functions



5 object-oriented and functional programming structure



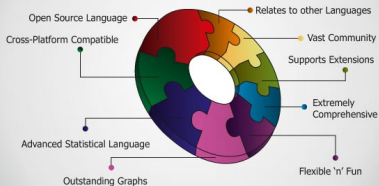
6 your data sets are saved between sessions, so you don't have to reload each time



7 open-software nature means it's easy to get help from the user community, and lots of new functions get contributed by users, many of which are prominent statisticians



Why Learn R?





R jest zaawansowanym pakietem statystycznym jak również językiem programowania istniejącym na platformy Windows, Linux oraz MacOS. Objęty jest licencją GNU GPL i oparty na języku S. Język R jest językiem **interpretowanym**, a nie kompilowanym. Z tego względu program w nim napisany nie będzie tak szybki jak np. program napisany w C++. Można jednak wykorzystać funkcje z bibliotek napisanych np. w C, C++ czy Fortran. Można również używać funkcji R w Javie, Pythonie czy C++, jak również w VB for Statistica.

Pakiet umożliwia również tworzenie zaawansowanych wykresów, które mogą zostać zapisane w formatach takich jak PDF i PNG czy JPG. O sile R stanowi ponad 20 000 bibliotek (główne repozytorium CRAN; istnieją również inne repozytoria np. Bioconductor zawierający blisko 2200 bibliotek do analizy danych genomicznych), przeznaczonych do najróżniejszych zastosowań. Całość wzbogacona jest kompleksową dokumentacją. Dodatkowo bardzo użyteczną cechą R jest dostępność zbiorów danych praktycznie dla każdego zagadnienia.



Wczesna historia R:

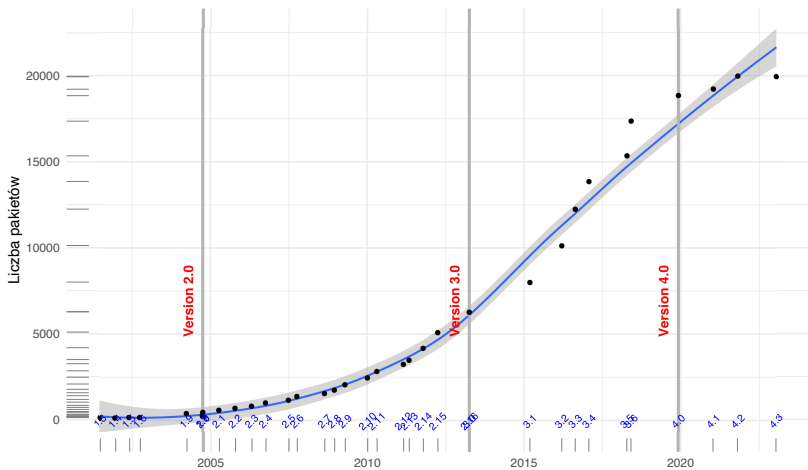


R. Gentleman (1959-) & R. Ihaka (1954-)

1994: Ross Ihaka i Robert Gentleman tworzą pierwszą wersję R (Robert and Ross) na Uniwersytecie w Auckland.

1997: Kurt Hornik i Fritz Leisch tworzą CRAN na Uniwersytecie Technicznym w Wiedniu.

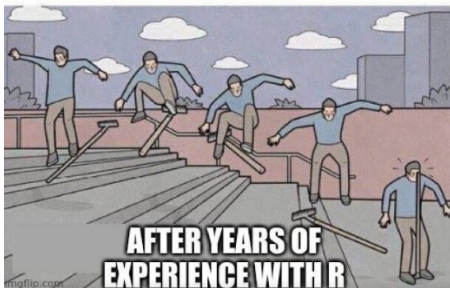
29 lutego 2000: wersja 1.0.





They use(d) R





RStudio jest zintegrowanym środowiskiem programowania (IDE) dla R. Dostępne jest w dwóch edycjach: RStudio Desktop oraz RStudio Server. Dostępne jest dla Windows, macOS oraz Linuxa. Występuje wersja darmowa (open source) oraz komercyjna. Zostało napisane głównie w C++. Pierwsza wersja beta powstała lutym 2011 roku. Wersja 1.0 została wydana 1 listopada 2016 roku.



BuzzR

RStudio anatomy

<https://buzzrheadline.blog/>

Emma Rand

Script file

Write code here

To run code put your cursor on the line and click the **run button**
 Edit to correct errors

⇒ record of commands that worked
 Save scripts with the **.R** extension
 ⇒ syntax will be highlighted
 ⇒ good practice

<- is the assignment operator
 ⇒ puts what is on the right in to the object on the left
 ⇒ Assign results if you want to use them again

Console

When you click run, code is sent to the console and executed

> is the prompt

⇒ do not type it
 ⇒ appears when R is ready for next command

Command output goes here by default

⇒ output is in a different colour
 ⇒ [1] indicates 3.4 is the first element of the output
 ⇒ many commands will not have output, the prompt just reappears

The screenshot shows the RStudio interface with several panes and annotations:

- Script file:** The main editor window contains R code with comments and syntax highlighting. Annotations point to the 'run' button and the script content.
- Environment:** The Environment pane shows the current workspace with variables 'mx' (value 3.4) and 'x' (value [1] 2 4 1 4 6). An annotation points to this pane with the text 'Environment: where saved output goes'.
- Console:** The Console pane shows the execution of the script, with output color-coded. Annotations point to the prompt and the output.
- Packages:** The Packages pane shows a list of installed and available packages, including 'abind', 'abind64', 'acpack', 'ade4', 'agricolae', and 'AlgDesign'. An annotation points to the 'Install' button.

Script: where you write code

Console: where output goes

Environment

Name objects by assignment to use them again

All the **objects** you created in your session

Saving the environment saves all the objects, but not the code with **.RData** extension

History

A history of every command you sent to the console, mistakes included.

File can be saved but usually you just need the script

Packages

Many functions come with R
 A huge amount of extra functionality is available in packages

Packages can be installed by clicking the Install button

Help

Access to manual pages for all installed packages

Plots

Figure output appears here

RStudio IDE Cheat Sheet

learn more at www.rstudio.com



The RStudio IDE is an Integrated Development Environment in R that comes in three versions



Desktop IDE

A local version of the IDE for your desktop



Open Source Server

for larger compute resources and remote access



Professional Server

for teams that share large compute resources, large data, and uniform environments for collaboration

Download all at www.rstudio.com. Each provides the same useful interface:

Documents and Apps

Annotations for Documents and Apps pane:

- Open Shiny, R Markdown, knitr, Sweave, LaTeX, .Rd files and more in Source Pane
- Check spelling
- Render output chunk
- Choose output format
- Choose output format
- Insert output chunk
- Jump to previous to next chunk
- Jump to next chunk
- Run selected chunk
- Publish to server
- Show file outline
- Access markdown guide at **Help > Markdown Quick Reference**
- Jump to chunk
- Set knitr options
- Run this and all previous code chunks
- Run this code chunk
- RStudio recognizes that files named **app.R**, **server.R**, **ui.R**, and **global.R** belong to a Shiny app
- Run app
- Choose app location to shinyapps.io, view app, or server
- Manage accounts

Write Code

Annotations for Write Code pane:

- Navigate tabs
- Open in new window
- Save
- Find and replace
- Compile and run selected code
- Run selected code
- Import data file with wizard
- History of past commands to run/add to source
- Display R/RShiny slideshows
- File > New File > R Presentation**
- Load workspace
- Delete all saved objects
- Search inside environment
- Choose environment to display from list of parent environments
- Display objects as fat or grid
- Displays saved objects by type with short description
- View in data viewer
- View function source code
- Path to displayed directory
- A file browser keyed to your working directory. Click on file or directory name to open.
- Change file type
- Working directory
- Maximize, minimize panels
- Drag command history
- Drag pane boundaries

R Support

Annotations for R Support pane:

- Import data file with wizard
- History of past commands to run/add to source
- Display R/RShiny slideshows
- File > New File > R Presentation**
- Load workspace
- Delete all saved objects
- Search inside environment
- Choose environment to display from list of parent environments
- Display objects as fat or grid
- Displays saved objects by type with short description
- View in data viewer
- View function source code
- Path to displayed directory
- A file browser keyed to your working directory. Click on file or directory name to open.
- Change file type
- Working directory
- Maximize, minimize panels
- Drag command history
- Drag pane boundaries

RStudio Pro Features

Annotations for RStudio Pro Features pane:

- Share Project with Collaborators
- Active shared collaborators
- Start new R Session in current project
- Close R Session in project
- Select R Version
- R Version 3.1.3
- R Version 3.1.2
- R Version 3.1.1
- R Version 3.0.3
- R Version 3.0.2
- R Version 3.0.1
- Project System
- File > New Project**
- RStudio saves the call history, workspace, and working directory associated with a project. It reloads each when you re-open a project.

Plots

Annotations for Plots pane:

- RStudio opens plots in a dedicated Plots pane
- Navigate recent plots
- Open in plot
- Export plot
- Delete all plots
- Delete recent plots
- Delete plot
- Click to load package with library() or detach package from installed library
- Click to load package with library() or detach package from installed library

Debug Mode

Annotations for Debug Mode pane:

- Open with **debug()**, **browser()**, or a breakpoint. RStudio will open the debugger mode when it encounters a breakpoint while executing code.
- Launch debugger mode from origin of error
- Open traceback to examine the functions that R called before the error occurred
- Click next to line number to add/remove a breakpoint.
- Highlighted line shows where execution has paused
- Run commands in environment where execution has paused
- Examine variables in executing environment
- Select function in traceback to debug
- Step through code one line at a time
- Step into and out of functions to run
- Resume execution
- Quit debug mode

Version Control with Git or SVN

Annotations for Version Control pane:

- Turn on at **Tools > Project Options > Git/SVN**
- Stage files
- Show file diff
- Commit staged files
- Push/Pull
- View history
- Added
- Deleted
- Modified
- Renamed
- Untracked
- Open shell to type commands
- current branch

Package Writing

Annotations for Package Writing pane:

- File > New Project > New Directory > R Package**
- Turn project into package
- Enable roxygen documentation with **Tools > Project Options > Build Tools**
- Roxygen guide at **Help > Roxygen Quick Reference**

	Windows/Linux	Mac		Windows/Linux	Mac		Windows/Linux	Mac
1 LAYOUT			4 WRITE CODE			5 DEBUG CODE		
Move focus to Source Editor	Ctrl+1	Ctrl+1	Attempt completion	Tab or Ctrl+Space	Tab or Cmd+Space	Toggle Breakpoint	Shift+F9	Shift+F9
Move focus to Console	Ctrl+2	Ctrl+2	Navigate candidates	↑/↓	↑/↓	Execute Next Line	F10	F10
Move focus to Help	Ctrl+3	Ctrl+3	Accept candidate	Enter, Tab, or →	Enter, Tab, or →	Step Into Function	Shift+F4	Shift+F4
Show History	Ctrl+4	Ctrl+4	Dismiss candidates	Esc	Esc	Finish Function/Loop	Shift+F6	Shift+F6
Show Files	Ctrl+5	Ctrl+5	Undo	Ctrl+Z	Cmd+Z	Continue	Shift+F5	Shift+F5
Show Plots	Ctrl+6	Ctrl+6	Redo	Ctrl+Shift+Z	Cmd+Shift+Z	Stop Debugging	Shift+F8	Shift+F8
Show Packages	Ctrl+7	Ctrl+7	Cut	Ctrl+X	Cmd+X	6 VERSION CONTROL		
Show Environment	Ctrl+8	Ctrl+8	Copy	Ctrl+C	Cmd+C	Show diff	Ctrl+Alt+D	Ctrl+Option+D
Show Git/SVN	Ctrl+9	Ctrl+9	Paste	Ctrl+V	Cmd+V	Commit changes	Ctrl+Alt+M	Ctrl+Option+M
Show Build	Ctrl+0	Ctrl+0	Select All	Ctrl+A	Cmd+A	Scroll diff view	Ctrl+↕	Ctrl+↕
2 RUN CODE			Delete Line	Ctrl+D	Cmd+D	Stage/Unstage (Git)	Spacebar	Spacebar
Search command history	Ctrl+↕	Cmd+↕	Select	Shift+(Arrow)	Shift+(Arrow)	Stage/Unstage and move to next	Enter	Enter
Navigate command history	↑/↓	↑/↓	Select Word	Ctrl+Shift+↔	Option+Shift+↔	7 MAKE PACKAGES		
Move cursor to start of line	Home	Cmd+←	Select to Line Start	Alt+Shift+←	Cmd+Shift+←	Build and Reload	Ctrl+Shift+B	Cmd+Shift+B
Move cursor to end of line	End	Cmd+→	Select to Line End	Alt+Shift+→	Cmd+Shift+→	Load All (devtools)	Ctrl+Shift+L	Cmd+Shift+L
Change working directory	Ctrl+Shift+H	Ctrl+Shift+H	Select Page Up/Down	Shift+PageUp/Down	Shift+PageUp/Down	Test Package (Desktop)	Ctrl+Shift+T	Cmd+Shift+T
Interrupt current command	Esc	Esc	Select to Start/End	Shift+Alt+↔	Cmd+Shift+↔	Test Package (Web)	Ctrl+Alt+F7	Cmd+Alt+F7
Clear console	Ctrl+Q	Ctrl+H	Delete Word Left	Ctrl+Backspace	Ctrl+Opt+Backspace	Check Package	Ctrl+Shift+E	Cmd+Shift+E
Quit Session (desktop only)	Ctrl+Q	Cmd+Q	Delete Word Right	Option+Delete	Option+Delete	Document Package	Ctrl+Shift+D	Cmd+Shift+D
Restart R Session	Ctrl+Shift+F10	Cmd+Shift+F10	Delete to Line End	Ctrl+K	Ctrl+K	8 DOCUMENTS AND APPS		
Run current line/selection	Ctrl+Enter	Cmd+Shift+Enter	Delete to Line Start	Option+Backspace	Option+Backspace	Preview HTML (Markdown, etc.)	Ctrl+Shift+K	Cmd+Shift+K
Run current (retain cursor)	Alt+Enter	Option+Enter	Indent	Tab (at start of line)	Tab (at start of line)	Knit Document (knitr)	Ctrl+Shift+K	Cmd+Shift+K
Run from current to end	Ctrl+Alt+E	Cmd+Option+E	Outdent	Shift+Tab	Shift+Tab	Compile Notebook	Ctrl+Shift+K	Cmd+Shift+K
Run the current function	Ctrl+Alt+F	Cmd+Option+F	Yank line up to cursor	Ctrl+U	Ctrl+U	Compile PDF (TeX and Sweave)	Ctrl+Shift+K	Cmd+Shift+K
Source a file	Ctrl+Shift+O	Cmd+Shift+O	Yank line after cursor	Ctrl+K	Ctrl+K	Insert chunk (Sweave and Knitr)	Ctrl+Alt+I	Cmd+Option+I
Source the current file	Ctrl+Shift+S	Cmd+Shift+S	Insert %>%	Ctrl+Shift+M	Cmd+Shift+M	Insert code section	Ctrl+Shift+R	Cmd+Shift+R
Source with echo	Ctrl+Shift+Enter	Cmd+Shift+Enter	Show help for function	F1	F1	Re-run previous region	Ctrl+Shift+P	Cmd+Option+P
3 NAVIGATE CODE			Show source code	F2	F2	Run current document	Ctrl+Alt+R	Cmd+Option+R
Goto File/Function	Ctrl+↕	Cmd+↕	New document	Ctrl+Shift+N	Cmd+Shift+N	Run from start to current line	Ctrl+Alt+B	Cmd+Option+B
Fold Selected	Alt+L	Cmd+Option+L	New document (Chrome)	Ctrl+Shift+N	Cmd+Shift+N	Run the current code section	Ctrl+Alt+T	Cmd+Option+T
Unfold Selected	Shift+Alt+L	Cmd+Shift+Option+L	Open document	Ctrl+O	Cmd+O	Run previous Sweave/Rmd code	Ctrl+Alt+P	Cmd+Option+P
Fold All	Alt+O	Cmd+Option+O	Save document	Ctrl+S	Cmd+S	Run the current chunk	Ctrl+Alt+C	Cmd+Option+C
Unfold All	Shift+Alt+O	Cmd+Shift+Option+O	Close document	Ctrl+W	Cmd+W	Run the next chunk	Ctrl+Alt+N	Cmd+Option+N
Go to line	Shift+Alt+G	Cmd+Shift+Option+G	Close all documents (Chrome)	Ctrl+Alt+W	Cmd+Shift+W	Sync Editor & PDF Preview	Ctrl+FB	Cmd+FB
Jump to	Shift+Alt+J	Cmd+Shift+Option+J	Extract function	Ctrl+Alt+F	Cmd+Option+F	Previous plot	Ctrl+Alt+F11	Cmd+Option+F11
Switch to tab	Ctrl+Shift+.	Ctrl+Shift+.	Extract variable	Ctrl+Alt+V	Cmd+Option+V	Next plot	Ctrl+Alt+F12	Cmd+Option+F12
Previous tab	Ctrl+F11	Ctrl+F11	Reindent lines	Ctrl+I	Cmd+I	Show Keyboard Shortcuts		
Next tab	Ctrl+F12	Ctrl+F12	(Un)Comment lines	Ctrl+Shift+C	Cmd+Shift+C	Alt+Shift+K	Option+Shift+K	
First tab	Ctrl+Shift+F11	Ctrl+Shift+F11	Reflex Comment	Ctrl+Shift+V	Cmd+Shift+V	Why RStudio Server Pro?		
Last tab	Ctrl+Shift+F12	Ctrl+Shift+F12	Reformat Selection	Ctrl+Shift+A	Cmd+Shift+A	Do everything you would do with the open source server with a commercial license, support, and more.		
Navigate back	Ctrl+F9	Cmd+F9	Select within braces	Ctrl+Shift+E	Cmd+Shift+E	• edit the same project at the same time as others		
Navigate forward	Ctrl+F10	Cmd+F10	Show Diagnostics	Ctrl+Shift+P	Cmd+Shift+P	• switch easily from one version of R to a different version		
Jump to Brace	Ctrl+P	Cmd+P	Transpose Letters	Ctrl+T	Cmd+T	• open and run multiple R sessions simultaneously		
Select within Braces	Ctrl+Shift+Alt+E	Ctrl+Shift+Alt+E	Move Lines Up/Down	Alt+↕	Option+↕	• see what you and others are doing on your server		
Use Selection for Find	Ctrl+F3	Cmd+F3	Copy Lines Up/Down	Shift+Alt+↕	Cmd+Option+↕	• tune your resources to improve performance		
Find in Files	Ctrl+Shift+F	Cmd+Shift+F	Add New Cursor Above	Ctrl+Alt+Up	Ctrl+Alt+Up	• integrate with your authentication, authorization, and audit practices		
Find Next	Win: F3, Linux: Cmd+F	Cmd+F	Add New Cursor Below	Ctrl+Alt+Down	Ctrl+Alt+Down	Download a free 45 day evaluation at www.rstudio.com/products/rstudio-server-pro/		
Find Previous	W: Shift+F3, L: Ctrl+Shift	Cmd+Shift+G	Move Active Cursor Up	Ctrl+Alt+Shift+Down	Ctrl+Alt+Shift+Down			
Jump to Start/End	Ctrl+↕	Option+↕	Find and Replace	Ctrl+F	Cmd+F			
Jump to Start/End	Ctrl+↕	Cmd+↕	Use Selection for Find	Ctrl+Shift+F	Cmd+Shift+F			
Jump to Start/End	Ctrl+↕	Cmd+↕	Replace and Find	Ctrl+Shift+J	Cmd+Shift+J			

BuzzR

Code explained by a novice

with notes

<https://buzzrbeeeline.blog/>

Emma Rand
and Elliot Spray

It's a challenge for an experienced user to remember what it was like to be totally new to R and come up with explanations that don't draw on understanding, and jargon, developed subsequently. So I asked a novice, Elliot, to explain a piece of code in his own words. I've add notes because in the long term, it does help you understand things more quickly if you become familiar with the terminology.

Assignment

<- is the assignment operator puts what is on the right in to the object on the left.

We assign the output of functions to an object so we can use them again

The action of giving the name

The data that will be named

Functions

We use functions in R for everything. All functions take arguments and have the form: `functionname(arg1, arg2...)`

The name that is being given

```
freq <- c(4, 13, 15, 13, 5)
```

Brackets are used to contain the data that you wish to name

Object names

The name you use is your choice except that names must:

- consist of letters, numbers and the dot or underline
- start with a letter or the dot not followed by a number

The c is used as more than one value is used inside the brackets

Arguments

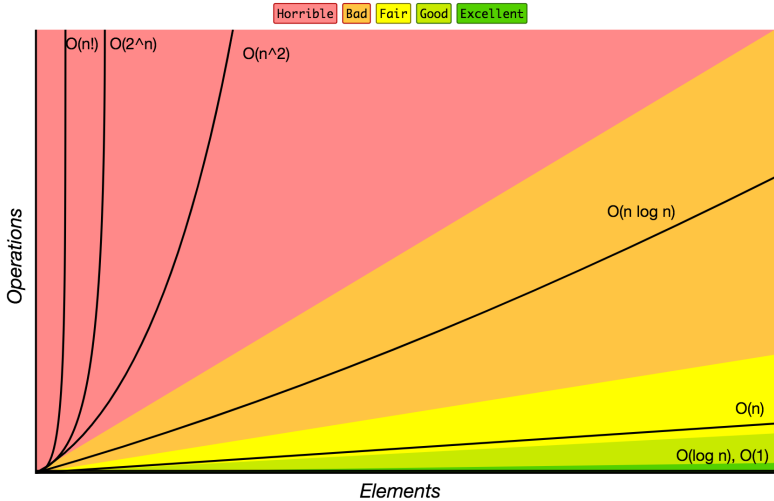
What you want to do the function on and how you want to do it

- go inside the brackets
- are separated by commas
- may have default values

The c() function

Is one of the most used functions in R. It is used to create a vector with values - it combines values. A vector is just a group of values in single column or row

Big-O Complexity Chart



Array Sorting Algorithms

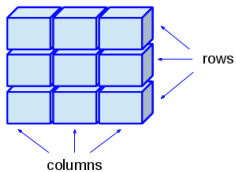
Algorithm	Time Complexity			Space Complexity
	Best	Average	Worst	Worst
<u>Quicksort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n^2)$	$O(\log(n))$
<u>Mergesort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Timsort</u>	$\Omega(n)$	$\theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Heapsort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n \log(n))$	$O(1)$
<u>Bubble Sort</u>	$\Omega(n)$	$\theta(n^2)$	$O(n^2)$	$O(1)$
<u>Insertion Sort</u>	$\Omega(n)$	$\theta(n^2)$	$O(n^2)$	$O(1)$
<u>Selection Sort</u>	$\Omega(n^2)$	$\theta(n^2)$	$O(n^2)$	$O(1)$
<u>Tree Sort</u>	$\Omega(n \log(n))$	$\theta(n \log(n))$	$O(n^2)$	$O(n)$
<u>Shell Sort</u>	$\Omega(n \log(n))$	$\theta(n(\log(n))^2)$	$O(n(\log(n))^2)$	$O(1)$
<u>Bucket Sort</u>	$\Omega(n+k)$	$\theta(n+k)$	$O(n^2)$	$O(n)$
<u>Radix Sort</u>	$\Omega(nk)$	$\theta(nk)$	$O(nk)$	$O(n+k)$
<u>Counting Sort</u>	$\Omega(n+k)$	$\theta(n+k)$	$O(n+k)$	$O(k)$
<u>Cubesort</u>	$\Omega(n)$	$\theta(n \log(n))$	$O(n \log(n))$	$O(n)$

Obiekty oraz typy

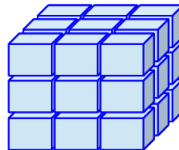
Vector



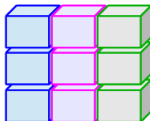
Matrix



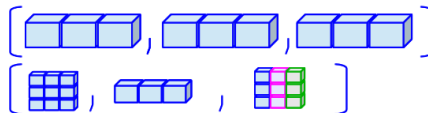
Array



Data Frame
(Table)



Lists



EDA odpowiada za proces eksploracji danych w celu późniejszego wykonania analiz. Skupiona jest głównie na obejrzeniu danych, podsumowaniu ich, znalezieniu głównych charakterystyk i wizualizacji.

Check the summary of each of your attributes (mean, median, range, inter-quartile range)

Check for Missing values in Data

Check for outliers in data

Check for the correlation among attributes

Plot individual variables and check the distribution of data

Eksploracyjna analiza danych (EDA)



dplyr jest pakietem zaprojektowanym do wydajnego (napisany w C++) i efektywnego manipulowania danymi. Dodatkowo wszystko co robimy na ramce danych możemy też zrobić na innych obiektach np. tabelach w bazie danych. Jest on obecnie częścią większej grupy pakietów zwanej tidyverse.



Data transformation with dplyr : : CHEAT SHEET



dplyr functions work with pipes and expect tidy data. In tidy data:



Each variable is in its own column



Each observation, or case, is in its own row



x %>% f(y) becomes f(x, y)

Summarise Cases

Apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

summary function



summarise(data, ...) Compute table of summaries. summarise(mtcars, avg = mean(mpg))



count(data, ..., wt = NULL, sort = FALSE, name = NULL) Count number of rows in each group defined by the variables in ... Also **tally**() count(mtcars, cyl)

Group Cases

Use **group_by**(data, ..., add = FALSE, drop = TRUE) to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.



mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))

Use **rowwise**(data, ...) to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidy cheat sheet for list-column workflow.



starwars %>%
rowwise() %>%
mutate(films_count = length(films))

ungroup(x, ...) Returns ungrouped copy of table. ungroup(mtcars)

Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.



filter(data, ..., preserve = FALSE) Extract rows that meet logical criteria. filter(mtcars, mpg > 20)



distinct(data, ..., keep_all = FALSE) Remove rows with duplicate values. distinct(mtcars, gear)



slice(data, ..., preserve = FALSE) Select rows by position. slice(mtcars, 10:15)



slice_sample(data, ..., n, prop, weight_by = NULL, replace = FALSE) Randomly select rows. Use n to select a number of rows and prop to select a fraction of rows. slice_sample(mtcars, n = 5, replace = TRUE)



slice_min(data, order_by, ..., n, prop, with_ties = TRUE) and **slice_max**() Select rows with the lowest and highest values. slice_min(mtcars, mpg, prop = 0.25)



slice_head(data, ..., n, prop) and **slice_tail**() Select the first or last rows. slice_head(mtcars, n = 5)

Logical and boolean operators to use with filter()

`==` `<` `<=` `is.na()` `%in%` `|` `xor()`
`!=` `>` `>=` `!is.na()` `!` `&`

See ?base::Logic and ?Comparison for help.

ARRANGE CASES



arrange(data, ..., by_group = FALSE) Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low. arrange(mtcars, mpg)
arrange(mtcars, desc(mpg))

ADD CASES



add_row(data, ..., before = NULL, after = NULL) Add one or more rows to a table. add_row(cars, speed = 1, dist = 1)

Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.



pull(data, var = 1, name = NULL, ...) Extract column(s) as a vector, by name or index. pull(mtcars, wt)



select(data, ...) Extract columns as a table. select(mtcars, mpg, wt)



relocate(data, ..., before = NULL, after = NULL) Move columns to new position. relocate(mtcars, mpg, cyl, after = last_col())

Use these helpers with select() and across()

e.g. select(mtcars, mpg:cyl)

contains(match) **num_range**(prefix, range) e.g. mpg:cyl
ends_with(match) **all_of**()/any_of() vars e.g. gear
starts_with(match) **matches**(match) everything()

MANIPULATE MULTIPLE VARIABLES AT ONCE



across(cols, funs, ..., names = NULL) Summarise or mutate multiple columns in the same way. summarise(mtcars, across(everything(), mean))



c_across(cols) Compute across columns in row-wise data. transmute(rowwise(kggas, sum) = total[c_across(1:2)])

MAKE NEW VARIABLES

Apply **vectorized functions** to columns. Vectorized functions take vectors as input and return vectors of the same length as output (see back).

vectorized function



mutate(data, ..., keep = "all", before = NULL, after = NULL) Compute new column(s), also **add_column**(), **add_count**(), and **add_tally**() mutate(mtcars, gpm = 1 / mpg)



transmute(data, ...) Compute new column(s), drop others. transmute(mtcars, gpm = 1 / mpg)



rename(data, ...) Rename columns. Use **rename_with**() to rename with a function. rename(cars, distance = dist)



Vectorized Functions

TO USE WITH MUTATE ()

mutate() and **transmute()** apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function

OFFSET

dplyr::lag() - offset elements by 1
dplyr::lead() - offset elements by -1

CUMULATIVE AGGREGATE

dplyr::cumall() - cumulative all()
dplyr::cumany() - cumulative any()
dplyr::cummax() - cumulative max()
dplyr::cummean() - cumulative mean()
dplyr::cummin() - cumulative min()
dplyr::cumprod() - cumulative prod()
dplyr::cumsum() - cumulative sum()

RANKING

dplyr::cume_dist() - proportion of all values <=
dplyr::dense_rank() - rank w ties = min, no gaps
dplyr::min_rank() - rank with ties = min
dplyr::ntile() - bins into n bins
dplyr::percent_rank() - min_rank scaled to [0,1]
dplyr::row_number() - rank with ties = "first"

MATH

+, **-**, *****, **/**, **%**, **%%**, **%%%** - arithmetic ops
log(), **log2()**, **log10()** - logs
<, **<=**, **>**, **>=**, **==**, **!=** - logical comparisons
dplyr::between() - x => left & x <= right
dplyr::near() - safe == for floating point numbers

MISCELLANEOUS

dplyr::case_when() - multi-case if_else()
 starwars %>%
 mutate(type = case_when(
 height > 200 ~ "tall", # "tall"
 species == "Droid" ~ "robot", # "robot"
 TRUE ~ "other"
))
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
dplyr::pmax() - element-wise max()
dplyr::pmin() - element-wise min()

Summary Functions

TO USE WITH SUMMARISE ()

summarise() applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

COUNT

dplyr::n() - number of values/rows
dplyr::n_distinct() - # of uniques
dplyr::sum(is.na()) - # of non-NA's

POSITION

dplyr::mean() - mean, also **dplyr::mean(is.na())**
dplyr::median() - median

LOGICAL

dplyr::mean() - proportion of TRUE's
dplyr::sum() - # of TRUE's

ORDER

dplyr::first() - first value
dplyr::last() - last value
dplyr::nth() - value in nth location of vector

RANK

dplyr::quantile() - nth quantile
dplyr::min() - minimum value
dplyr::max() - maximum value

SPREAD

dplyr::IQR() - Inter-Quartile Range
dplyr::mad() - median absolute deviation
dplyr::sd() - standard deviation
dplyr::var() - variance

Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

```

#> tibble: 2 x 2
  # A tibble: 2 x 2
  #   rownames column
  #   <chr>   <dbl>
1  "row1"  1
2  "row2"  2

#> tibble: 2 x 2
  # A tibble: 2 x 2
  #   column rownames
  #   <dbl> <chr>
1     1    "row1"
2     2    "row2"
  
```

Also tibble: **has_rownames()** and tibble: **remove_rownames()**.

Combine Tables

COMBINE VARIABLES

x + **y** = **combined**

bind_cols(..., name_repair) Returns tables placed side by side as a single table. Column lengths must be equal. Columns will NOT be matched by id (to do that look at Relational Data below), so be sure to check that both tables are ordered the way you want before binding.

RELATIONAL DATA

Use a "Mutating Join" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

dplyr::left_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ..., keep = FALSE, na_matches = "na") Join matching values from y to x.

dplyr::right_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ..., keep = FALSE, na_matches = "na") Join matching values from x to y.

dplyr::inner_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ..., keep = FALSE, na_matches = "na") Join data. Retain only rows with matches.

dplyr::full_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ..., keep = FALSE, na_matches = "na") Join data. Retain all values, all rows.

COLUMN MATCHING FOR JOINS

dplyr::usecols(c("col1", "col2", ...)) Use by = c("col1", "col2", ...) to specify one or more common columns to match on.
dplyr::left_join(x, y, by = "x")

dplyr::left_join(x, y, by = c("col1" = "col2")) Use a named vector, by = c("col1" = "col2"), to match on columns that have different names in each table.
dplyr::left_join(x, y, by = c("C" = "D"))

dplyr::left_join(x, y, by = c("C" = "D"), suffix = c("x", "y")) Use suffix to specify the suffix to give to unmatched columns that have the same name in both tables.
dplyr::left_join(x, y, by = c("C" = "D"), suffix = c("x", "y"))

COMBINE CASES

x + **y** = **combined**

bind_rows(..., id = NULL) Returns tables one on top of the other as a single table. Set id to a column name to add a column of the original table names (as pictured).

Use a "Filtering Join" to filter one table against the rows of another.




dplyr::semi_join(x, y, by = NULL, copy = FALSE, ..., na_matches = "na") Return rows of x that have a match in y. Use to see what will be included in a join.

dplyr::anti_join(x, y, by = NULL, copy = FALSE, ..., na_matches = "na") Return rows of x that do not have a match in y. Use to see what will not be included in a join.

Use a "Nest Join" to inner join one table to another into a nested data frame.

dplyr::nest_join(x, y, by = NULL, copy = FALSE, keep = FALSE, name = NULL, ...) Join data, nesting matches from y in a single new data frame column.

SET OPERATIONS

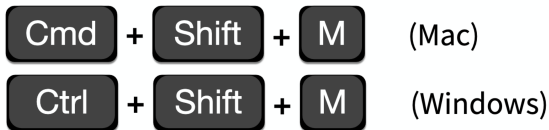
dplyr::intersect(x, y, ...) Rows that appear in both x and y. 
dplyr::setdiff(x, y, ...) Rows that appear in x but not y. 
dplyr::union(x, y, ...) Rows that appear in x or y. (Duplicates removed) 
dplyr::union_all(x, y, ...) (Duplicates removed) retains duplicates.

Use **dplyr::setequal()** to test whether two data sets contain the exact same rows (in any order).

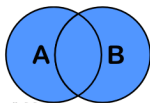


- filter,
- slice, top_n,
- arrange,
- %>% – strumienie.

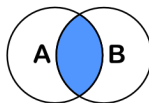
Shortcut to type %>%



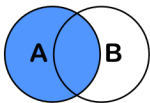
- `select`,
- `mutate`,
- `group_by`, `count`,
- `summarise`.



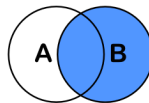
full_join(A,B)



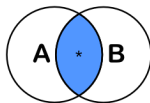
inner_join(A,B)



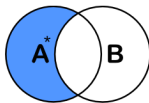
left_join(A,B)



right_join(A,B)



semi_join(A,B)



anti_join(A,B)

ID	X1
1	a1
2	a2

ID	X2
2	b1
3	b2

inner_join

ID	X1	X2
2	a2	b1

left_join

ID	X1	X2
1	a1	NA
2	a2	b1

right_join

ID	X1	X2
2	a2	b1
3	NA	b2

full_join

ID	X1	X2
1	a1	NA
2	a2	b1
3	NA	b2

semi_join

ID	X1
2	a2

anti_join

ID	X1
1	a1

A data set is **tidy** iff:

1. Each **variable** is in its own **column**
2. Each **case** is in its own **row**
3. Each **value** is in its own **cell**

tame data

baker	cinnamon_1	cardamom_2	nutmeg_3
Emma	1	0	1
Harry	1	1	1
Ruby	1	0	1
Zainab	0	NA	0

trial	Emma	Harry	Ruby	Zainab
cinnamon_1	1	1	1	0
cardamom_2	0	1	0	NA
nutmeg_3	1	1	1	0

tidy data

baker	spice	order	correct
Emma	Cinnamon	1	1
Harry	Cinnamon	1	1
Ruby	Cinnamon	1	1
Zainab	Cinnamon	1	0
Emma	Cardamom	2	0
Harry	Cardamom	2	1
Ruby	Cardamom	2	0
Zainab	Cardamom	2	NA
Emma	Nutmeg	3	1
Harry	Nutmeg	3	1
Ruby	Nutmeg	3	1
Zainab	Nutmeg	3	0

wide

id	x	y	z
1	a	c	e
2	b	d	f

long

id	key	val
1	x	a
2	x	b
1	y	c
2	y	d
1	z	e
2	z	f

Data tidying with tidyr : : CHEAT SHEET

Tidy data is a way to organize tabular data in a consistent data structure across packages.

A table is tidy if:



Each **variable** is in its own **column**



Each **observation**, or **case**, is in its own row



Access **variables** as **vectors**



Preserve **cases** in vectorized operations

Tibbles

AN ENHANCED DATA FRAME

Tibbles are a table format provided by the **tibble** package. They inherit the data frame class, but have improved behaviors:

- **Subset** a new tibble with `[]`, a vector with `[]` and `$`.
- **No partial matching** when subsetting columns.
- **Display** concise views of the data on one screen.

`options(tibble.print_max = n, tibble.print_min = m, tibble.width = inf)` Control default display settings.

`View()` or `glimpse()` View the entire data set.

CONSTRUCT A TIBBLE

`tibble(...)` Construct by columns.

`tibble(x = 1:3, y = c("a", "b", "c"))`

`tibble(...)` Construct by rows.

`tibble(x = "a", y = 1, z = 2), tibble(x = "b", y = 2, z = 3), tibble(x = "c", y = 3, z = 1))`



`as_tibble(x, ...)` Convert a data frame to a tibble.

`enframe(x, name = "name", value = "value")`

Convert a named vector to a tibble. Also `deframe()`.

`is_tibble(x)` Test whether x is a tibble.



Reshape Data - Pivot data to reorganize values into a new layout.

table4a

country	year	cases
A	1999	29K
B	2000	29K
C	2000	29K

country	year	cases
A	1999	29K
B	2000	29K
C	2000	29K

`pivot_longer(data, cols, names_to = "name", values_to = "value", values_drop_na = FALSE)`

"Lengthen" data by collapsing several columns into two. Column names move to a new names_to column and values to a new values_to column.

`pivot_longer(table4a, cols = 2:3, names_to = "year", values_to = "cases")`

table3

country	year	cases	count
A	1999	29K	19M
A	1999	29K	19M
A	2000	29K	19M
B	1999	29K	17M
B	1999	29K	17M
B	2000	29K	17M
C	1999	29K	17M
C	1999	29K	17M
C	2000	29K	17M

`pivot_wider(data, names_from = "name", values_from = "value")`

The inverse of `pivot_longer()`. "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

`pivot_wider(table2, names_from = type, values_from = count)`

Split Cells

- Use these functions to split or combine cells into individual, isolated values.

table5

country	century	year	cases
A	20	03	29K
B	19	03	29K
C	19	03	29K

`unite(data, col, ..., sep = "_", remove = TRUE, na.rm = FALSE)` Collapse cells across several columns into a single column.

`unite(table5, century, year, col = "year", sep = "_")`

table1

country	year	rate
A	1999	0.7K/19M
A	2000	0.7K/19M
B	1999	0.7K/17M
B	2000	0.7K/17M

`separate(data, col, into, sep = "[[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", ...)` Separate each cell in a column into several columns. Also `extract()`.

`separate(table3, rate, sep = "/", into = c("cases", "pop"))`

table1

country	year	rate
A	1999	0.7K
A	1999	0.7K
A	2000	0.7K
B	1999	0.7K
B	1999	0.7K
B	2000	0.7K

`separate_rows(data, ..., sep = "[[:alnum:]]+", convert = FALSE)` Separate each cell in a column into several rows.

`separate_rows(table3, rate, sep = "/")`

Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

`expand(data, ...)` Create a new tibble with all possible combinations of the values of the variables listed in ... Drop other variables.

`expand(mtcars, cyl, gear, carb)`

`complete(data, ..., fill = NA)` Add missing combinations of values of variables listed in ... Fill remaining variables with NA.

`complete(mtcars, cyl, gear, carb)`

Handle Missing Values

Drop or replace explicit missing values (NA).

table6

country	year	cases
A	1999	29K
A	2000	29K
B	1999	29K
B	2000	29K

`drop_na(data, ...)` Drop rows containing NA's in ... columns.

`drop_na(x, y)`

table6

country	year	cases
A	1999	29K
A	2000	29K
B	1999	29K
B	2000	29K

`fill(data, ..., direction = "down")` Fill in NA's in ... columns using the next or previous value.

`fill(x, y)`

table6

country	year	cases
A	1999	29K
A	2000	29K
B	1999	29K
B	2000	29K

`replace_na(data, replace)` Specify a value to replace NA in selected columns.

`replace_na(x, list(x2 = 2))`