

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)

Return x sorted.

rev(x)

Return x reversed.

table(x)

See counts of values.

unique(x)

See unique 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.

	<code>m[2,]</code> - Select a row	<code>t(m)</code> Transpose
	<code>m[, 1]</code> - Select a column	<code>m %*% n</code> Matrix Multiplication
	<code>m[2, 3]</code> - Select an element	<code>solve(m, n)</code> Find x in: $m \cdot x = n$

Lists

`l <- list(x = 1:5, y = c('a', 'b'))`
A list is a collection of elements which can be of different types.

<code>l[[2]]</code>	<code>l[1]</code>	<code>l\$x</code>	<code>l['y']</code>
Second element of l.	New list with only the first element.	Element named x.	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

<code>df[, 2]</code>	
<code>df[2,]</code>	
<code>df[2, 2]</code>	

List subsetting	
<code>df\$x</code>	<code>df[[2]]</code>

<i>Understanding a data frame</i>
<code>View(df)</code>
See the full data frame.

<i>head(df)</i>
See the first 6 rows.

`nrow(df)`
Number of rows.

`cbind` - Bind columns.

`ncol(df)`
Number of columns.

`rbind` - Bind rows.

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

Strings

<code>paste(x, y, sep = ' ')</code>	Join multiple vectors together.
<code>paste(x, collapse = ' ')</code>	Join elements of a vector together.
<code>grep(pattern, x)</code>	Find regular expression matches in x.
<code>gsub(pattern, replace, x)</code>	Replace matches in x with a string.
<code>toupper(x)</code>	Convert to uppercase.
<code>tolower(x)</code>	Convert to lowercase.
<code>nchar(x)</code>	Number of characters in a string.

Factors

<code>factor(x)</code>	
Turn a vector into a factor. Can set the levels of the factor and the order.	Turn a numeric vector into a factor by 'cutting' into sections.

Statistics

<code>lm(y ~ x, data=df)</code>	Linear model.
<code>glm(y ~ x, data=df)</code>	Generalised linear model.
<code>summary</code>	Get more detailed information out a model.
<code>pairwise.t.test</code>	Perform a t-test for paired data.

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>unif</code>	<code>qunif</code>

Plotting

<code>plot(x)</code>	Values of x in order.
<code>plot(x, y)</code>	Values of x against y.
<code>hist(x)</code>	Histogram of x.

Dates

See the `lubridate` package.

Get an Overview with overviewR: : CHEAT SHEET



Generate Tables

overview_tab generates a data frame that collapses the time condition for each id by taking into account potential gaps in the time frame

id	time	Var1	Var2
A	1990		
A	1991		
A	1992		
B	1990		

id	time
A	1990 - 1992
B	1990

```
output_table <-  
  overview_tab(  
    dat = toydata,  
    id = ccode,  
    time = year)
```

add data frame
define your time and scope variables

overview_crosstab generates a cross table that divides the data based on two conditions

id	time		
		Yes	No
		Yes	
		No	

```
output_crosstab <-  
  overview_crosstab(  
    dat = toydata,  
    cond1 = gdp,  
    cond2 = population,  
    threshold1 = 25000,  
    threshold2 = 27000,  
    id = ccode,  
    time = year  
)
```

define your conditions with cond1 and cond2
set your thresholds

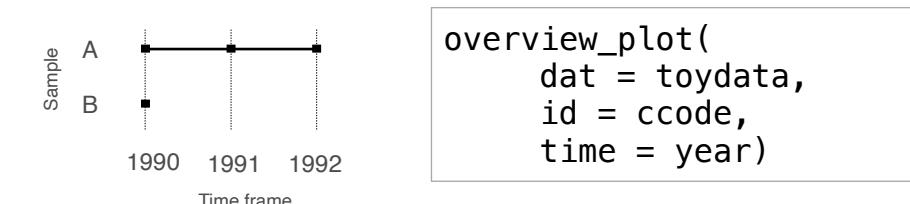
Note, if a data set is used that has multiple observations on the id-time unit, the function automatically aggregates the data set using the mean of condition 1 (**cond1**) and condition 2 (**cond2**).

If you store your results in an object, you can use **overview_print** to export them to a LaTeX output.

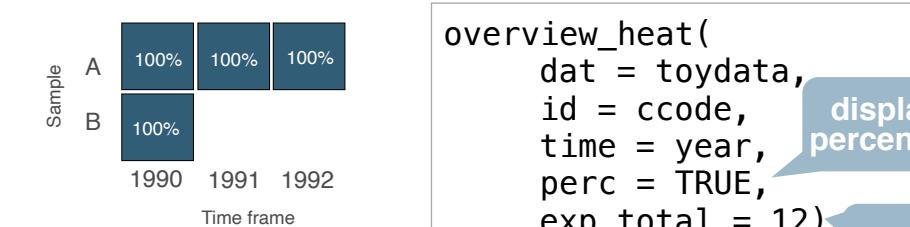
Generate Plots

Sample overview

overview_plot illustrates the information that is generated in **overview_table** in a ggplot2 graphic



overview_heat is similar to **overview_plot** but presents the frequency of data points by id-time-unit in a heat map



Missing values (NAs)

overview_na returns a horizontal ggplot2 bar plot that indicates the amount of missing data (NAs) for each variable



```
overview_na(toydata_with_na,  
            perc = FALSE)
```

FALSE gives total number

Export Results

Tables

overview_print generates a LaTeX output (works with both **overview_tab** and **overview_crosstab** output)

```
overview_print(  
  obj = output_table)  
  
overview_print(  
  obj = output_crosstab,  
  crosstab = TRUE)
```

TRUE for cross tables

The table can be modified with the **title**, **id**, **time**, **cond1**, and **cond2** arguments to replace default names

It also allows to save your output in a .tex file

```
overview_print(  
  obj = output_table,  
  save_out = TRUE,  
  path = "SET-YOUR-PATH",  
  file = "output.tex")
```

define where your output should be stored

The outputs of **overview_tab** and **overview_crosstab** are also compatible with other packages and functions such as **xtable**, **flextable**, or **kable** from **knitr**.

To generate a table in Rmarkdown with **knitr::kable**:

```
knitr::kable(output_table)
```

Plots

As the plots are based on ggplot2, plots can be stored with **ggplot2::ggsave**

```
ggplot2::ggsave(  
  output_plot,  
  filename = "FILENAME.png")
```

add plot object
add filename

Alternatively, storing the object also works this way:

```
png("FILENAME.png")  
  
output_plot  
  
dev.off()
```

R Syntax Comparison :: CHEAT SHEET

Dollar sign syntax

```
goal(data$x, data$y)
```

SUMMARY STATISTICS:

one continuous variable:

```
mean(mtcars$mpg)
```

one categorical variable:

```
table(mtcars$cyl)
```

two categorical variables:

```
table(mtcars$cyl, mtcars$am)
```

one continuous, one categorical:

```
mean(mtcars$mpg [mtcars$cyl==4])
```

```
mean(mtcars$mpg [mtcars$cyl==6])
```

```
mean(mtcars$mpg [mtcars$cyl==8])
```

PLOTTING:

one continuous variable:

```
hist(mtcars$disp)
```

```
boxplot(mtcars$disp)
```

one categorical variable:

```
barplot(table(mtcars$cyl))
```

two continuous variables:

```
plot(mtcars$disp, mtcars$mpg)
```

two categorical variables:

```
mosaicplot(table(mtcars$am, mtcars$cyl))
```

one continuous, one categorical:

```
histogram(mtcars$disp[mtcars$cyl==4])
```

```
histogram(mtcars$disp[mtcars$cyl==6])
```

```
histogram(mtcars$disp[mtcars$cyl==8])
```

```
boxplot(mtcars$disp[mtcars$cyl==4])
boxplot(mtcars$disp[mtcars$cyl==6])
boxplot(mtcars$disp[mtcars$cyl==8])
```

WRANGLING:

subsetting:

```
mtcars[mtcars$mpg>30, ]
```

making a new variable:

```
mtcars$efficient[mtcars$mpg>30] <- TRUE
```

```
mtcars$efficient[mtcars$mpg<30] <- FALSE
```

Formula syntax

```
goal(y~x|z, data=data, group=w)
```

SUMMARY STATISTICS:

one continuous variable:

```
mosaic::mean(~mpg, data=mtcars)
```

one categorical variable:

```
mosaic::tally(~cyl, data=mtcars)
```

two categorical variables:

```
mosaic::tally(cyl~am, data=mtcars)
```

one continuous, one categorical:

```
mosaic::mean(mpg~cyl, data=mtcars)
```

tilde

PLOTTING:

one continuous variable:

```
lattice::histogram(~disp, data=mtcars)
```

```
lattice::bwplot(~disp, data=mtcars)
```

one categorical variable:

```
mosaic::bargraph(~cyl, data=mtcars)
```

two continuous variables:

```
lattice::xyplot(mpg~disp, data=mtcars)
```

two categorical variables:

```
mosaic::bargraph(~am, data=mtcars, group=cyl)
```

one continuous, one categorical:

```
lattice::histogram(~disp|cyl, data=mtcars)
```

```
lattice::bwplot(cyl~disp, data=mtcars)
```

The variety of R syntaxes give you many ways to “say” the same thing

read across the cheatsheet to see how different syntaxes approach the same problem

Tidyverse syntax

```
data %>% goal(x)
```

SUMMARY STATISTICS:

one continuous variable:

```
mtcars %>% dplyr::summarize(mean(mpg))
```

one categorical variable:

```
mtcars %>% dplyr::group_by(cyl) %>%  
dplyr::summarize(n())
```

two categorical variables:

```
mtcars %>% dplyr::group_by(cyl, am) %>%  
dplyr::summarize(n())
```

one continuous, one categorical:

```
mtcars %>% dplyr::group_by(cyl) %>%  
dplyr::summarize(mean(mpg))
```

the pipe

PLOTTING:

one continuous variable:

```
ggplot2::qplot(x=mpg, data=mtcars, geom = "histogram")
```

```
ggplot2::qplot(y=disp, x=1, data=mtcars, geom="boxplot")
```

one categorical variable:

```
ggplot2::qplot(x=cyl, data=mtcars, geom="bar")
```

two continuous variables:

```
ggplot2::qplot(x=disp, y=mpg, data=mtcars, geom="point")
```

two categorical variables:

```
ggplot2::qplot(x=factor(cyl), data=mtcars, geom="bar") +  
facet_grid(.~am)
```

one continuous, one categorical:

```
ggplot2::qplot(x=disp, data=mtcars, geom = "histogram") +  
facet_grid(.~cyl)
```

```
ggplot2::qplot(y=disp, x=factor(cyl), data=mtcars,  
geom="boxplot")
```

WRANGLING:

subsetting:

```
mtcars %>% dplyr::filter(mpg>30)
```

making a new variable:

```
mtcars <- mtcars %>%  
dplyr::mutate(efficient = if_else(mpg>30, TRUE, FALSE))
```

R Syntax Comparison :: CHEAT SHEET

Syntax is the set of rules that govern what code works and doesn't work in a programming language. Most programming languages offer one standardized syntax, but R allows package developers to specify their own syntax. As a result, there is a large variety of (equally valid) R syntaxes.

The three most prevalent R syntaxes are:

1. The **dollar sign syntax**, sometimes called **base R syntax**, expected by most base R functions. It is characterized by the use of `dataset$variablename`, and is also associated with square bracket subsetting, as in `dataset[1, 2]`. Almost all R functions will accept things passed to them in dollar sign syntax.
2. The **formula syntax**, used by modeling functions like `lm()`, lattice graphics, and `mosaic` summary statistics. It uses the tilde (~) to connect a response variable and one (or many) predictors. Many base R functions will accept formula syntax.
3. The **tidyverse syntax** used by `dplyr`, `tidyverse`, and more. These functions expect data to be the first argument, which allows them to work with the "pipe" (%>%) from the `magrittr` package. Typically, `ggplot2` is thought of as part of the tidyverse, although it has its own flavor of the syntax using plus signs (+) to string pieces together. `ggplot2` author Hadley Wickham has said the package would have had different syntax if he had written it after learning about the pipe.

Educators often try to teach within one unified syntax, but most R programmers use some combination of all the syntaxes.

Internet research tip:

If you are searching on google, StackOverflow, or another favorite online source and see code in a syntax you don't recognize:

- Check to see if the code is using one of the three common syntaxes listed on this cheatsheet
- Try your search again, using a keyword from the syntax name ("tidyverse") or a relevant package ("mosaic")



Sometimes particular syntaxes work, but are considered dangerous to use, because they are so easy to get wrong. For example, passing variable names without assigning them to a named argument.

Even more ways to say the same thing

Even within one syntax, there are often variations that are equally valid. As a case study, let's look at the `ggplot2` syntax. `ggplot2` is the plotting package that lives within the `tidyverse`. If you read down this column, all the code here produces the same graphic.

quickplot

`qplot()` stands for quickplot, and allows you to make quick plots. It doesn't have the full power of `ggplot2`, and it uses a slightly different syntax than the rest of the package.

```
ggplot2::qplot(x=disp, y=mpg, data=mtcars, geom="point")
```

```
ggplot2::qplot(x=disp, y=mpg, data=mtcars) 
```

```
ggplot2::qplot(disp, mpg, data=mtcars)  
```

read down this column for many pieces of code in one syntax that look different but produce the same graphic

ggplot

To unlock the power of `ggplot2`, you need to use the `ggplot()` function (which sets up a plotting region) and add geoms to the plot.

```
ggplot2::ggplot(mtcars) +  
  geom_point(aes(x=disp, y=mpg))
```

```
ggplot2::ggplot(data=mtcars) +  
  geom_point(mapping=aes(x=disp, y=mpg))
```

plus adds layers

```
ggplot2::ggplot(mtcars, aes(x=disp, y=mpg)) +  
  geom_point()
```

```
ggplot2::ggplot(mtcars, aes(x=disp)) +  
  geom_point(aes(y=mpg))
```

ggformula

The "third and a half way" to use the formula syntax, but get `ggplot2`-style graphics

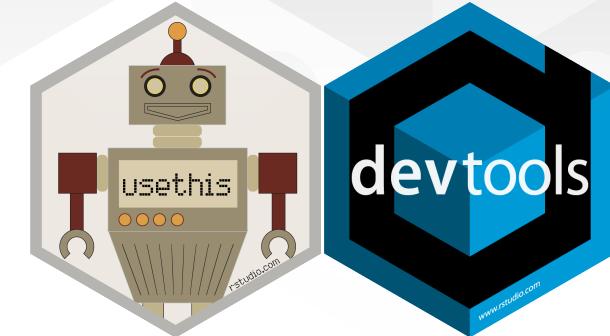
```
ggformula::gf_point(mpg~disp, data= mtcars)
```

formulas in base plots

Base R plots will also take the formula syntax, although it's not as commonly used

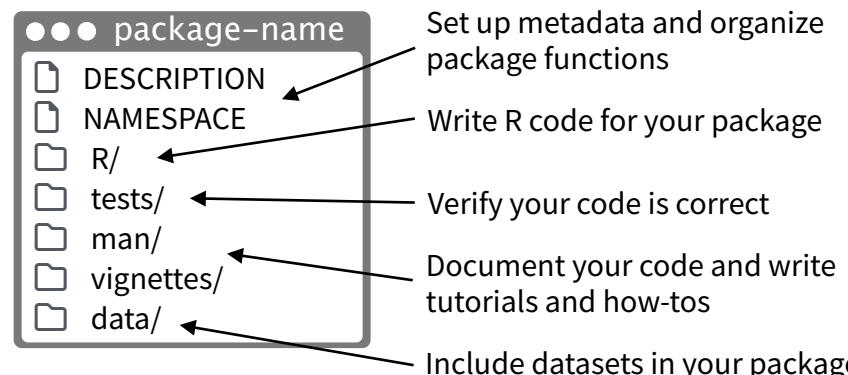
```
plot(mpg~disp, data=mtcars)
```

Package Development :: CHEAT SHEET



Package Structure

A package is a convention for organizing files into directories. This cheat sheet shows how to work with the 7 most common parts of an R package:



There are multiple packages useful to package development, including **usethis** which handily automates many of the more repetitive tasks. Install and load **devtools**, which wraps together several of these packages to access everything in one step.

Getting Started

Once per machine:

- Get set up with **use_r_profile()** so **devtools** is always loaded in interactive R sessions

```
if (interactive()) {  
  require("devtools", quietly = TRUE)  
  # automatically attaches usethis  
}
```

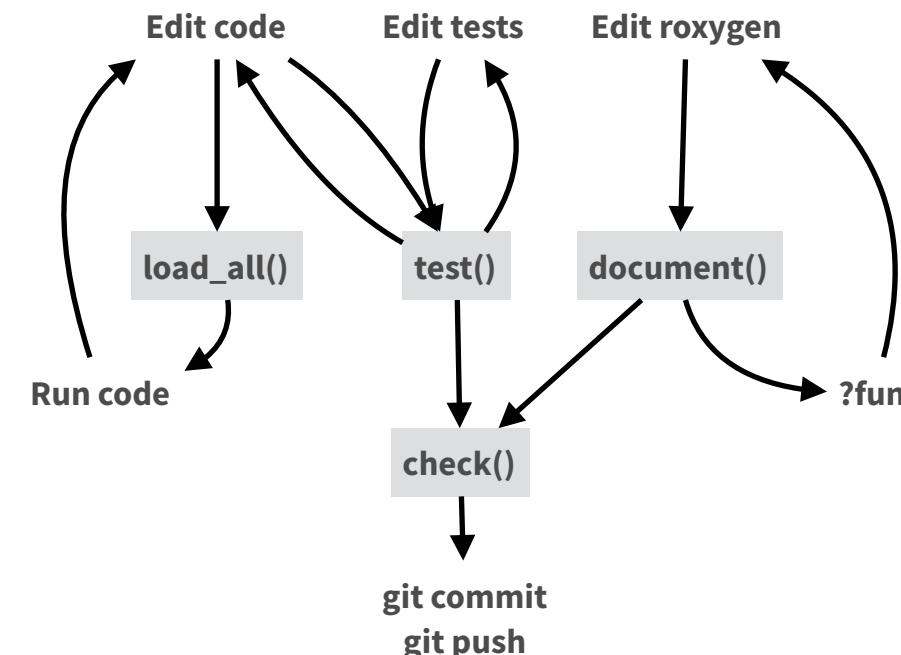
- create_github_token()** — Set up GitHub credentials
- git_vaccinate()** — Ignores common special files

Once per package:

- create_package()** — Create a project with package scaffolding
- use_git()** — Activate git
- use_github()** — Connect to GitHub
- use_github_action()** — Set up automated package checks

Having problems with git? Get a situation report with **git_sitrep()**.

Workflow



- load_all()** (Ctrl/Cmd + Shift + L) — Load code
- document()** (Ctrl/Cmd + Shift + D) — Rebuild docs and NAMESPACE
- test()** (Ctrl/Cmd + Shift + T) — Run tests
- check()** (Ctrl/Cmd + Shift + E) — Check complete package

R/

All of the R code in your package goes in **R/**. A package with just an R/ directory is still a very useful package.

- Create a new package project with **create_package("path/to/name")**.
- Create R files with **use_r("file-name")**.

- Follow the tidyverse style guide at style.tidyverse.org
- Click on a function and press **F2** to go to its definition
- Find a function or file with **Ctrl + .**

DESCRIPTION

The **DESCRIPTION** file describes your work, sets up how your package will work with other packages, and applies a license.

- Pick a license with **use_mit_license()**, **use_gpl3_license()**, **use_proprietary_license()**.
- Add packages that you need with **use_package()**.

Import packages that your package requires to work. R will install them when it installs your package.

use_package(x, type = "imports")

Suggest packages that developers of your package need. Users can install or not, as they like.

use_package(x, type = "suggests")

NAMESPACE

The **NAMESPACE** file helps you make your package self-contained: it won't interfere with other packages, and other packages won't interfere with it.

- Export functions for users by placing **@export** in their roxygen comments.
- Use objects from other packages with **package::object** or **@importFrom package object** (recommended) or **@import package** (use with caution).
- Call **document()** to generate NAMESPACE and **load_all()** to reload.

DESCRIPTION

Makes **packages** available
Mandatory
use_package()

NAMESPACE

Makes **function** available
Optional (can use `::` instead)
use_import_from()

man/

The documentation will become the help pages in your package.

- Document each function with a roxygen block above its definition in R/. In RStudio, Code > Insert Roxygen Skeleton helps (Ctrl/Cmd + Alt + Shift + R).
- Document each dataset with roxygen block above the name of the dataset in quotes.
- Document the package with `use_package_doc()`.
- Build documentation in `man/` from Roxygen blocks with `document()`.

vignettes/

- Create a vignette that is included with your package with `use_vignette()`.
- Create an article that only appears on the website with `use_article()`.
- Write the body of your vignettes in R Markdown.

Websites with pkgdown



- Use GitHub and `use_pkdown_github_pages()` to set up pkgdown and configures an automated workflow using GitHub Actions and Pages.
- If you're not using GitHub, call `use_pkdown()` to configure pkgdown. Then build locally with `pkdown::build_site()`.

tests/



- Set up test infrastructure with `use_testthat()`.
- Create a test file with `use_test()`.
- Write tests with `test_that()` and `expect_()`.
- Run all tests with `test()` and run tests for current file with `test_active_file()`.
- See coverage of all files with `test_coverage()` and see coverage of current file with `test_coverage_active_file()`.

ROXYGEN2

The **roxygen2** package lets you write documentation inline in your .R files with shorthand syntax.

- Add roxygen documentation as comments beginning with `#'`.
- Place a roxygen `@` tag (right) after `#'` to supply a specific section of documentation.
- Untagged paragraphs will be used to generate a title, description, and details section (in that order).

```
#' Add together two numbers
#'
#' @param x A number.
#' @param y A number.
#' @returns The sum of `x` and `y`.
#' @export
#' @examples
#' add(1, 1)
add <- function(x, y) {
  x + y
}
```



COMMON ROXYGEN TAGS

<code>@description</code>	<code>@family</code>	<code>@returns</code>
<code>@examples</code>	<code>@inheritParams</code>	<code>@seealso</code>
<code>@examplesIf</code>	<code>@param</code>	
<code>@export</code>	<code>@rdname</code>	

README.Rmd + NEWS.md

- Create a README and NEWS markdown files with `use_readme_rmd()` and `use_news_md()`.

Expect statement

`expect_equal()`

Is equal? (within numerical tolerance)

`expect_error()`

Throws specified error?

`expect_snapshot()`

Output is unchanged?

```
test_that("Math works", {
  expect_equal(1 + 1, 2)
  expect_equal(1 + 2, 3)
  expect_equal(1 + 3, 4)
})
```

Tests

data/

- Record how a data set was prepared as an R script and save that script to `data/` with `use_data_raw()`.
- Save a prepared data object to `data/` with `use_data()`.

Package States

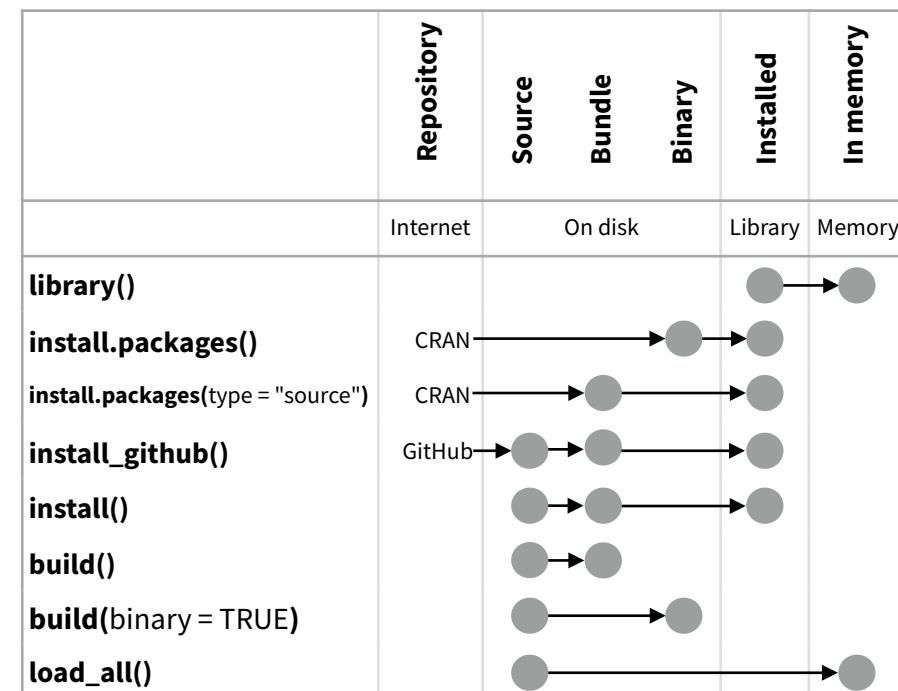
The contents of a package can be stored on disk as a:

- **source** - a directory with sub-directories (as shown in Package structure)
- **bundle** - a single compressed file (.tar.gz)
- **binary** - a single compressed file optimized for a specific OS

Packages exist in those states locally or remotely, e.g. on CRAN or on GitHub.

From those states, a package can be installed into an R library and then loaded into memory for use during an R session.

Use the functions below to move between these states.



Visit r-pkgs.org to learn much more about writing and publishing packages for R.

Apply functions with purrr :: CHEAT SHEET



Map Functions

ONE LIST

map(.x, .f, ...) Apply a function to each element of a list or vector, and return a list.
`x <- list(a = 1:10, b = 11:20, c = 21:30)
l1 <- list(x = c("a", "b"), y = c("c", "d"))
map(l1, sort, decreasing = TRUE)`



map_dbl(.x, .f, ...)
Return a double vector.
`map_dbl(x, mean)`

map_int(.x, .f, ...)
Return an integer vector.
`map_int(x, length)`

map_chr(.x, .f, ...)
Return a character vector.
`map_chr(l1, paste, collapse = "")`

map_lgl(.x, .f, ...)
Return a logical vector.
`map_lgl(x, is.integer)`

map_dfc(.x, .f, ...)
Return a data frame created by column-binding.
`map_dfc(l1, rep, 3)`

map_dfr(.x, .f, ..., .id = NULL)
Return a data frame created by row-binding.
`map_dfr(x, summary)`

walk(.x, .f, ...) Trigger side effects, return invisibly.
`walk(x, print)`

Function Shortcuts

Use `~.` with functions like **map()** that have single arguments.

`map(l, ~ . + 2)`
becomes
`map(l, function(x) x + 2)`

Use `~ .x .y` with functions like **map2()** that have two arguments.

`map2(l, p, ~ .x +.y)`
becomes
`map2(l, p, function(l, p) l + p)`

Use `~ ..1 ..2 ..3` etc with functions like **pmap()** that have many arguments.

`pmap(list(a, b, c), ~ ..3 + ..1 - ..2)`
becomes
`pmap(list(a, b, c), function(a, b, c) c + a - b)`

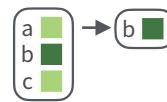
Use `~ .x .y` with functions like **imap()**. `.x` will get the list value and `.y` will get the index, or name if available.

`imap(list(a, b, c), ~ paste0(.y, ": ", .x))`
outputs "index: value" for each item

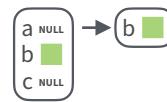


Work with Lists

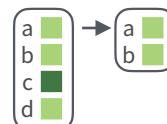
Filter



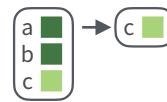
keep(.x, .p, ...)
Select elements that pass a logical test.
Conversely, **discard()**.
`keep(x, is.numeric)`



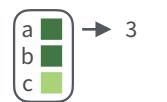
compact(.x, .p = identity)
Drop empty elements.
`compact(x)`



head_while(.x, .p, ...)
Return head elements until one does not pass.
Also **tail_while()**.
`head_while(x, is.character)`



detect(.x, .f, ..., dir = c("forward", "backward"), .right = NULL, .default = NULL)
Find first element to pass.
`detect(x, is.character)`



detect_index(.x, .f, ..., dir = c("forward", "backward"), .right = NULL)
Find index of first element to pass.
`detect_index(x, is.character)`



every(.x, .p, ...)
Do all elements pass a test?
`every(x, is.character)`



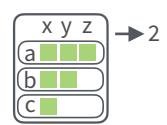
some(.x, .p, ...)
Do some elements pass a test?
`some(x, is.character)`



none(.x, .p, ...)
Do no elements pass a test?
`none(x, is.character)`

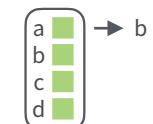


has_element(.x, .y)
Does a list contain an element?
`has_element(x, "foo")`

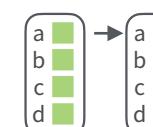


pluck_depth(x)
Return depth (number of levels of indexes).
`pluck_depth(x)`

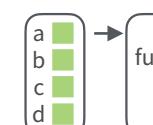
Index



pluck(.x, ..., .default=NULL)
Select an element by name or index. Also **attr_getter()** and **chuck()**.
`pluck(x, "b")`
`x > pluck("b")`

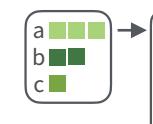


assign_in(x, where, value)
Assign a value to a location using pluck selection.
`assign_in(x, "b", 5)`
`x > assign_in("b", 5)`



modify_in(.x, .where, .f)
Apply a function to a value at a selected location.
`modify_in(x, "b", abs)`
`x > modify_in("b", abs)`

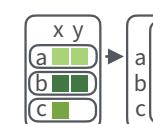
Reshape



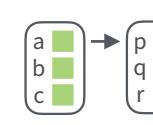
flatten(.x) Remove a level of indexes from a list.
Also **flatten_chr()** etc.
`flatten(x)`



array_tree(array, margin = NULL) Turn array into list.
Also **array_branch()**.
`z <- array(1:12, c(2,2,2))`
`array_tree(x, margin = 3)`

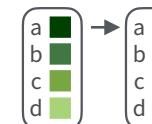


transpose(.l, .names = NULL)
Transposes the index order in a multi-level list.
`transpose(x)`

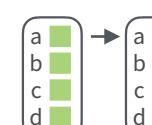


set_names(x, nm = x)
Set the names of a vector/list directly or with a function.
`set_names(x, c("p", "q", "r"))`
`set_names(x, tolower)`

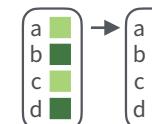
Modify



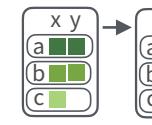
modify(.x, .f, ...) Apply a function to each element. Also **modify2()**, and **imodify()**.
`modify(x, ~.+ 2)`



modify_at(.x, .at, .f, ...) Apply a function to selected elements. Also **map_at()**.
`modify_at(x, "b", ~.+ 2)`



modify_if(.x, .p, .f, ...) Apply a function to elements that pass a test. Also **map_if()**.
`modify_if(x, is.numeric, ~.+ 2)`



modify_depth(.x, .depth, .f, ...) Apply function to each element at a given level of a list. Also **map_depth()**.
`modify_depth(x, 1, ~.+ 2)`

List-Columns

List-columns are columns of a data frame where each element is a list or vector instead of an atomic value. Columns can also be lists of data frames. See **tidyverse** for more about nested data and list columns.

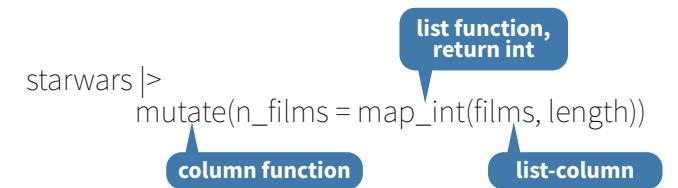
WORK WITH LIST-COLUMNS

Manipulate list-columns like any other kind of column, using **dplyr** functions like **mutate()** and **transmute()**. Because each element is a list, use **map functions** within a column function to manipulate each element.

map(), **map2()**, or **pmap()** return lists and will create new list-columns.



Suffixed map functions like **map_int()** return an atomic data type and will simplify list-columns into regular columns.



R For Data Science

data.table Cheat Sheet

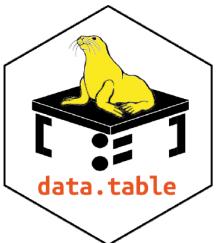
[Learn data.table online at www.DataCamp.com](http://www.DataCamp.com)

data.table

data.table is an R package that provides a high-performance version of base R's data.frame with syntax and feature enhancements for ease of use, convenience and programming speed.

Load the package:

```
> library(data.table)
```



Creating A data.table

```
> set.seed(45L) #Create a data.table and call it DT
> DT <- data.table(V1=c(1L,2L),
V2=LETTERS[1:3],
V3=rnorm(4),4),
V4=1:12)
```

Subsetting Rows Using i

```
> DT[3:5,] #Select 3rd to 5th row
> DT[3:5] #Select 3rd to 5th row
> DT[V2=="A"] #Select all rows that have value A in column V2
> DT[V2 %in% c("A","C")] #Select all rows that have value A or C in column V2
```

Manipulating on Columns in j

```
> DT[,V2] Return V2 as a vector
[1] "A" "B" "C" "A" "B" "C" ...
> DT[,.(V2,V3)] #Return V2 and V3 as a data.table
> DT[,sum(V1)] #Return the sum of all elements of V1 in a vector
[1] 18
#Return the sum of all elements of V1 and the std. dev. of V3 in a data.table
> DT[,.(sum(V1),sd(V3))]
V1 V2
1: 18 0.4546055
> DT[,.(Aggregate=sum(V1), #The same as the above, with new names
Sd.V3=sd(V3))]
Aggregate Sd.V3
1: 18 0.4546055
#Select column V2 and compute std. dev. of V3, which returns a single value & gets recycled
> DT[,(V1,Sd.V3=sd(V3))]
> DT[,.(print(V2), #Print column V2 and plot V3
plot(V3),
NULL)]
```

Chaining

```
> DT <- DT[,.(V4.Sum=sum(V4)), by=V1] #Calculate sum of V4, grouped by V1
V1 V4.Sum
1: 1 36
2: 2 42
> DT[V4.Sum>40] #Select that group of which the sum is >40
> DT[,(V4.Sum=sum(V4)), #Select that group of which the sum is >40 (chaining)
by=V1][V4.Sum>40]
V1 V4.Sum
1: 2 42
> DT[,(V4.Sum=sum(V4)), by=V1][order(-V1)] Calculate sum of V4, grouped by ordered on V1
V1 V4.Sum
1: 2 42
2: 1 36
```

> Doing j by Group

```
> DT[,(V4.Sum=sum(V4)), by=V1] #Calculate sum of V4 for every group in V1 Sum
V1 V4.
1: 1 36
2: 2 42
> DT[,(V4.Sum=sum(V4)), by=(V1,V2)] #Calculate sum of V4 for every group in V1 and V2
> DT[,(V4.Sum=sum(V4)), by=sign(V1-1)] #Calculate sum of V4 for every group in sign(V1-1)
sign V4.Sum
1: 0 36
2: 1 42
#The same as the above, with new name for the variable you're grouping by
> DT[,(V4.Sum=sum(V4)), by=(V1.01=sign(V1-1))]
#Calculate sum of V4 for every group in V1 after subsetting on the first 5 rows
> DT[1:5,(V4.Sum=sum(V4)), by=V1]
> DT[,N,by=V1] #Count number of rows for every group in V1

General form: DT[i,j,by] "Take DT, subset rows using i, then calculate j grouped by by"
```

> Advanced Data Table Operations

```
> DT[,N-1] #Return the penultimate row of the DT
> DT[,N] #Return the number of rows
> DT[,.(V2,V3)] #Return V2 and V3 as a data.table
> DT[,list(V2,V3)] #Return V2 and V3 as a data.table
#Return the result of j, grouped by all possible combinations of groups specified in by
> DT[,mean(V3),by=(V1,V2)]
V1 V2 V1
1: 1 A 0.4053
2: 1 B 0.4053
3: 1 C 0.4053
4: 2 A -0.6443
5: 2 B -0.6443
6: 2 C -0.6443
```

.SD & .SDcols

```
> DT[,print(.SD),by=V2] #Look at what .SD contains
> DT[,SD[c(1,N)],by=V2] #Select the first and last row grouped by V2
> DT[,lapply(.SD,sum),by=V2] #Calculate sum of columns in .SD grouped by V2
> DT[,lapply(.SD,sum),by=V2, #Calculate sum of V3 and V4 in .SD grouped by V2
.SDcols=c("V3","V4")]
V2 V3 V4
1: A -0.478 22
2: B -0.478 26
3: C -0.478 30
> DT[,lapply(.SD,sum),by=V2, #Calculate sum of V3 and V4 in .SD grouped by V2
.SDcols=paste0("V",3:4)]
```

> Indexing And Keys

```
> setkey(DT,V2) #A key is set on V2; output is returned invisibly
> DT["A"] #Return all rows where the key column (set to V2) has the value A
V1 V2 V3 V4
1: 1 A -0.2392 1
2: 2 A -1.6148 4
3: 1 A 1.0498 7
4: 2 A 0.3262 10
> DT[c("A","C")] #Return all rows where the key column (V2) has value A or C
> DT["A",mult="first"] #Return first row of all rows that match value A in key column V2
> DT["A",mult="last"] #Return last row of all rows that match value A in key column V2
> DT[c("A","D")] #Return all rows where key column V2 has value A or D
V1 V2 V3 V4
1: 1 A -0.2392 1
2: 2 A -1.6148 4
3: 1 A 1.0498 7
4: 2 A 0.3262 10
5: NA D NA NA
> DT[c("A","D"),nomatch=0] #Return all rows where key column V2 has value A or D
V1 V2 V3 V4
1: 1 A -0.2392 1
2: 2 A -1.6148 4
3: 1 A 1.0498 7
4: 2 A 0.3262 10
#Return total sum of V4, for rows of key column V2 that have values A or C
> DT[c("A","C"),sum(V4)]
#Return sum of column V4 for rows of V2 that have value A, sum(V4),
and another sum for rows of V2 that have value C
> DT[c("A","C"), by=.EACH]
V2 V1
1: A 22
2: C 30
> setkey(DT,V1,V2) #Sort by V1 and then by V2 within each group of V1 (invisible)
#Select rows that have value 2 for the first key (V1) &
the value C for the second key (V2)
> DT[(,2,"C")]
V1 V2 V3 V4
1: 2 C 0.3262 6
2: 2 C -1.6148 12
Select rows that have value 2 for the first key (V1) &
within those rows the value A or C for the second key (V2)
> DT[(,2,c("A","C"))]
V1 V2 V3 V4
1: 2 A -1.6148 4
2: 2 A 0.3262 10
3: 2 C 0.3262 6
4: 2 C -1.6148 12
```



Data import with the tidyverse :: CHEATSHEET

Read Tabular Data with readr

```
read_*(file, col_names = TRUE, col_types = NULL, col_select = NULL, id = NULL, locale, n_max = Inf,
skip = 0, na = c("", "NA"), guess_max = min(1000, n_max), show_col_types = TRUE) See ?read_delim
```

A B C	1 2 3	4 5 NA
A	B	C
1	2	3
4	5	NA

read_delim("file.txt", delim = "|") Read files with any delimiter. If no delimiter is specified, it will automatically guess.

To make file.txt, run: `write_file("A|B|C\n1|2|3\n4|5|NA", file = "file.txt")`

A,B,C	1,2,3	4,5,NA
A	B	C
1	2	3
4	5	NA

read_csv("file.csv") Read a comma delimited file with period decimal marks.

`write_file("A,B,C\n1,2,3\n4,5,NA", file = "file.csv")`

A;B;C	1;5;2;3	4;5;5;NA
A	B	C
1	2	3
4	5	NA

read_csv2("file2.csv") Read semicolon delimited files with comma decimal marks.

`write_file("A;B;C\n1,5;2;3\n4,5;5;NA", file = "file2.csv")`

A B C	1 2 3	4 5 NA
A	B	C
1	2	3
4	5	NA

read_tsv("file.tsv") Read a tab delimited file. Also **read_table()**.

read_fwf("file.tsv", fwf_widths(c(2, 2, NA))) Read a fixed width file.

`write_file("A\tB\tC\n1\t2\t3\n4\t5\tNA", file = "file.tsv")`

USEFUL READ ARGUMENTS

A	B	C
1	2	3
4	5	NA

No header

`read_csv("file.csv", col_names = FALSE)`

x	y	z
A	B	C
1	2	3
4	5	NA

Provide header

`read_csv("file.csv", col_names = c("x", "y", "z"))`



Read multiple files into a single table

`read_csv(c("f1.csv", "f2.csv", "f3.csv"), id = "origin_file")`

1	2	3
4	5	NA

Skip lines

`read_csv("file.csv", skip = 1)`

A	B	C
1	2	3

Read a subset of lines

`read_csv("file.csv", n_max = 1)`

A	B	C
NA	2	3
4	5	NA

Read values as missing

`read_csv("file.csv", na = c("1"))`

A;B;C	1;5;2;3;0	
A	B	C
1	5	2
4	5	NA

Specify decimal marks

`read_delim("file2.csv", locale = locale(decimal_mark = ","))`

One of the first steps of a project is to import outside data into R. Data is often stored in tabular formats, like csv files or spreadsheets.



The front page of this sheet shows how to import and save text files into R using **readr**.



The back page shows how to import spreadsheet data from Excel files using **readxl** or Google Sheets using **googlesheets4**.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files:

- **haven** - SPSS, Stata, and SAS files
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)
- **readr::read_lines()** - text data

Column Specification with readr

Column specifications define what data type each column of a file will be imported as. By default **readr** will generate a column spec when a file is read and output a summary.

spec(x) Extract the full column specification for the given imported data frame.

```
spec(x)
# cols(
#   age = col_integer(),
#   edu = col_character(),
#   earn = col_double()
# )
```

age is an integer

edu is a character

earn is a double (numeric)

COLUMN TYPES

Each column type has a function and corresponding string abbreviation.

- **col_logical()** - "l"
- **col_integer()** - "i"
- **col_double()** - "d"
- **col_number()** - "n"
- **col_character()** - "c"
- **col_factor(levels, ordered = FALSE)** - "f"
- **col_datetime(format = "")** - "T"
- **col_date(format = "")** - "D"
- **col_time(format = "")** - "t"
- **col_skip()** - "-", "_"
- **col_guess()** - "?"

DEFINE COLUMN SPECIFICATION

Set a default type

```
read_csv(
  file,
  col_type = list(.default = col_double())
)
```

Use column type or string abbreviation

```
read_csv(
  file,
  col_type = list(x = col_double(), y = "l", z = "_")
)
```

Use a single string of abbreviations

```
# col types: skip, guess, integer, logical, character
read_csv(
  file,
  col_type = "_?ilc"
)
```



Import Spreadsheets with readxl

READ EXCEL FILES

	A	B	C	D	E
1	x1	x2	x3	x4	x5
2	x		z	8	
3	y	7		9	10

```
read_excel(path, sheet = NULL, range = NULL)  
Read a .xls or .xlsx file based on the file extension.  
See front page for more read arguments. Also  
read_xls() and read_xlsx().  
read_excel("excel_file.xlsx")
```

READ SHEETS

A	B	C	D	E
s1	s2	s3		

s1	s2	s3
----	----	----

A	B	C	D	E
A	B	C	D	E
A	B	C	D	E

- To **read multiple sheets**:
1. Get a vector of sheet names from the file path.
 2. Set the vector names to be the sheet names.
 3. Use purrr::map() and purrr::list_rbind() to read multiple files into one data frame.

```
path <- "your_file_path.xlsx"  
path >  
  excel_sheets() |>  
  set_names() |>  
  map(read_excel, path = path) |>  
  list_rbind()
```

OTHER USEFUL EXCEL PACKAGES

For functions to write data to Excel files, see:

- **openxlsx**
- **writexl**

For working with non-tabular Excel data, see:

- **tidyxl**



with googlesheets4

READ SHEETS

	A	B	C	D	E
1	x1	x2	x3	x4	x5
2	x		z	8	
3	y	7		9	10

```
read_sheet(ss, sheet = NULL, range = NULL)  
Read a sheet from a URL, a Sheet ID, or a dribble  
from the googledrive package. See front page for  
more read arguments. Same as range_read().
```

SHEETS METADATA

URLs are in the form:
<https://docs.google.com/spreadsheets/d/>
SPREADSHEET_ID/edit#gid=**SHEET_ID**

gs4_get(ss) Get spreadsheet meta data.

gs4_find(...) Get data on all spreadsheet files.

sheet_properties(ss) Get a tibble of properties
for each worksheet. Also **sheet_names()**.

WRITE SHEETS

1	x	4
2	y	5
3	z	6

1	A	B	C	D
2				

x1	x2	x3
2	y	5
3	z	6

write_sheet(data, ss = NULL, sheet = NULL)
Write a data frame into a new or existing Sheet.

gs4_create(name, ..., sheets = NULL) Create a new Sheet with a vector of names, a data frame, or a (named) list of data frames.

sheet_append(ss, data, sheet = 1) Add rows to the end of a worksheet.



GOOGLESHEETS4 COLUMN SPECIFICATION

Column specifications define what data type each column of a file will be imported as.

Use the **col_types** argument of **read_sheet()**/**range_read()** to set the column specification.

Guess column types

To guess a column type, **read_excel()** looks at the first 1000 rows of data. Increase with the **guess_max** argument.
read_excel(path, guess_max = Inf)

Set all columns to same type, e.g. character
read_sheet(path, col_types = "c")

Set each column individually

col types: skip, guess, integer, logical, character
read_sheets(ss, col_types = "?ilc")

COLUMN TYPES

I	n	c	D	L
TRUE	2	hello	1947-01-08	hello
FALSE	3.45	world	1956-10-21	1

- skip - "_" or "-"
- guess - "?"
- logical - "l"
- integer - "i"
- double - "d"
- numeric - "n"
- date - "D"
- datetime - "T"
- character - "c"
- list-column - "L"
- cell - "C" Returns list of raw cell data.

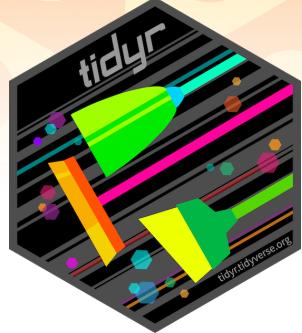
Use list for columns that include multiple data types. See **tidyxl** and **purrr** for list-column data.

FILE LEVEL OPERATIONS

googlesheets4 also offers ways to modify other aspects of Sheets (e.g. freeze rows, set column width, manage (work)sheets). Go to googlesheets4.tidyverse.org to read more.

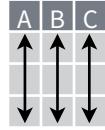
For whole-file operations (e.g. renaming, sharing, placing within a folder), see the tidyverse package **googledrive** at googledrive.tidyverse.org.

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**

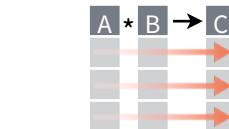
&



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"))`

Both make this tibble

A tibble: 3 × 2
`x` <int> <chr>
 1 1 a
 2 2 b
 3 3 c

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	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K



country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

table2

country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M
C	1999	212K	1T
C	2000	213K	1T

Split Cells

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

table5

country	century	year
A	19	99
A	20	00
B	19	99
B	20	00



country	year
A	1999
A	2000
B	1999
B	2000

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M



country	year	rate
A	1999	0.7K
A	1999	19M
A	2000	2K
A	2000	20M
B	1999	37K
B	1999	172M
B	2000	80K
B	2000	174M

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")
```

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)
```

Expand Tables

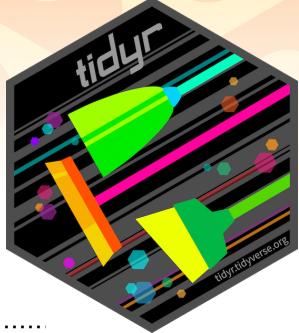
Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

x	x1	x2	x3
A	1	3	
B	1	4	
B	2	3	

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)
```

x	x1	x2	x3
A	1	3	
B	1	4	
B	2	3	
			NA



Nested Data

A **nested data frame** stores individual tables as a list-column of data frames within a larger organizing data frame. List-columns can also be lists of vectors or lists of varying data types.

Use a nested data frame to:

- Preserve relationships between observations and subsets of data. Preserve the type of the variables being nested (factors and datetimes aren't coerced to character).
- Manipulate many sub-tables at once with **purrr** functions like `map()`, `map2()`, or `pmap()` or with **dplyr** `rowwise()` grouping.

CREATE NESTED DATA

nest(data, ...) Moves groups of cells into a list-column of a data frame. Use alone or with `dplyr::group_by()`:

1. Group the data frame with `group_by()` and use `nest()` to move the groups into a list-column.

```
n_storms <- storms %>%
  group_by(name) %>%
  nest()
```

2. Use `nest(new_col = c(x, y))` to specify the columns to group using `dplyr::select()` syntax.

```
n_storms <- storms %>%
  nest(data = c(year:long))
```

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

Index list-columns with `[[[]]]`. `n_storms$data[[1]]`

CREATE TIBBLES WITH LIST-COLUMNS

tibble::tribble(...) Makes list-columns when needed.

```
tibble(~max, ~seq,
      3, 1:3,
      4, 1:4,
      5, 1:5)
```

max	seq
3	<int [3]>
4	<int [4]>
5	<int [5]>

tibble::tibble(...) Saves list input as list-columns.

```
tibble(max = c(3, 4, 5), seq = list(1:3, 1:4, 1:5))
```

tibble::enframe(x, name="name", value="value")

Converts multi-level list to a tibble with list-cols.
`enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')`

OUTPUT LIST-COLUMNS FROM OTHER FUNCTIONS

dplyr::mutate(), transmute(), and summarise() will output list-columns if they return a list.

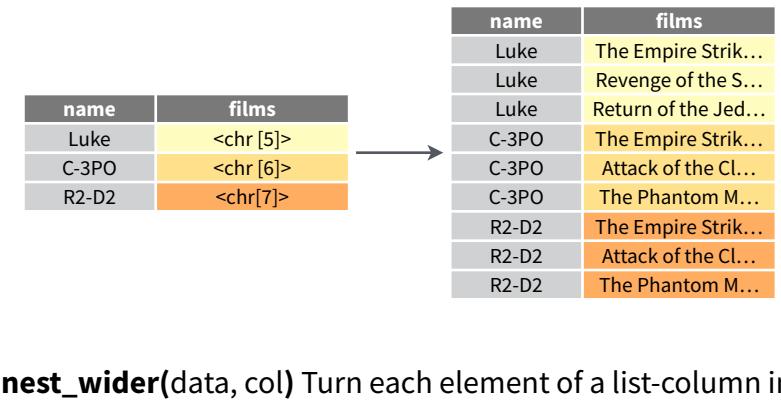
```
mtcars %>%
  group_by(cyl) %>%
  summarise(q = list(quantile(mpg)))
```

RESHAPE NESTED DATA

unnest(data, cols, ..., keep_empty = FALSE) Flatten nested columns back to regular columns. The inverse of `nest()`.
`n_storms %>% unnest(data)`

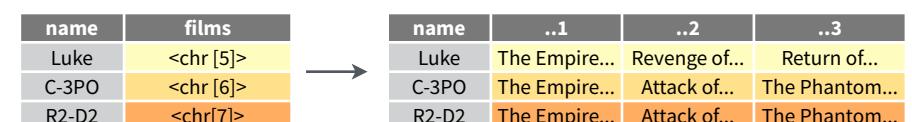
unnest_longer(data, col, values_to = NULL, indices_to = NULL)
Turn each element of a list-column into a row.

```
starwars %>%
  select(name, films) %>%
  unnest_longer(films)
```



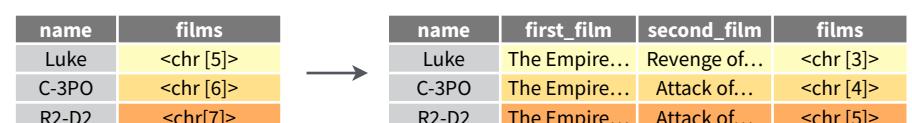
unnest_wider(data, col) Turn each element of a list-column into a regular column.

```
starwars %>%
  select(name, films) %>%
  unnest_wider(films)
```



hoist(.data, .col, ..., .remove = TRUE) Selectively pull list components out into their own top-level columns. Uses `purrr::pluck()` syntax for selecting from lists.

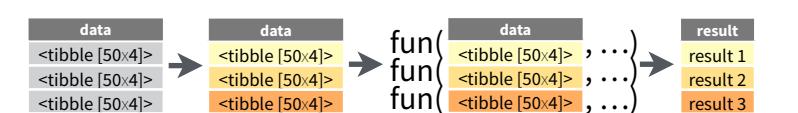
```
starwars %>%
  select(name, films) %>%
  hoist(films, first_film = 1, second_film = 2)
```



TRANSFORM NESTED DATA

A vectorized function takes a vector, transforms each element in parallel, and returns a vector of the same length. By themselves vectorized functions cannot work with lists, such as list-columns.

dplyr::rowwise(.data, ...) Group data so that each row is one group, and within the groups, elements of list-columns appear directly (accessed with `[]`, not as lists of length one. **When you use `rowwise()`, dplyr functions will seem to apply functions to list-columns in a vectorized fashion.**



Apply a function to a list-column and **create a new list-column**.

`n_storms %>% rowwise() %>% mutate(n = list(dim(data)))`

dim() returns two values per row
wrap with `list` to tell `mutate` to create a list-column

Apply a function to a list-column and **create a regular column**.

`n_storms %>% rowwise() %>% mutate(n = nrow(data))`

nrow() returns one integer per row

Collapse **multiple list-columns** into a single list-column.

`starwars %>% rowwise() %>% mutate(transport = list(append(vehicles, starships)))`

append() returns a list for each row, so col type must be list

Apply a function to **multiple list-columns**.

`starwars %>% rowwise() %>% mutate(n_transports = length(c(vehicles, starships)))`

length() returns one integer per row

See **purrr** package for more list functions.

Data transformation with dplyr :: CHEATSHEET



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**

pipes

$x |> f(y)$ becomes $f(x, y)$

Summarize 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).



summarize(.data, ...)
Compute table of summaries.
mtcars |> summarize(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()**, **add_count()**, **add_tally()**.
mtcars |> count(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) |>
summarize(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 tidyverse cheat sheet for list-column workflow.

starwars |>
rowwise() |>
mutate(film_count = length(films))

ungroup(x, ...) Returns ungrouped copy of table.
g_mtcars <- mtcars |> group_by(cyl)
ungroup(g_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.
mtcars |> filter(mpg > 20)



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



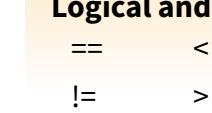
slice(.data, ..., .preserve = FALSE) Select rows by position.
mtcars |> slice(10:15)



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.
mtcars |> slice_sample(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.
mtcars |> slice_min(mpg, prop = 0.25)



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

Logical and boolean operators to use with filter()

<code>==</code>	<code><</code>	<code><=</code>	<code>is.na()</code>	<code>%in%</code>	<code> </code>	<code>xor()</code>
<code>!=</code>	<code>></code>	<code>>=</code>	<code>!is.na()</code>	<code>!</code>	<code>&</code>	

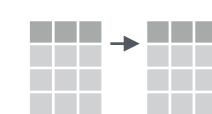
See [?base::Logic](#) and [?Comparison](#) for help.

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.
mtcars |> arrange(mpg)
mtcars |> arrange(desc(mpg))

ADD CASES



add_row(.data, ..., .before = NULL, .after = NULL)
Add one or more rows to a table.
cars |> add_row(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 values as a vector, by name or index.
mtcars |> pull(wt)



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



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

Use these helpers with select() and across()

e.g. mtcars |> select(mpg:cyl)

contains(match)

ends_with(match)

starts_with(match)

num_range(prefix, range)

all_of(x)/any_of(x, ..., vars)

matches(match)

; e.g., mpg:cyl

!, e.g., !gear

everything()

MANIPULATE MULTIPLE VARIABLES AT ONCE

df <- tibble(x_1 = c(1, 2), x_2 = c(3, 4), y = c(4, 5))



across(.cols, .funs, ..., .names = NULL) Summarize or mutate multiple columns in the same way.
df |> summarize(across(everything(), mean))



c_across(.cols) Compute across columns in row-wise data.
df |>
rowwise() |>
mutate(x_total = sum(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).



mutate(.data, ..., .keep = "all", .before = NULL, .after = NULL) Compute new column(s). Also **add_column()**.
mtcars |> mutate(gpm = 1 / mpg)
mtcars |> mutate(gpm = 1 / mpg, .keep = "none")



rename(.data, ...) Rename columns. Use **rename_with()** to rename with a function.
mtcars |> rename(miles_per_gallon = mpg)



Vectorized Functions

TO USE WITH MUTATE ()

mutate() applies 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()
cummax() - cumulative max()
dplyr::cummean() - cumulative mean()
cummin() - cumulative min()
cumprod() - cumulative prod()
cumsum() - cumulative sum()

RANKING

dplyr::cume_dist() - proportion of all values <= 1
dense_rank() - rank w ties = min, no gaps
min_rank() - rank with ties = min
ntile() - bins into n bins
percent_rank() - min_rank scaled to [0,1]
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 | mass > 200 ~ "large",
 species == "Droid" ~ "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
pmax() - element-wise max()
pmin() - element-wise min()

Summary Functions

TO USE WITH SUMMARIZE ()

summarize() 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
sum(!is.na()) - # of non-NAs

POSITION

mean() - mean, also mean(!is.na())
median() - median

LOGICAL

mean() - proportion of TRUEs
sum() - # of TRUEs

ORDER

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

RANK

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

SPREAD

IQR() - Inter-Quartile Range
mad() - median absolute deviation
sd() - standard deviation
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.

A B → C A B tibble::rownames_to_column()
1 a t 1 a Move row names into col.
2 b u 2 b a <- mtcars |>
3 c v 3 c rownames_to_column(var = "C")

A B C → A B tibble::column_to_rownames()
1 a t 1 a Move col into row names.
2 b u 2 b a |> column_to_rownames(var = "C")
3 c v 3 c

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

Combine Tables

COMBINE VARIABLES

x	y
A B C	E F G
a t 1	a t 3
b u 2	b u 2
c v 3	d w 1

=

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.

A B C D	left_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")
a t 1 3	Join matching values from y to x.

A B C D	right_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")
a t 1 3	Join matching values from x to y.

A B C D	inner_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")
a t 1 3	Join data. Retain only rows with matches.

A B C D	full_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na")
a t 1 3	Join data. Retain all values, all rows.

COLUMN MATCHING FOR JOINS

A B.x C.B.y D	Use by = c("col1", "col2", ...) to specify one or more common columns to match on.
a t 1 t 3	left_join(x, y, by = "A")

A.x B.x C.A.y B.y	Use a named vector, by = c("col1" = "col2") , to match on columns that have different names in each table.
a t 1 d w	left_join(x, y, by = c("C" = "D"))

A1 B1 C A2 B2	Use suffix to specify the suffix to give to unmatched columns that have the same name in both tables.
a t 1 d w	left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

COMBINE CASES

x	y
A B C	A B C
a t 1	a t 1
b u 2	b u 2
c v 3	d w 4

=

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.

x	y
A B C	A B C
a t 1	a t 3
b u 2	b u 2
c v 3	d w 4

=

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.

=

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.

A B C	nest_join(x, y, by = NULL, copy = FALSE, keep = FALSE, name = NULL, ...)
a t 1 <tibble [1x2]>	Join data, nesting matches from y in a single new data frame column.

SET OPERATIONS

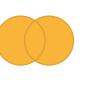
intersect(x, y, ...) Rows that appear in both x and y.



setdiff(x, y, ...) Rows that appear in x but not y.



union(x, y, ...) Rows that appear in x or y, duplicates removed). **union_all()** retains duplicates.



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

Data Wrangling with dplyr and tidyr

Cheat Sheet



Syntax - Helpful conventions for wrangling

dplyr::tbl_df(iris)

Converts data to `tbl` class. `tbl`'s are easier to examine than data frames. R displays only the data that fits onscreen:

```
Source: local data frame [150 x 5]
  Sepal.Length Sepal.Width Petal.Length
1          5.1        3.5         1.4
2          4.9        3.0         1.4
3          4.7        3.2         1.3
4          4.6        3.1         1.5
5          5.0        3.6         1.4
...
Variables not shown: Petal.Width (dbl), Species (fctr)
```

dplyr::glimpse(iris)

Information dense summary of `tbl` data.

utils::View(iris)

View data set in spreadsheet-like display (note capital V).

iris x					
<input type="button" value="Filter"/> <input type="button" value="Print"/> <input type="button" value="Copy"/> <input type="button" value="Save"/>					
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

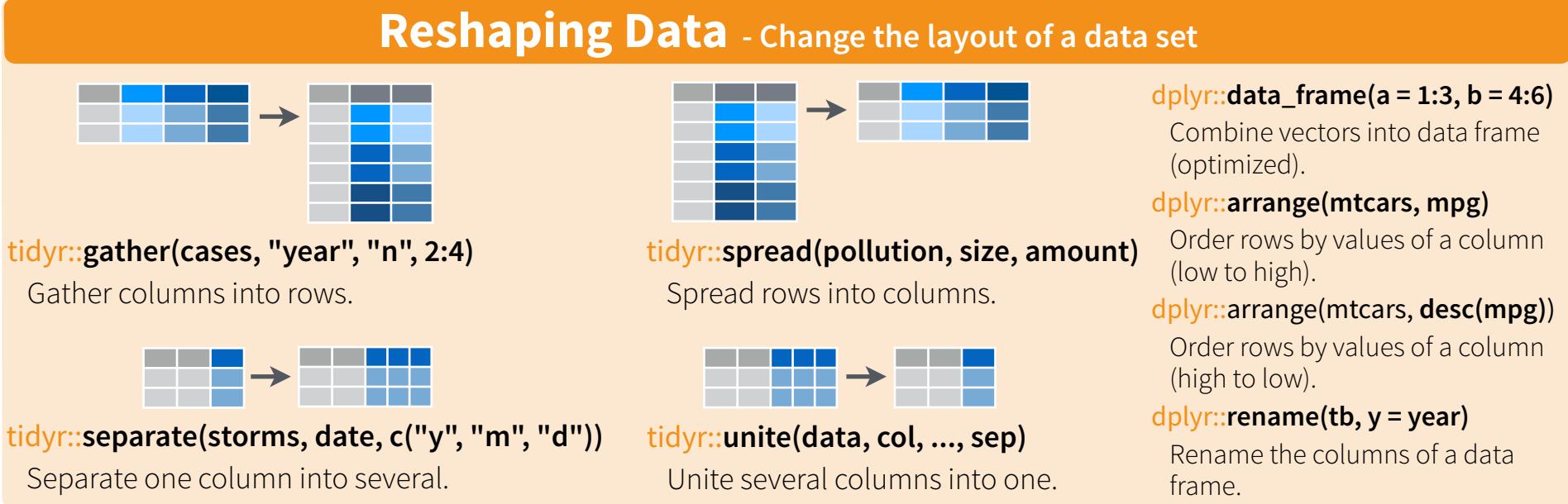
dplyr::%>%

Passes object on left hand side as first argument (or . argument) of function on righthand side.

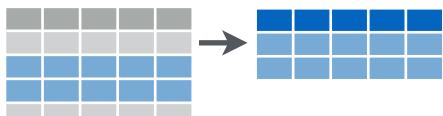
x %>% f(y) is the same as **f(x, y)**
y %>% f(x, ., z) is the same as **f(x, y, z)**

"Piping" with `%>%` makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(avg)
```



Subset Observations (Rows)



dplyr::filter(iris, Sepal.Length > 7)

Extract rows that meet logical criteria.

dplyr::distinct(iris)

Remove duplicate rows.

dplyr::sample_frac(iris, 0.5, replace = TRUE)

Randomly select fraction of rows.

dplyr::sample_n(iris, 10, replace = TRUE)

Randomly select n rows.

dplyr::slice(iris, 10:15)

Select rows by position.

dplyr::top_n(storms, 2, date)

Select and order top n entries (by group if grouped data).

Logic in R - ?Comparison, ?base::Logic			
<	Less than	!=	Not equal to
>	Greater than	%in%	Group membership
==	Equal to	is.na	Is NA
<=	Less than or equal to	!is.na	Is not NA
>=	Greater than or equal to	&, , !, xor, any, all	Boolean operators

Subset Variables (Columns)



dplyr::select(iris, Sepal.Width, Petal.Length, Species)

Select columns by name or helper function.

Helper functions for select - ?select

select(iris, contains("."))

Select columns whose name contains a character string.

select(iris, ends_with("Length"))

Select columns whose name ends with a character string.

select(iris, everything())

Select every column.

select(iris, matches(".t.))

Select columns whose name matches a regular expression.

select(iris, num_range("x", 1:5))

Select columns named x1, x2, x3, x4, x5.

select(iris, one_of(c("Species", "Genus")))

Select columns whose names are in a group of names.

select(iris, starts_with("Sepal"))

Select columns whose name starts with a character string.

select(iris, Sepal.Length:Petal.Width)

Select all columns between Sepal.Length and Petal.Width (inclusive).

select(iris, -Species)

Select all columns except Species.

Summarise Data



dplyr::summarise(iris, avg = mean(Sepal.Length))

Summarise data into single row of values.

dplyr::summarise_each(iris, funs(mean))

Apply summary function to each column.

dplyr::count(iris, Species, wt = Sepal.Length)

Count number of rows with each unique value of variable (with or without weights).



Summarise uses **summary functions**, functions that take a vector of values and return a single value, such as:

dplyr::first

First value of a vector.

dplyr::last

Last value of a vector.

dplyr::nth

Nth value of a vector.

dplyr::n

of values in a vector.

dplyr::n_distinct

of distinct values in a vector.

IQR

IQR of a vector.

min

Minimum value in a vector.

max

Maximum value in a vector.

mean

Mean value of a vector.

median

Median value of a vector.

var

Variance of a vector.

sd

Standard deviation of a vector.

Make New Variables



dplyr::mutate(iris, sepal = Sepal.Length + Sepal.Width)

Compute and append one or more new columns.

dplyr::mutate_each(iris, funs(min_rank))

Apply window function to each column.

dplyr::transmute(iris, sepal = Sepal.Length + Sepal.Width)

Compute one or more new columns. Drop original columns.



Mutate uses **window functions**, functions that take a vector of values and return another vector of values, such as:

dplyr::lead

Copy with values shifted by 1.

dplyr::lag

Copy with values lagged by 1.

dplyr::dense_rank

Ranks with no gaps.

dplyr::min_rank

Ranks. Ties get min rank.

dplyr::percent_rank

Ranks rescaled to [0, 1].

dplyr::row_number

Ranks. Ties got to first value.

dplyr::ntile

Bin vector into n buckets.

dplyr::between

Are values between a and b?

dplyr::cume_dist

Cumulative distribution.

dplyr::cumall

Cumulative **all**

dplyr::cumany

Cumulative **any**

dplyr::cummean

Cumulative **mean**

cumsum

Cumulative **sum**

cummax

Cumulative **max**

cummin

Cumulative **min**

cumprod

Cumulative **prod**

pmax

Element-wise **max**

pmin

Element-wise **min**

Group Data

dplyr::group_by(iris, Species)

Group data into rows with the same value of Species.

dplyr::ungroup(iris)

Remove grouping information from data frame.

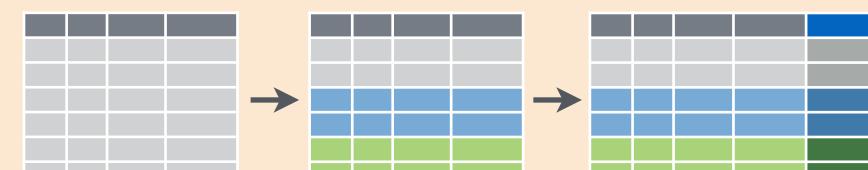
iris %>% group_by(Species) %>% summarise(...)

Compute separate summary row for each group.



iris %>% group_by(Species) %>% mutate(...)

Compute new variables by group.



Combine Data Sets

a	b		
x1	x2	x1	x3
A	1	A	T
B	2	B	F
C	3	D	T

$$+ =$$

Mutating Joins

x1	x2	x3
A	1	T
B	2	F
C	3	NA

x1	x3	x2
A	T	1
B	F	2
D	NA	T

x1	x2	x3
A	1	T
B	2	F
C	3	NA

Filtering Joins

x1	x2
A	1
B	2

x1	x2
C	3

dplyr::semi_join(a, b, by = "x1")

All rows in a that have a match in b.

dplyr::anti_join(a, b, by = "x1")

All rows in a that do not have a match in b.

y	z		
x1	x2	x1	x2
A	1	B	2
B	2	C	3
C	3	D	4

x1	x2
B	2
C	3
D	4

x1	x2
A	1

dplyr::intersect(y, z)

Rows that appear in both y and z.

dplyr::union(y, z)

Rows that appear in either or both y and z.

dplyr::setdiff(y, z)

Rows that appear in y but not z.

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

x1	x2	x1</

String manipulation with stringr :: CHEATSHEET



The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.

Detect Matches

	str_detect(string, pattern, negate = FALSE) Detect the presence of a pattern match in a string. Also str_like() . str_detect(fruit, "a")
	str_starts(string, pattern, negate = FALSE) Detect the presence of a pattern match at the beginning of a string. Also str_ends() . str_starts(fruit, "a")
	str_which(string, pattern, negate = FALSE) Find the indexes of strings that contain a pattern match. str_which(fruit, "a")
	str_locate(string, pattern) Locate the positions of pattern matches in a string. Also str_locate_all() . str_locate(fruit, "a")
	str_count(string, pattern) Count the number of matches in a string. str_count(fruit, "a")

Subset Strings

	str_sub(string, start = 1L, end = -1L) Extract substrings from a character vector. str_sub(fruit, 1, 3); str_sub(fruit, -2)
	str_subset(string, pattern, negate = FALSE) Return only the strings that contain a pattern match. str_subset(fruit, "p")
	str_extract(string, pattern) Return the first pattern match found in each string, as a vector. Also str_extract_all() to return every pattern match. str_extract(fruit, "[aeiou]")
	str_match(string, pattern) Return the first pattern match found in each string, as a matrix with a column for each () group in pattern. Also str_match_all() . str_match(sentences, "(a the) ([^ +])")

Manage Lengths

	str_length(string) The width of strings (i.e. number of code points, which generally equals the number of characters). str_length(fruit)
	str_pad(string, width, side = c("left", "right", "both"), pad = " ") Pad strings to constant width. str_pad(fruit, 17)
	str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...") Truncate the width of strings, replacing content with ellipsis. str_trunc(sentences, 6)
	str_trim(string, side = c("both", "left", "right")) Trim whitespace from the start and/or end of a string. str_trim(str_pad(fruit, 17))
	str_squish(string) Trim whitespace from each end and collapse multiple spaces into single spaces. str_squish(str_pad(fruit, 17, "both"))

Mutate Strings

	str_sub() <- value. Replace substrings by identifying the substrings with str_sub() and assigning into the results. str_sub(fruit, 1, 3) <- "str"
	str_replace(string, pattern, replacement) Replace the first matched pattern in each string. Also str_remove() . str_replace(fruit, "p", "-")
	str_replace_all(string, pattern, replacement) Replace all matched patterns in each string. Also str_remove_all() . str_replace_all(fruit, "p", "-")
	str_to_lower(string, locale = "en")¹ Convert strings to lower case. str_to_lower(sentences)
	str_to_upper(string, locale = "en")¹ Convert strings to upper case. str_to_upper(sentences)
	str_to_title(string, locale = "en")¹ Convert strings to title case. Also str_to_sentence() . str_to_title(sentences)

Join and Split

	str_c(..., sep = "", collapse = NULL) Join multiple strings into a single string. str_c(letters, LETTERS)
	str_flatten(string, collapse = "") Combines into a single string, separated by collapse. str_flatten(fruit, ",")
	str_dup(string, times) Repeat strings times times. Also str_unique() to remove duplicates. str_dup(fruit, times = 2)
	str_split_fixed(string, pattern, n) Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also str_split() to return a list of substrings and str_split_n() to return the nth substring. str_split_fixed(sentences, " ", n=3)
	str_glue(..., .sep = "", .envir = parent.frame()) Create a string from strings and {expressions} to evaluate. str_glue("Pi is {pi}")
	str_glue_data(.x, ..., .sep = "", .envir = parent.frame(), .na = "NA") Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate. str_glue_data(mtcars, "{rownames(mtcars)} has {hp} hp")

Order Strings

	str_order(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)¹ Return the vector of indexes that sorts a character vector. fruit[str_order(fruit)]
	str_sort(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)¹ Sort a character vector. str_sort(fruit)

Helpers

	str_conv(string, encoding) Override the encoding of a string. str_conv(fruit, "ISO-8859-1")
	str_view(string, pattern, match = NA) View HTML rendering of all regex matches. str_view(sentences, "[aeiou])")
	str_equal(x, y, locale = "en", ignore_case = FALSE, ...)¹ Determine if two strings are equivalent. str_equal(c("a", "b"), c("a", "c"))
	str_wrap(string, width = 80, indent = 0, exdent = 0) Wrap strings into nicely formatted paragraphs. str_wrap(sentences, 20)

¹ See bit.ly/ISO639-1 for a complete list of locales.

Need to Know

Pattern arguments in string are interpreted as regular expressions *after any special characters have been parsed*.

In R, you write regular expressions as *strings*, sequences of characters surrounded by quotes ("") or single quotes('').

Some characters cannot be represented directly in an R string. These must be represented as **special characters**, sequences of characters that have a specific meaning., e.g.

Special Character	Represents
\\\	\
'"	"
\n	new line

Run `?""` to see a complete list

Because of this, whenever a \ appears in a regular expression, you must write it as \\ in the string that represents the regular expression.

Use `writeLines()` to see how R views your string after all special characters have been parsed.

```
writeLines("\\\")
#\\"
```

```
writeLines("\\\" is a backslash")
#\\" is a backslash
```

INTERPRETATION

Patterns in stringr are interpreted as regexs. To change this default, wrap the pattern in one of:

`regex(pattern, ignore_case = FALSE, multiline = FALSE, comments = FALSE, dotall = FALSE, ...)`
Modifies a regex to ignore cases, match end of lines as well of end of strings, allow R comments within regex's , and/or to have . match everything including \n.
`str_detect("i", regex("i", TRUE))`

`fixed()` Matches raw bytes but will miss some characters that can be represented in multiple ways (fast). `str_detect("\u0130", fixed("i"))`

`coll()` Matches raw bytes and will use locale specific collation rules to recognize characters that can be represented in multiple ways (slow). `str_detect("\u0130", coll("i", TRUE, locale = "tr"))`

`boundary()` Matches boundaries between characters, line_breaks, sentences, or words. `str_split(sentences, boundary("word"))`



Regular Expressions -

Regular expressions, or *regexp*s, are a concise language for describing patterns in strings.

MATCH CHARACTERS

string (type this)	regexp (to mean this)	matches (which matches this)	example
a (etc.)	a (etc.)	a (etc.)	see("a")
\.	\.	.	see("\.")
\!	\!	!	see("\!")
\?	\?	?	see("\?")
\\\	\\\	\	see("\\\\")
\(\((see("\()")
\)	\))	see("\)")
\{	\{	{	see("\{")
\}	\}	}	see("\}")
\n	\n	new line (return)	see("\n")
\t	\t	tab	see("\t")
\s	\s	any whitespace (\\$ for non-whitespaces)	see("\s")
\d	\d	any digit (\D for non-digits)	see("\d")
\w	\w	any word character (\W for non-word chars)	see("\w")
\b	\b	word boundaries	see("\b")
[:digit:] ¹	[:digit:] ¹	digits	see("[:digit:]")
[:alpha:] ¹	[:alpha:] ¹	letters	see("[:alpha:]")
[:lower:] ¹	[:lower:] ¹	lowercase letters	see("[:lower:]")
[:upper:] ¹	[:upper:] ¹	uppercase letters	see("[:upper:]")
[:alnum:] ¹	[:alnum:] ¹	letters and numbers	see("[:alnum:]")
[:punct:] ¹	[:punct:] ¹	punctuation	see("[:punct:]")
[:graph:] ¹	[:graph:] ¹	letters, numbers, and punctuation	see("[:graph:]")
[:space:] ¹	[:space:] ¹	space characters (i.e. \s)	see("[:space:]")
[:blank:] ¹	[:blank:] ¹	space and tab (but not new line)	see("[:blank:]")
.	.	every character except a new line	see(".")

¹ Many base R functions require classes to be wrapped in a second set of [], e.g. [[:digit:]]

ALTERNATES

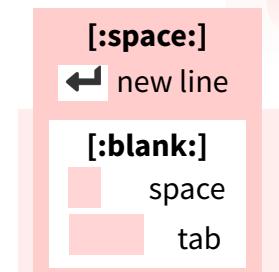
regexp	matches	example
ab d	or	alt("ab d")
[abe]	one of	alt("[abe]")
[^abe]	anything but	alt("[^abe]")
[a-c]	range	alt("[a-c]")

ANCHORS

regexp	matches	example
^a	start of string	anchor("^a")
a\$	end of string	anchor("a\$")

LOOK AROUNDS

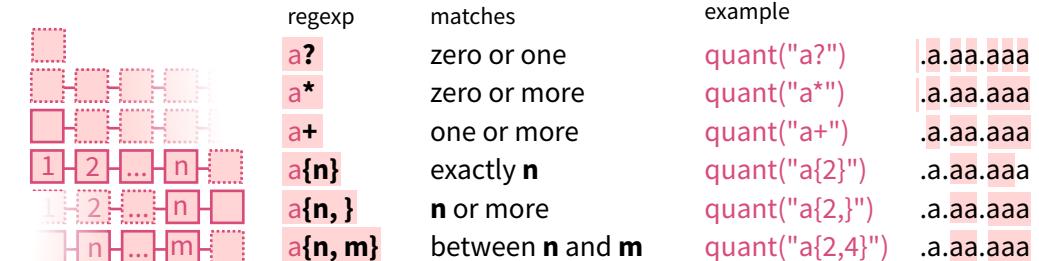
regexp	matches	example
a(?=c)	followed by	look("a(?=c)")
a(?!c)	not followed by	look("a(?!c)")
(?<=b)a	preceded by	look("(?<=b)a")
(?<!b)a	not preceded by	look("(?<!b)a")



[:alnum:]	
[:digit:]	
0	1
1	2
2	3
3	4
4	5
5	6
6	7
7	8
8	9

[:alpha:]	
[:lower:]	[:upper:]
a	A
b	B
c	C
d	D
e	E
f	F
g	G
h	H
i	I
j	J
k	K
l	L
m	M
n	N
o	O
p	P
q	Q
r	R
s	S
t	T
u	U
v	V
w	W
x	X
y	Y
z	Z

QUANTIFIERS



`quant <- function(rx) str_view(".a.aa.aaa", rx)`

quant	regexp	matches	example
zero or one	a?	aaa	quant("a?")
zero or more	a*	aaa	quant("a*")
one or more	a+	aaa	quant("a+")
exactly n	a{n}	aaa	quant("a{2}")
n or more	a{n,}	aaa	quant("a{2,}")
between n and m	a{n, m}	aaa	quant("a{2,4}")

GROUPS

Use parentheses to set precedent (order of evaluation) and create groups

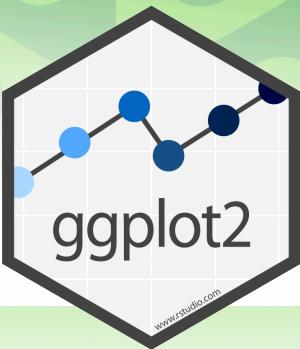
regexp	matches	example
(ab d)e	sets precedence	alt("(ab d)e")

Use an escaped number to refer to and duplicate parentheses groups that occur earlier in a pattern. Refer to each group by its order of appearance

string	regexp	matches	example
(type this)	(to mean this)	(which matches this)	(the result is the same as ref("abba"))

`ref <- function(rx) str_view("abbaab", rx)`

Data visualization with ggplot2 :: CHEATSHEET



Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data** set, a **coordinate system**, and **geoms**—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot (data = <DATA>) +
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),
  stat = <STAT>, position = <POSITION>) +
  <COORDINATE_FUNCTION> +
  <FACET_FUNCTION> +
  <SCALE_FUNCTION> +
  <THEME_FUNCTION>
```

required: **GEOM_FUNCTION**, **mapping**, **stat**, **position**
Not required, sensible defaults supplied: **COORDINATE_FUNCTION**, **FACET_FUNCTION**, **SCALE_FUNCTION**, **THEME_FUNCTION**

`ggplot(data = mpg, aes(x = cty, y = hwy))` Begins a plot that you finish by adding layers to. Add one geom function per layer.

`last_plot()` Returns the last plot.

`ggsave("plot.png", width = 5, height = 5)` Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

Aes Common aesthetic values.

color and **fill** - string ("red", "#RRGGBB")

linetype - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

size - integer (line width in mm)

shape - integer/shape name or a single character ("a")

0	1	2	3	4	5	6	7	8	9	10	11	12
□	○	△	+	×	◇	▽	■	*	⊕	⊖	⊗	⊗⊗
13	14	15	16	17	18	19	20	21	22	23	24	25
☒	▢	○	△	◇	○	○	●	●	●	●	●	△△



Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables.
Each function returns a layer.

GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
```

- a + geom_blank()** and **a + expand_limits()**
Ensure limits include values across all plots.
- b + geom_curve(aes(yend = lat + 1, xend = long + 1, curvature = 1))** - x, yend, alpha, angle, color, curvature, linetype, size
- a + geom_path(lineend = "butt", linejoin = "round", linemitre = 1)** - x, y, alpha, color, group, linetype, size
- a + geom_polygon(aes(alpha = 50))** - x, y, alpha, color, fill, group, subgroup, linetype, size
- b + geom_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1))** - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size
- a + geom_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

- b + geom_abline(aes(intercept = 0, slope = 1))**
- b + geom_hline(aes(yintercept = lat))**
- b + geom_vline(aes(xintercept = long))**
- b + geom_segment(aes(yend = lat + 1, xend = long + 1))**
- b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

- ```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```
- c + geom\_area(stat = "bin")**  
x, y, alpha, color, fill, linetype, size
  - c + geom\_density(kernel = "gaussian")**  
x, y, alpha, color, fill, group, linetype, size, weight
  - c + geom\_dotplot()**  
x, y, alpha, color, fill
  - c + geom\_freqpoly()**  
x, y, alpha, color, group, linetype, size
  - c + geom\_histogram(binwidth = 5)**  
x, y, alpha, color, fill, linetype, size, weight
  - c2 + geom\_qq(aes(sample = hwy))**  
x, y, alpha, color, fill, linetype, size, weight

### discrete

- ```
d <- ggplot(mpg, aes(fct))
```
- d + geom_bar()**
x, alpha, color, fill, linetype, size, weight

TWO VARIABLES both continuous

```
e <- ggplot(mpg, aes(cty, hwy))
```

- e + geom_label(aes(label = cty), nudge_x = 1, nudge_y = 1)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust
- e + geom_point()**
x, y, alpha, color, fill, shape, size, stroke
- e + geom_quantile()**
x, y, alpha, color, group, linetype, size, weight
- e + geom_rug(sides = "bl")**
x, y, alpha, color, linetype, size
- e + geom_smooth(method = lm)**
x, y, alpha, color, fill, group, linetype, size, weight
- e + geom_text(aes(label = cty), nudge_x = 1, nudge_y = 1)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

one discrete, one continuous

```
f <- ggplot(mpg, aes(class, hwy))
```

- f + geom_col()**
x, y, alpha, color, fill, group, linetype, size
- f + geom_boxplot()**
x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight
- f + geom_dotplot(binaxis = "y", stackdir = "center")**
x, y, alpha, color, fill, group
- f + geom_violin(scale = "area")**
x, y, alpha, color, fill, group, linetype, size, weight

both discrete

```
g <- ggplot(diamonds, aes(cut, color))
```

- g + geom_count()**
x, y, alpha, color, fill, shape, size, stroke
- e + geom_jitter(height = 2, width = 2)**
x, y, alpha, color, fill, shape, size

THREE VARIABLES

```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)); l <- ggplot(seals, aes(long, lat))
```

- l + geom_contour(aes(z = z))**
x, y, z, alpha, color, group, linetype, size, weight
- l + geom_contour_filled(aes(fill = z))**
x, y, alpha, color, fill, group, linetype, size, subgroup

continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))
```

- h + geom_bin2d(binwidth = c(0.25, 500))**
x, y, alpha, color, fill, linetype, size, weight
- h + geom_density_2d()**
x, y, alpha, color, group, linetype, size
- h + geom_hex()**
x, y, alpha, color, fill, size

continuous function

```
i <- ggplot(economics, aes(date, unemploy))
```

- i + geom_area()**
x, y, alpha, color, fill, linetype, size
- i + geom_line()**
x, y, alpha, color, group, linetype, size
- i + geom_step(direction = "hv")**
x, y, alpha, color, group, linetype, size

visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

- j + geom_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size
- j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width
Also **geom_errorbarh()**.
- j + geom_linerange()**
x, ymin, ymax, alpha, color, group, linetype, size
- j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

maps

```
data <- data.frame(murder = USArrests$Murder,
  state = tolower(rownames(USArrests)))
```

```
map <- map_data("state")
```

```
k <- ggplot(data, aes(fill = murder))
```

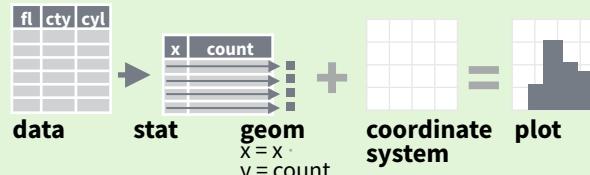
- k + geom_map(aes(map_id = state), map = map)**
+ **expand_limits(x = map\$long, y = map\$lat)**
map_id, alpha, color, fill, linetype, size

- l + geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)**
x, y, alpha, fill
- l + geom_tile(aes(fill = z))**
x, y, alpha, color, fill, linetype, size, width

Stats

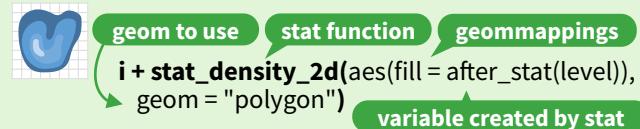
An alternative way to build a layer.

A stat builds new variables to plot (e.g., count, prop).



Visualize a stat by changing the default stat of a geom function, `geom_bar(stat="count")` or by using a stat function, `stat_count(geom="bar")`, which calls a default geom to make a layer (equivalent to a geom function).

Use `after_stat(name)` syntax to map the stat variable `name` to an aesthetic.



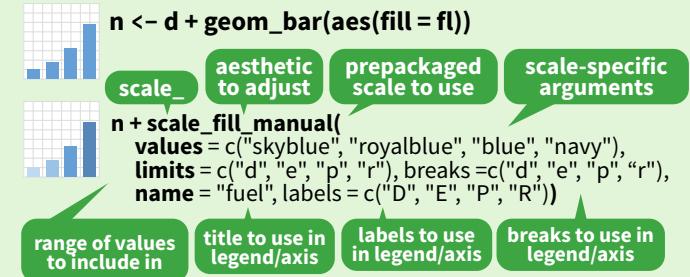
```

c + stat_bin(binwidth = 1, boundary = 10)
x, y | count, ncount, density, ndensity
c + stat_count(width = 1) x, y | count, prop
c + stat_density(adjust = 1, kernel = "gaussian")
x, y | count, density, scaled
e + stat_bin_2d(bins = 30, drop = T)
x, y, fill | count, density
e + stat_bin_hex(bins = 30) x, y, fill | count, density
e + stat_density_2d(contour = TRUE, n = 100)
x, y, color, size | level
e + stat_ellipse(level = 0.95, segments = 51, type = "t")
l + stat_contour(aes(z = z)) x, y, z, order | level
l + stat_summary_hex(aes(z = z), bins = 30, fun = max)
x, y, z, fill | value
l + stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | value
f + stat_boxplot(coef = 1.5)
x, y | lower, middle, upper, width, ymin, ymax
f + stat_ydensity(kernel = "gaussian", scale = "area") x, y |
density, scaled, count, n, violinwidth, width
e + stat_ecdf(n = 40) x, y | x, y
e + stat_quantile(quantiles = c(0.1, 0.9),
formula = y ~ log(x), method = "rq") x, y | quantile
e + stat_smooth(method = "lm", formula = y ~ x, se = T,
level = 0.95) x, y | se, x, y, ymin, ymax
ggplot() + xlim(-5, 5) + stat_function(fun = dnorm,
n = 20, geom = "point") x | x, y
ggplot() + stat_qq(aes(sample = 1:100))
x, y, sample | sample, theoretical
e + stat_sum() x, y, size | n, prop
e + stat_summary(fun.data = "mean_cl_boot")
h + stat_summary_bin(fun = "mean", geom = "bar")
e + stat_identity()
e + stat_unique()
  
```

Scales

Override defaults with `scales` package.

Scales map data values to the visual values of an aesthetic. To change a mapping, add a new scale.



GENERAL PURPOSE SCALES

Use with most aesthetics

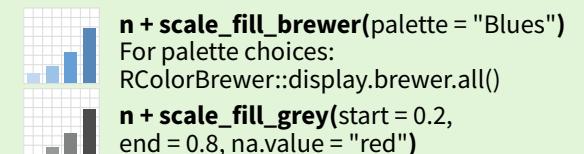
`scale_*_continuous()` - Map cont' values to visual ones.
`scale_*_discrete()` - Map discrete values to visual ones.
`scale_*_binned()` - Map continuous values to discrete bins.
`scale_*_identity()` - Use data values as visual ones.
`scale_*_manual(values = c())` - Map discrete values to manually chosen visual ones.
`scale_*_date(date_labels = "%m/%d")`,
`date_breaks = "2 weeks"` - Treat data values as dates.
`scale_*_datetime()` - Treat data values as date times.
 Same as `scale_*_date()`. See `?strptime` for label formats.

X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

`scale_x_log10()` - Plot x on log10 scale.
`scale_x_reverse()` - Reverse the direction of the x axis.
`scale_x_sqrt()` - Plot x on square root scale.

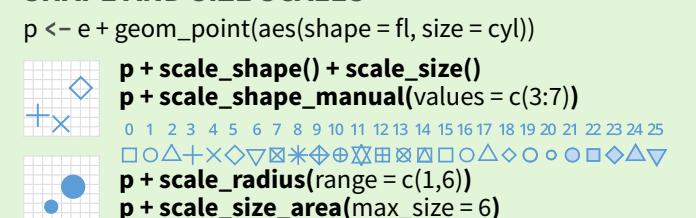
COLOR AND FILL SCALES (DISCRETE)



COLOR AND FILL SCALES (CONTINUOUS)



SHAPE AND SIZE SCALES



Coordinate Systems

`r <- d + geom_bar()`

`r + coord_cartesian(xlim = c(0, 5))` - xlim, ylim
The default cartesian coordinate system.

`r + coord_fixed(ratio = 1/2)`
ratio, xlim, ylim - Cartesian coordinates with fixed aspect ratio between x and y units.

`r + coord_flip()`
Flip cartesian coordinates by switching x and y aesthetic mappings.

`r + coord_polar(theta = "x", direction=1)`
theta, start, direction - Polar coordinates.

`r + coord_trans(y = "sqrt")` - x, y, xlim, ylim
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

`pi + coord_quickmap()`
`pi + coord_map(projection = "ortho", orientation = c(41, -74, 0))` - projection, xlim, ylim
Map projections from the mapproj package (mercator (default), azequalarea, lagrange, etc.).

Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

`s <- ggplot(mpg, aes(fl, fill = drv))`

`s + geom_bar(position = "dodge")`
Arrange elements side by side.

`s + geom_bar(position = "fill")`
Stack elements on top of one another, normalize height.

`e + geom_point(position = "jitter")`
Add random noise to X and Y position of each element to avoid overplotting.

`e + geom_label(position = "nudge")`
Nudge labels away from points.

`s + geom_bar(position = "stack")`
Stack elements on top of one another.

Each position adjustment can be recast as a function with manual `width` and `height` arguments:
`s + geom_bar(position = position_dodge(width = 1))`

Themes

`r + theme_bw()`
White background with grid lines.

`r + theme_gray()`
Grey background (default theme).

`r + theme_dark()`
Dark for contrast.

`r + theme_classic()`
`r + theme_light()`

`r + theme_linedraw()`
`r + theme_minimal()`

`r + theme_void()`
Empty theme.

`r + theme()` Customize aspects of the theme such as axis, legend, panel, and facet properties.

`r + ggtitle("Title") + theme(plot.title.position = "plot")`
`r + theme(panel.background = element_rect(fill = "blue"))`

Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t <- ggplot(mpg, aes(cty, hwy)) + geom_point()`

`t + facet_grid(. ~ fl)`
Facet into columns based on fl.

`t + facet_grid(year ~ .)`
Facet into rows based on year.

`t + facet_grid(year ~ fl)`
Facet into both rows and columns.

`t + facet_wrap(~ fl)`
Wrap facets into a rectangular layout.

Set `scales` to let axis limits vary across facets.

`t + facet_grid(drv ~ fl, scales = "free")`
x and y axis limits adjust to individual facets:
`"free_x"` - x axis limits adjust
`"free_y"` - y axis limits adjust

Set `labeler` to adjust facet label:

`t + facet_grid(. ~ fl, labeler = label_both)`

`fl: c` `fl: d` `fl: e` `fl: p` `fl: r`

`t + facet_grid(fl ~ ., labeler = label_bquote(alpha ^ .(fl)))`

`alpha^c` `alpha^d` `alpha^e` `alpha^p` `alpha^r`

Labels and Legends

Use `labs()` to label the elements of your plot.

`t + labs(x = "New x axis label", y = "New y axis label", title = "Add a title above the plot", subtitle = "Add a subtitle below title", caption = "Add a caption below plot", alt = "Add alt text to the plot", <AES> = "New <AES> legend title")`

`t + annotate(geom = "text", x = 8, y = 9, label = "A")`
Places a geom with manually selected aesthetics.

`p + guides(x = guide_axis(n.dodge = 2))` Avoid crowded or overlapping labels with `guide_axis(n.dodge` or `angle`).

`n + guides(fill = "none")` Set legend type for each aesthetic: colorbar, legend, or none (no legend).

`n + theme(legend.position = "bottom")`
Place legend at "bottom", "top", "left", or "right".

`n + scale_fill_discrete(name = "Title", labels = c("A", "B", "C", "D", "E"))`
Set legend title and labels with a scale function.

Zooming

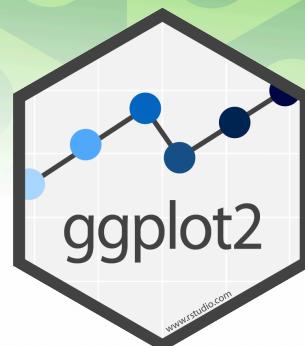
Without clipping (preferred):

`t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))`

With clipping (removes unseen data points):

`t + xlim(0, 100) + ylim(10, 20)`

`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`





Labelled data with labelled :: CHEAT SHEET

The **labelled** package provides a set of functions and methods to handle and to manipulate labelled data, as imported with **haven** package.

Basics

Labelled data is a common data structure in other statistical environment such as Stata, SAS or SPSS.

It consists of a set of additional attributes for numeric and character vectors (including columns of a data frame).

There are 3 types of attributes:

1. **Variable labels** (a short description of a variable)
2. **Value labels** (labels associated to specific values)
3. **Missing values**:
 - User-defined missing values (SPSS style)
 - Tagged NA (Stata and SAS style)

Variable labels

MANIPULATING A VECTOR

- var_label(x) or var_label(df\$v1)**
Get the variable label associated to a vector x
- var_label(x) <- "variable description"**
Add/modify a variable label to x
- var_label(x) <- NULL**
Remove the variable label associated to x

MANIPULATING A DATA.FRAME

- var_label(df)**
List all variable labels associated with columns of df
- var_label(df) <- list(v1 = "variable 1", v2 = "variable 2")**
Update variable labels of some columns of df
- df %>% set_variable_labels(v1 = "variable 1", v2 = "variable 2", v3 = NULL)**
Update variable labels using dplyr syntax
- df %>% look_for()**
Return a data frame with all variable names and labels
- df %>% look_for("s")**
Search variables containing "s" in their name or label
- df %>% look_for(details = TRUE)**
Return additional details on each variable

Value labels

When value labels are attached to a numeric or character vector, the vector's class becomes **haven_labelled**. A major difference with a factor is that values of the vector are not changed and it is not mandatory to attach a label to each value.

MANIPULATING A VECTOR

- val_label(x, value) or val_label(df\$v1, value)**
Get the label attached to a specific value of a vector
- val_label(x, value) <- "label"**
Set/Update the label attached to a specific value
- val_label(x, value) <- NULL**
Remove the label attached to a specific value
- val_labels(x)**
Get all value labels attached to a vector
- val_labels(x) <- c(no = 0, yes = 1, maybe = 9)**
Set/Update all value labels attached to a vector
- val_labels(x) <- NULL**
Remove all value labels attached to a vector
- labelled(c("F", "F", "M"), c(Female = "F", Male = "M"))**
Create a labelled vector
- sort_val_labels(x, according_to = "values")**
Sort value labels according to values (or labels)
- drop_unused_value_labels(x)**
Remove value labels not observed in the data

MANIPULATING A DATA.FRAME

- df %>% set_value_labels(v1 = c(Yes = 1, No = 2), v2 = c(Male = "M", Female = "F"))**
Define value labels of several variables
- df %>% add_value_labels(v1 = c(Unknown = 9))**
Add specific value labels to a variable (other already defined value labels remains unchanged)
- df %>% remove_value_labels(v1 = 9)**
Remove specific value labels to a variable
- df %>% set_value_labels(v1 = NULL)**
Remove all value labels attached to a variable
- df %>% drop_unused_value_labels()**
Remove value labels not observed in the data

Missing values

USER-DEFINED MISSING VALUES (SPSS STYLE)

Used to indicate that some values should be considered as missing. However, they will not be treated as NA as long as they are not converted to proper NA.

When missing values are attached to a numeric or character vector, the vector's class becomes **haven_labelled_spss**.

When importing a SPSS file, use the option *user_na = TRUE* to keep defined missing values (otherwise, they will be converted to NA).

na_values(x)

Get individual missing values attached to a vector

na_values(x) <- c(8, 9, 10)

df %>% set_na_values(v1 = c(8, 9, 10))
Set/Update individual missing values (NULL to remove)

na_range(x)

Get a range of missing values attached to a vector

na_range(x) <- c(8, 10)

df %>% set_na_range(v1 = c(8, 10))
Set/Update a range of missing values (NULL to remove)

user_na_to_na(x) or df %>% user_na_to_na()

Convert user-defined missing values to NA

is_na(x)

TRUE if NA or if a user-defined missing value

TAGGED NAs (STATA & SAS STYLE)

"Tagged" missing values work exactly like regular R missing values except that they store one additional byte of information: a tag, which is usually a letter ("a" to "z").

x <- c(1:5, tagged_na("a"), tagged_na("z"), NA)

tagged_na("a") generates a NA with a tag

is.na(x)

Tagged NAs work identically to regular NAs

is_tagged_na(x)

Test if it is a tagged NA

na_tag(x)

Display the tags associated to tagged NAs

format_tagged_na(x)

Convert x to a character vector showing the tagged NAs



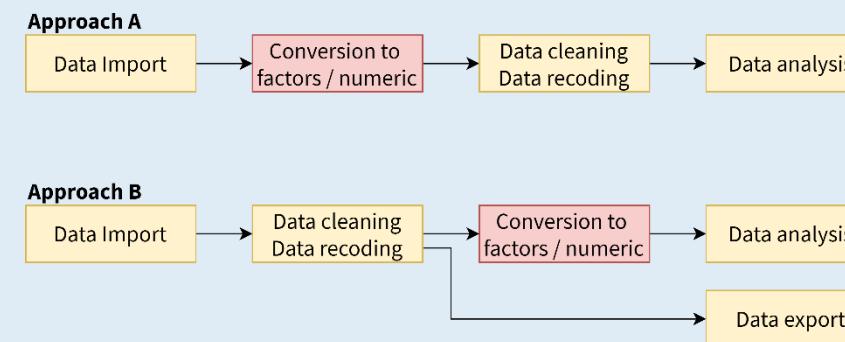
When using labelled data?

`haven_labelled` and `haven_labelled_spss` classes introduced in **haven** package allow to add metadata (variable labels, value labels and SPSS-style missing values) to vectors / data frame columns and to properly import these metadata from SAS, Stata or SPSS.

Functions and methods provided by **labelled** package are designed for easy manipulation of such labelled data.

It should be noted that **value labels** doesn't imply that your vectors should be considered as categorical or continuous. Therefore, value labels are not intended to be used for data analysis. For example, before performing modeling or plotting, you should convert vectors with value labels into factors or into classic numeric/character vectors.

Two main approaches could be considered:



In **approach A**, labelled vectors are converted into factors or into numeric/character vectors just after data import, using **unlabelled()**, **to_factor()** or **unclass()**. Then, data cleaning, data recoding and analysis are performed using classic R vector types.

In **approach B**, labelled vectors are kept for data cleaning and recoding, allowing to preserve original coding, in particular if data should be reexported after that step. Functions provided by **labelled** will be useful for managing value labels.

However, as in approach A, labelled vectors will have to be converted into classic factors or numeric vectors before data analysis as this is the way categorical and continuous variables should be coded for analysis.

Conversion

OF LABELLED VECTORS

If all value labels and user-defined missing values are removed from a labelled vector, the `haven_labelled` class will be removed and the vector will be transformed into a basic numeric or character vector. **Values of the vector will remain unchanged.**

`remove_val_labels(x)`

Remove value labels attached to a vector.

`remove_user_na(x)`

Remove user-defined (`na_values` and `na_range`) from a vector

`unclass(x)`

Remove the `haven_labelled` class. Therefore, the vector will be considered as a classical numeric or character vector. Value labels and user-defined missing values will still be visible as attributes attached to the vector.

When converting a labelled vector into a factor or a character vector of the value labels, be aware that **original values of the vector will be converted.**

`to_character(x)`

Convert into a character vector replacing values by their corresponding value label

`to_factor(x)`

Convert into a character vector replacing values by their corresponding value label

`to_factor(x, levels = "prefixed")`

Value labels will be prefixed with their original value

`to_factor(x, strict = TRUE)`

Convert into a factor only if all observed values have a value label

`df %>% to_factor()`

Convert all labelled vectors and only labelled vectors into factors

`df %>% to_factor(labelled_only = FALSE)`

Convert all columns (including non labelled vectors) into factors

`unlabelled(x)`

`df %>% unlabelled()`

Labelled vectors will be converted into factors only if all observed values have a value label. Otherwise, they will be unclassed. Similar to `df %>% to_factor(labelled_only = T, strict = T, unclass = T)`

`to_factor(x, drop_unused_labels = TRUE)`

`df %>% unlabelled(drop_unused_labels = TRUE)`

Unused value labels will be dropped before conversion into factors

INTO LABELLED VECTORS

`to_labelled(f)`

Convert a factor into a numeric labelled vector. Note that `to_labelled(to_factor(x))` and `x` will not be identical (original coding will be lost).

`to_labelled(df)`

If `df` was imported with the **foreign** package or if it is a data set created with **memisc** package, meta data (variable labels, value labels and user-defined missing values) will be converted into **labelled** format.

If any value label or user-defined missing value is added to a numeric or a character vector, it will be automatically converted into a labelled vector.

Values of the vector will remain unchanged.

Miscellaneous

`nolabel_to_na(x)` or `df %>% nolabel_to_na()`

For labelled vectors, values without a value label will be converted into NA

`val_labels_to_na(x)` or `df %>% val_labels_to_na()`

For labelled vectors, values with a value label will be converted into NA

`df2 %>% copy_labels_from(df1)`

Copy variable labels, values labels and user-defined missing values from `df1` to `df2` based on shared columns names. Useful when attributes are lost after some data manipulation.

`recode(x, `2` = 1, `3` = 2)`

Apply `dplyr::recode()` to a labelled vector. Attached value labels will remain unchanged.

`recode(x, `2` = 1, .combine_value_labels = TRUE)`

This option will combine value labels of original values merged together to produce new value labels. It is recommended to check that the result is appropriate.

`update_labelled(x)` or `df %>% update_labelled()`

If `x/df` was imported/created using an older version of **haven** or **labelled**, you may encounter some unexpected results. `update_labelled()` will update all labelled vectors to be consistent with the current implementation.

Data Science in Spark with sparklyr :: CHEAT SHEET



Intro

sparklyr is an R interface for Apache Spark™. It enables us to write all of our analysis code in R, but have the actual processing happen inside Spark clusters. Easily manipulate and model large-scale using R and Spark via **sparklyr**.

Import



READ A FILE INTO SPARK

Arguments that apply to all functions:
sc, name, path, options=list(), repartition=0,
memory=TRUE, overwrite=TRUE

CSV `spark_read_csv(header=TRUE, columns=NULL, infer_schema=TRUE, delimiter = "", quote = "\\"", escape = "\\\\", charset = "UTF-8", null_value = NULL)`

JSON `spark_read_json()`

PARQUET `spark_read_parquet()`

TEXT `spark_read_text()`

ORC `spark_read_orc()`

LIBSVM `spark_read_libsvm()`

DELTA `spark_read_delta()`

AVRO `spark_read_avro()`

R DATA FRAME INTO SPARK

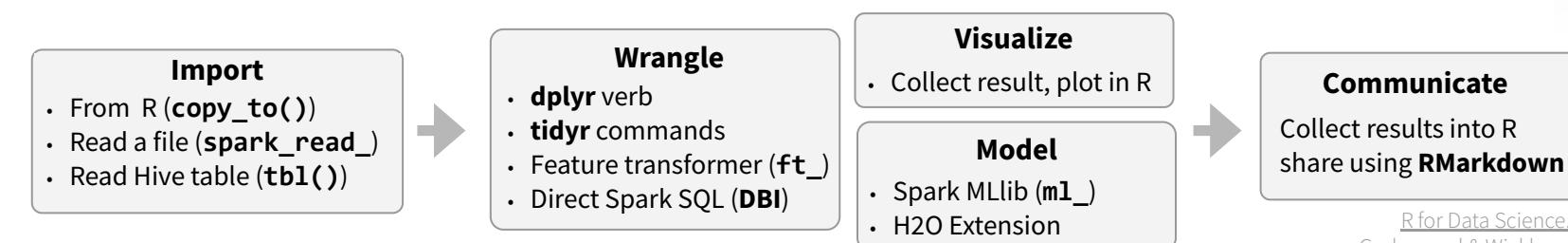
`dplyr::copy_to(dest, df, name)`

Apache Arrow accelerates data transfer between R and Spark. To use, simply load the library

`library(sparklyr)`
`library(arrow)`

FROM A TABLE IN HIVE

`dplyr::tbl(scr, ...)` - Creates a reference to the table without loading it into memory



*R for Data Science,
Golemund & Wickham*

Wrangle

DPLYR VERBS

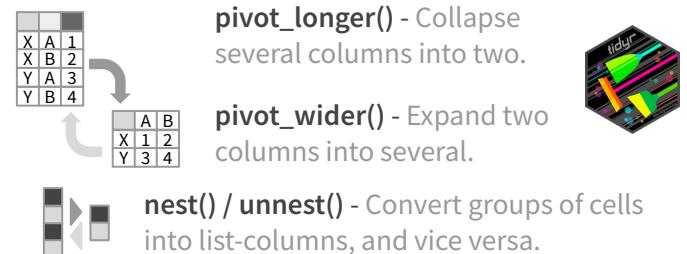
Translates into Spark SQL statements

```

dplyr
copy_to(sc, mtcars) %>%
  mutate(trm = ifelse(am == 0,
                      "auto", "man")) %>%
  group_by(trm) %>%
  summarise_all(mean)
  
```

TIDYR

`pivot_longer()` - Collapse several columns into two.



`pivot_wider()` - Expand two columns into several.



`nest()` / `unnest()` - Convert groups of cells into list-columns, and vice versa.

`unite()` / `separate()` - Split a single column into several columns, and vice versa.

`fill()` - Fill NA with the previous value

FEATURE TRANSFORMERS

`ft_binarizer()` - Assigns values based on threshold

`ft_bucketizer()` - Numeric column to discretized column

`ft_count_vectorizer()` - Extracts a vocabulary from document

`ft_discrete_cosine_transform()` - 1D discrete cosine transform of a real vector

`ft_elementwise_product()` - Element-wise product between 2 cols

`ft_hashing_tf()` - Maps a sequence of terms to their term frequencies using the hashing trick.

`ft_idf()` - Compute the Inverse Document Frequency (IDF) given a collection of documents.

`ft_imputer()` - Imputation estimator for completing missing values, uses the mean or the median of the columns.

`ft_index_to_string()` - Index labels back to label as strings

`ft_interaction()` - Takes in Double and Vector columns and outputs a flattened vector of their feature interactions.

`ft_max_abs_scaler()` - Rescale each feature individually to range [-1, 1]

`ft_min_max_scaler()` - Rescale each feature to a common range [min, max] linearly

`ft_ngram()` - Converts the input array of strings into an array of n-grams

`ft_bucketed_random_projection_lsh()`
`ft_minhash_lsh()` - Locality Sensitive Hashing functions for Euclidean distance and Jaccard distance (MinHash)

`ft_normalizer()` - Normalize a vector to have unit norm using the given p-norm

`ft_one_hot_encoder()` - Continuous to binary vectors

`ft_pca()` - Project vectors to a lower dimensional space of top k principal components.

`ft_quantile_discretizer()` - Continuous to binned categorical values.

`ft_regex_tokenizer()` - Extracts tokens either by using the provided regex pattern to split the text.

`ft_robust_scaler()` - Removes the median and scales according to standard scale.

`ft_standard_scaler()` - Removes the mean and scaling to unit variance using column summary statistics

`ft_stop_words_remover()` - Filters out stop words from input

`ft_string_indexer()` - Column of labels into a column of label indices.

`ft_tokenizer()` - Converts to lowercase and then splits it by white spaces

`ft_vectorAssembler()` - Combine vectors into single row-vector

`ft_vector_indexer()` - Indexing categorical feature columns in a dataset of Vector

`ft_vector_slicer()` - Takes a feature vector and outputs a new feature vector with a subarray of the original features

`ft_word2vec()` - Word2Vec transforms a word into a code

Visualize



DPLYR + GGPLOT2

```

copy_to(sc, mtcars) %>%
  group_by(cyl) %>%
  summarise(mpg_m = mean(mpg)) %>%
  collect() %>%
  ggplot() +
  geom_col(aes(cyl, mpg_m))
  
```

Summarize in Spark
Collect results in R
Create plot

Data Science in Spark with sparklyr :: CHEAT SHEET



Modeling

REGRESSION

`ml_linear_regression()` - Linear regression.

`ml_aft_survival_regression()` - Parametric survival regression model named accelerated failure time (AFT) model

`ml_generalized_linear_regression()` - GLM

`ml_isotonic_regression()` - Currently implemented using parallelized pool adjacent violators algorithm. Only univariate (single feature) algorithm supported

`ml_random_forest_regressor()` - Regression using random forests.

CLASSIFICATION

`ml_linear_svc()` - Classification using linear support vector machines

`ml_logistic_regression()` - Logistic regression

`ml_multilayer_perceptron_classifier()` - Classification model based on the Multilayer Perceptron.

`ml_naive_bayes()` - It supports Multinomial NB which can handle finitely supported discrete data

`ml_one_vs_rest()` - Reduction of Multiclass Classification to Binary Classification. Performs reduction using one against all strategy.

TREE

`ml_decision_tree_classifier()` | `ml_decision_tree()` | `ml_decision_tree_regressor()` - Classification and regression using decision trees

`ml_gbt_classifier()` | `ml_gradient_boosted_trees()` | `ml_gbt_regressor()` - Binary classification and regression using gradient boosted trees

`ml_random_forest_classifier()` - Classification and regression using random forests.

`ml_feature_importances()` | `ml_tree_feature_importance()` - Feature Importance for Tree Models

CLUSTERING

`ml_bisecting_kmeans()` - A bisecting k-means algorithm based on the paper

`ml_lda()` | `ml_describe_topics()` | `ml_log_likelihood()` | `ml_log_perplexity()` | `ml_topics_matrix()` - LDA topic model designed for text documents.

`ml_gaussian_mixture()` - Expectation maximization for multivariate Gaussian Mixture Models (GMMs)

`ml_kmeans()` | `ml_compute_cost()` | `ml_compute_silhouette_measure()` - Clustering with support for k-means

`ml_power_iteration()` - For clustering vertices of a graph given pairwise similarities as edge properties.

FEATURE

`ml_chisquare_test(x,features,label)` - Pearson's independence test for every feature against the label

`ml_default_stop_words()` - Loads the default stop words for the given language

STATS

`ml_summary()` - Extracts a metric from the summary object of a Spark ML model

`ml_corr()` - Compute correlation matrix

RECOMMENDATION

`ml_als()` | `ml_recommend()` - Recommendation using Alternating Least Squares matrix factorization

EVALUATION

`ml_clustering_evaluator()` - Evaluator for clustering

`ml_evaluate()` - Compute performance metrics

`ml_binary_classification_evaluator()` | `ml_binary_classification_eval()` | `ml_classification_eval()` - A set of functions to calculate performance metrics for prediction models.

FREQUENT PATTERN

`ml_fpgrowth()` | `ml_association_rules()` |

`ml_freq_itemsets()` - A parallel FP-growth algorithm to mine frequent itemsets.

`ml_freq_seq_patterns()` | `ml_prefixspan()` - PrefixSpan algorithm for mining frequent itemsets.

UTILITIES

`ml_call_constructor()` - Identifies the associated sparklyr ML constructor for the JVM

`ml_model_data()` - Extracts data associated with a Spark ML model

`ml_standardize_formula()` - Generates a formula string from user inputs, to be used in `ml_model` constructor

`ml_uid()` - Extracts the UID of an ML object.

ML Pipelines

Easily create a formal Spark Pipeline models using R. Save the Pipeline in native Scala. The saved model will have no dependencies on R.

INITIALIZE AND TRAIN

`ml_pipeline()` - Initializes a new Spark Pipeline

`ml_fit()` - Trains the model, outputs a Spark Pipeline Model.

SAVE AND RETRIEVE

`ml_save()` - Saves into a format that can be read by Scala and PySpark .

`ml_read()` - Reads Spark object into sparklyr.

SQL AND DPLYR

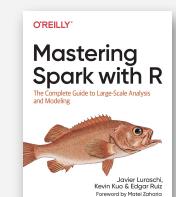
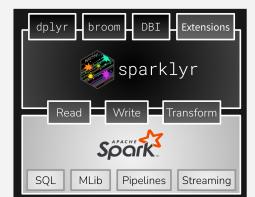
`ft_sql_transformer()` - Creates a Pipeline step based on the SQL statement passed to the command.

`ft_dplyr_transformer()` - Creates a Pipeline step based on one or several dplyr commands.



spark.rstudio.com/guides/pipelines

More Info



spark.rstudio.com

therinspark.com

Sessions

YARN CLIENT

1. Install RStudio Server on an edge node
2. Locate path to the cluster's Spark Home Directory, it normally is `"/usr/lib/spark"`
3. Basic configuration example

```
conf <- spark_config()  
conf$spark.executor.memory <- "300M"  
conf$spark.executor.cores <- 2  
conf$spark.executor.instances <- 3  
conf$spark.dynamicAllocation.enabled<-"false"
```
4. Open a connection

```
sc <- spark_connect(master = "yarn",  
                      spark_home = "/usr/lib/spark/",  
                      version = "2.1.0", config = conf)
```

YARN CLUSTER

1. Make sure to have copies of the `yarn-site.xml` and `hive-site.xml` files in the RStudio Server
2. Point environment variables to the correct paths

```
Sys.setenv(JAVA_HOME="[Path]")  
Sys.setenv(SPARK_HOME ="[Path]")  
Sys.setenv(YARN_CONF_DIR ="[Path]")
```
3. Open a connection

```
sc <- spark_connect(master = "yarn-cluster")
```

STANDALONE CLUSTER

1. Install RStudio Server on one of the existing nodes or a server in the same LAN
2. Open a connection

```
spark_connect(master="spark://host:port",  
              version = "2.0.1",  
              spark_home = [path to Spark])
```

LOCAL MODE

No cluster required. Use for learning purposes only

1. Install a local version of Spark: `spark_install()`
2. Open a connection

```
sc <- spark_connect(master="local")
```

KUBERNETES

1. Use the following to obtain the Host and Port

```
system2("kubectl", "cluster-info")
```
2. Open a connection

```
sc <- spark_connect(config =  
                     spark_config_kubernetes(  
                     "k8s://https://[HOST]:[PORT]",  
                     account = "default",  
                     image = "docker.io/owner/repo:version"))
```

CLOUD

- Databricks - `spark_connect(method = "databricks")`
- Qubole- `spark_connect(method = "qubole")`

Machine Learning Modelling in R :: CHEAT SHEET

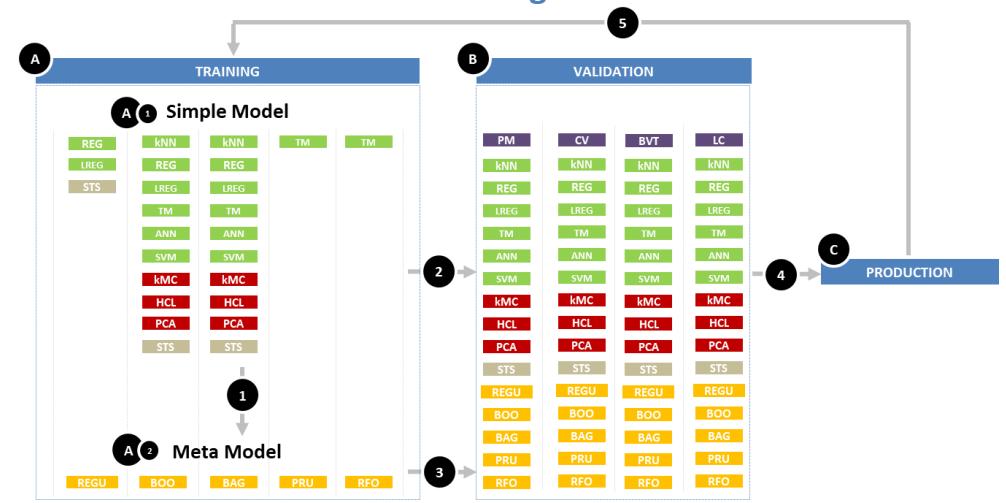
Supervised & Unsupervised Learning

ALGORITHM	DESCRIPTION	R PACKAGE::FUNCTION	SAMPLE CODE
NBC Naïve Bayes classifier	A classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naïve Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature	e1071::naiveBayes	naiveBayes(class ~ ., data = x)
kNN k-Nearest Neighbours	A non-parametric method used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether k-NN is used for classification or regression	class::knn	knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE)
LRG Linear Regression	Model the linear relationship between a scalar dependent variable Y and one or more explanatory variables (or independent variables) denoted X	stats::lm	lm(dist ~ speed, data=cars)
LRC Logistic Regression	Used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables.	stats::glm	glm(Y ~ ., family = binomial(link = 'logit'), data = X)
TM Tree-Based Models	The idea is to consecutively divide (branch) the training data into smaller and smaller features until an assignment criterion with respect to the target variable into a "data bucket" (leaf) is reached	rpart::rpart	rpart(Kyphosis ~ Age + Number + Start, data = kyphosis)
ANN Artificial Neural Network	Neural networks are built from units called perceptrons. Perceptrons have one or more inputs, an activation function and an output. An ANN model is built up by combining perceptrons in structured layers.	neuralnet::neuralnet	neuralnet(f,data=train_hidden=(5,3),linear.output=T)
SVM Support Vector Machine	A data classification method that separates data using hyperplanes	e1071::svm	svm(formula, data = NULL, ..., subset, na.action = na.omit, scale = TRUE)
PCA Principal Component Analysis	A procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.	stats::prcomp stats::princomp FactoMineR::PCA ade4::dudi.pca amap::acp	stats::prcomp(formula, data = NULL, subset, na.action, ...) stats::princomp(formula, data = NULL, subset, na.action, ...) FactoMineR::PCA(decatlon, quanti.sup = 11:12, quali.sup = 13) ade4::dudi.pca(deugStab, center = deugCent, scale = FALSE, scan = FALSE) amap::acp(lubisch)
HAC k-Mean Clustering	Aims at partitioning n observations into k clusters in which each observation belongs to the cluster with the nearest mean	stats::kmeans	kmeans(k, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"), trace = FALSE)
HCL Hierarchical Clustering	An approach which builds a hierarchy from the bottom-up, and doesn't require the number of clusters to be specified beforehand.	stats::hclust	hclust(d, method = "complete", members = NULL)

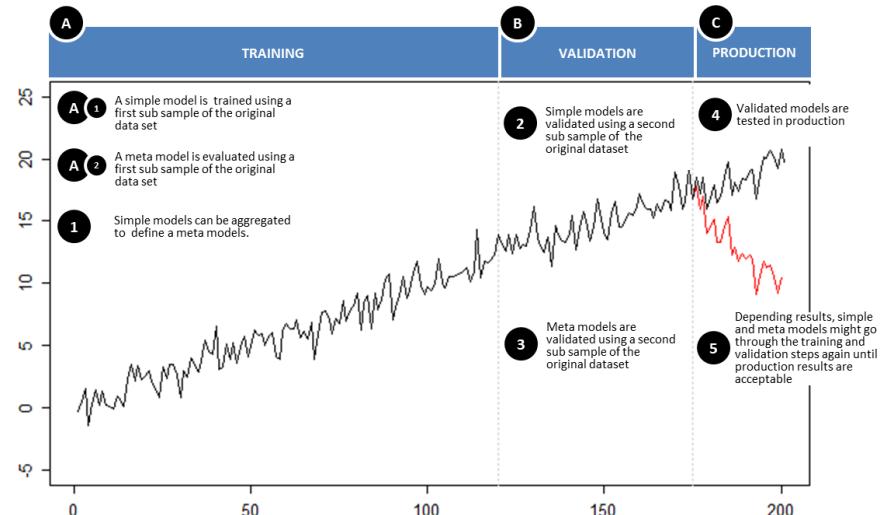
Meta-Algorithm, Time Series & Model Validation

ALGORITHM	DESCRIPTION	R PACKAGE::FUNCTION	SAMPLE CODE
REGU Regularisation L1 (Lasso) L2 (Ridge)	Regularisation adds a penalty on the different parameters of a model to reduce the freedom of the model. Hence, the model will be less likely to fit the noise of the training data and will improve the generalization abilities of the model	glmnet::glmnet	L1 : glmnet(myMatrixA, myMatrixB, family = "gaussian", alpha = 1) L2 : glmnet(myMatrixA, myMatrixB, family = "gaussian", alpha = 0)
BOO Boosting	A process of iteratively refining, e.g. by reweighting, of estimated regression and classification functions (though it has primarily been applied to the latter), in order to improve predictive ability.	gbm::gbm	gbmboost(Y ~ ., data = curr1[trnidxs,])
BAG Bagging	Bagging is a way to increase the power of a predictive statistical model by taking multiple random samples (with replacement) of the training data set, and using each of them to construct a separate model and separate predictions for the original test set	randomForest::randomForest	foreach : d <- data.frame(x=1:10, y=rnorm(10)) s <- foreach(d=i %in% d, by=row, .combine=rbind, .id=i, .d = d) ipred : bagging(formula, data, subset, na.action=na.rpart, \dots)
PRU Pruning	Pruning is a technique that reduces the size of decision tree by removing sections of the tree that provide little power to classify instances. Pruning reduces the complexity of the final classifier and hence improves predictive accuracy by reducing overfitting	rpart::rpart	prune(x, cp = 0.1)
RFO Random Forest	An ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression)	randomForest::randomForest	randomForest(X ~ ., data = Y, subset = mySub)
STS Time Series	Random sampling of observations for training and testing a model can be an issue when faced with a times dimension. Random sampling may either destroy serial correlation properties in the data which we would like to exploit	stats::xts forecast::spectral TTR	Auto-correlation: acf(x, lag.max = NULL, type = c("correlation", "covariance", "partial")) Spectral Analysis: spec.pgram(..., spans = NULL) Seasonal Decomposition of Time Series : stl(x, s.window = 7, t.window = 50, t.jump = 1)
PM Performance metrics	Depends on the problem: • Regression: squared errors, outliers, error rate... • Classification: Accuracy, precision, recall, F-score...	Regression::stats::outlierTest, stats::qqPlot ... Classification::ROCR::Tree: caret::confusionMatrix	Regression: stats::outlierTest, stats::qqPlot ... Classification: ROCR::Tree: caret::confusionMatrix
JVT Bias-Variance Tradeoff	• Simple models with few parameters are easier to compute but may lead to poorer fits (high bias). • Complex models may provide more accurate fits but may over-fit the data (high variance)	Tailored to the analysis	Tailored to the analysis
CV Cross validation	Cross validation compares the test performances of different model realisations with different sets or values of parameters	caret::createDataPartition caret::createFolds	createDataPartition(classes, p = 0.8, list = FALSE)
LC Learning Curves	Learning curves plot a model's training and test errors, or the chosen performance metric, depending on the training set size	caret::learning_curve_dat	learning_curve_dat(dat, outcome = NULL, proportion = (1:10)/10, test_prop = 0, verbose = TRUE, ...)

Standard Modelling Workflow



Time Series View



Using Git and GitHub with RStudio: : CHEATSHEET



Version control control, also known as **source control**, is the practice of tracking and managing changes to software code.

Version control systems are software tools that help software teams manage changes to source code over time.

Git is an **open-source** software for version control, originally developed in 2005 by Linus Torvalds, the creator of the Linux operating system kernel.

Git is a version control tool to track the changes in the source code of a project.

GitHub is the most popular hosting service for collaborating on code using Git.

Requirements

1. R and RStudio installed
2. Git installed
3. Register a free GitHub account



Check that Git is installed

In the Terminal of RStudio, enter `which git` to request the path to your Git executable:

```
which git
## /usr/bin/git
```

and `git --version` to see its version:

```
git --version
## git version 2.34.1
```

Introduce yourself to Git

Open a shell from RStudio *Tools > Shell* and type each line separately by substituting your name and the email associated with your GitHub account:

```
git config --global user.name 'Jane Doe'
git config --global user.email
'jane@example.com'
```

Github Glossary

This [glossary](#) introduces common Git and GitHub terminology.

Basics

<code>git init <directory></code>	Create empty Git repository in specified directory.
<code>git clone <repository></code>	Clone a repository located at <code><repository></code> on your local machine.
<code>git config user.name <username></code>	Define author name to be used for all commits in current repository.
<code>git add <directory></code>	Stage all changes in <code><directory></code> for the next commit.
<code>git commit -m <"message"></code>	Commit the staged snapshot, but instead of launching a text editor, use <code><"message"></code> as the commit message.
<code>git status</code>	List which files are staged, unstaged, and untracked.
<code>git log</code>	Display the entire commit history using the default format.
<code>git diff</code>	Show unstaged changes between your index and working directory.

Remote Repositories

<code>git remote add <name> <url></code>	Create a new connection to a remote repository. After adding a remote, you can use <code><name></code> as a shortcut for <code><url></code> in other commands.
<code>git fetch <remote> <branch></code>	Fetches a specific <code><branch></code> , from the repository. Leave off <code><branch></code> to fetch all remote refs.
<code>git pull <remote></code>	Fetch the specified remote's copy of current branch and immediately merge it into the local copy.
<code>git push <remote> <branch></code>	Push the branch to <code><remote></code> , along with necessary commits and objects. Creates named branch in the remote repository if it doesn't exist.

Undoing Changes

<code>git revert <commit></code>	Create new commit that undoes all of the changes made in <code><commit></code> , then apply it to the current branch.
<code>git reset <file></code>	Remove <code><file></code> from the staging area but leave the working directory unchanged. This unstages a file without overwriting any changes.
<code>git clean -n</code>	Shows which files would be removed from working directory. Use the <code>-f</code> flag in place of the <code>-n</code> flag to execute the clean.

Rewriting Git History

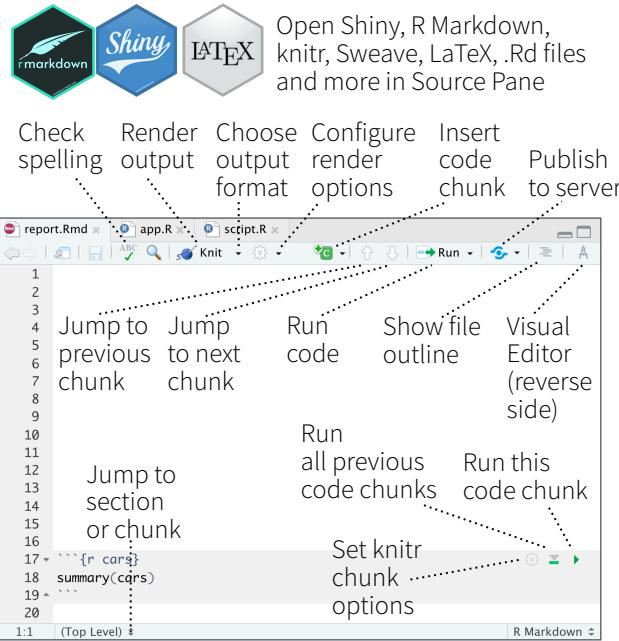
<code>git commit --amend</code>	Replace the last commit with the staged changes and last commit combined. Use with nothing staged to edit the last commit's message.
<code>git rebase <base></code>	Rebase the current branch onto <code><base></code> . <code><base></code> can be a commit ID, branch name, a tag, or a relative reference to HEAD.
<code>git reflog</code>	Show a log of changes to the local repository's HEAD. Add <code>--relative-date</code> flag to show date info or <code>--all</code> to show all refs.

Git Branches

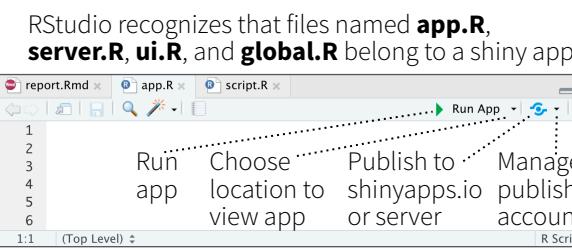
<code>git branch</code>	List all of the branches in your repo. Add a <code><branch></code> argument to create a new branch with the name <code><branch></code> .
<code>git checkout -b <branch></code>	Create and check out a new named <code><branch></code> . Drop the <code>-b</code> flag to checkout an existing branch.
<code>git merge <branch></code>	Merge <code><branch></code> into the current branch.

RStudio IDE :: CHEATSHEET

Documents and Apps



Access markdown guide at [Help > Markdown Quick Reference](#)
See reverse side for more on [Visual Editor](#)

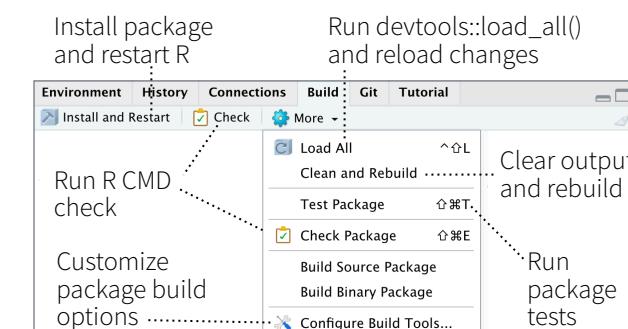


Package Development

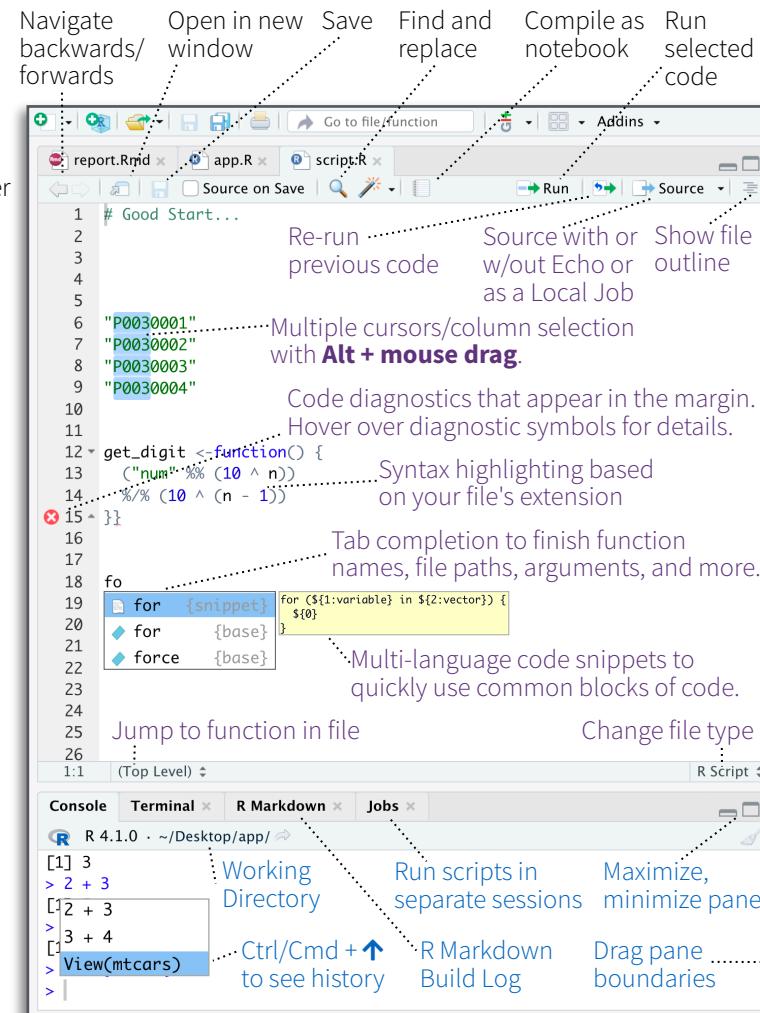
Create a new package with [File > New Project > New Directory > R Package](#)
Enable roxygen documentation with [Tools > Project Options > Build Tools](#)

Roxygen guide at [Help > Roxygen Quick Reference](#)

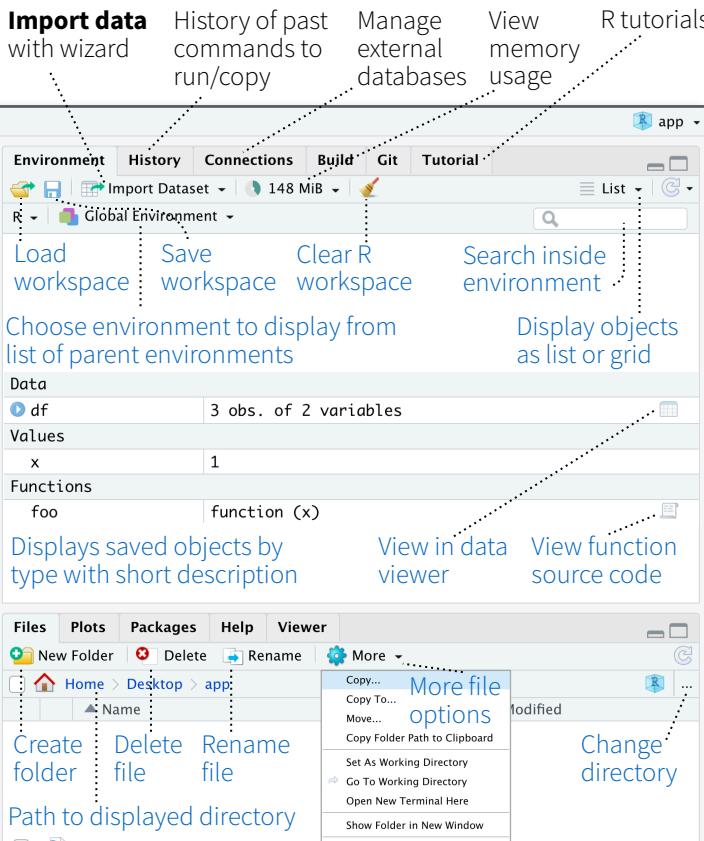
See package information in the [Build Tab](#)



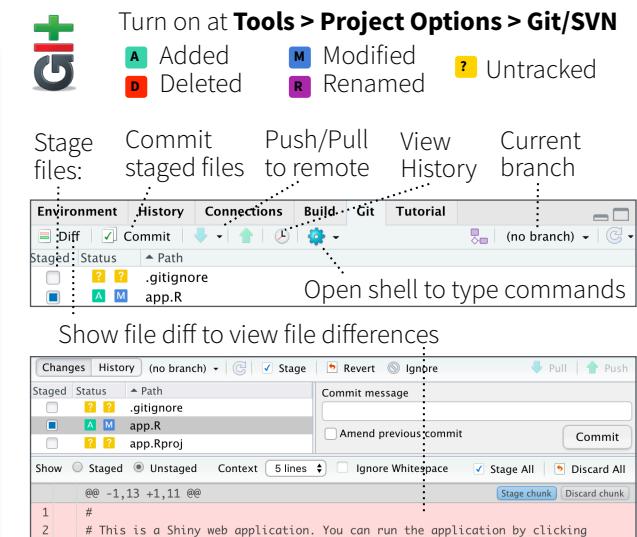
Source Editor



Tab Panes

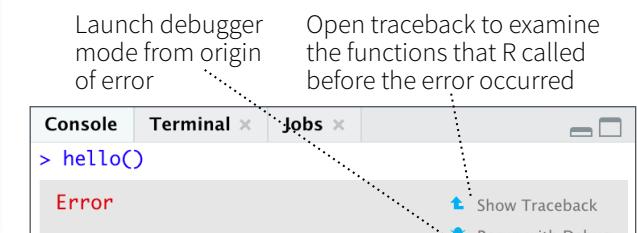


Version Control



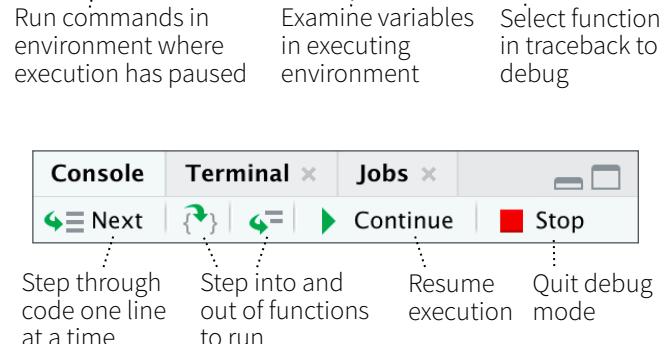
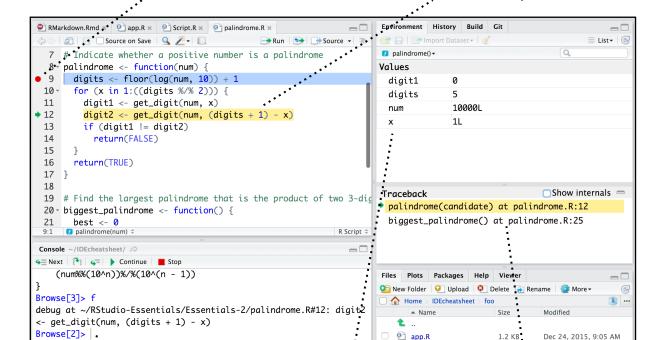
Debug Mode

Use **debug()**, **browser()**, or a breakpoint and execute your code to open the debugger mode.



Click next to line number to add/remove a breakpoint.

Highlighted line shows where execution has paused



Keyboard Shortcuts

RUN CODE

Search command history
Interrupt current command
Clear console

	Windows/Linux	Mac
Search command history	Ctrl+↑	Cmd+↑
Interrupt current command	Esc	Esc
Clear console	Ctrl+L	Ctrl+L

Navigate Code

Go to File/Function

Ctrl+. .

Ctrl+. .

Write Code

Attempt completion

Insert <- (assignment operator)
Insert |> or %>% (pipe operator)
(Un)Comment selection

	Windows/Linux	Mac
Insert <- (assignment operator)	Alt+-	Option+-
Insert > or %>% (pipe operator)	Ctrl+Shift+M	Cmd+Shift+M
(Un)Comment selection	Ctrl+Shift+C	Cmd+Shift+C

MAKE PACKAGES

Load All (devtools)
Test Package (Desktop)
Document Package

	Windows/Linux	Mac
Load All (devtools)	Ctrl+Shift+L	Cmd+Shift+L
Test Package (Desktop)	Ctrl+Shift+T	Cmd+Shift+T
Document Package	Ctrl+Shift+D	Cmd+Shift+D

DOCUMENTS AND APPS

Knit Document (knitr)
Insert chunk (Sweave & Knitr)
Run from start to current line

Ctrl+Shift+K
Ctrl+Alt+I
Ctrl+Alt+B

Cmd+Shift+K
Cmd+Option+I
Cmd+Option+B

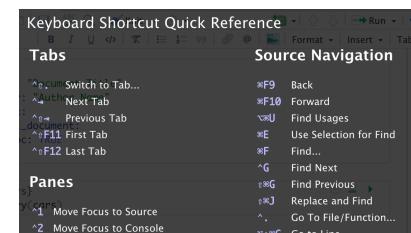
MORE KEYBOARD SHORTCUTS

Keyboard Shortcuts Help
Show Command Palette

Alt+Shift+K
Ctrl+Shift+P

Option+Shift+K
Cmd+Shift+P

View the Keyboard Shortcut Quick Reference with **Tools > Keyboard Shortcuts** or **Alt/Option + Shift + K**



Search for keyboard shortcuts with **Tools > Show Command Palette** or **Ctrl/Cmd + Shift + P**.

RStudio Workbench



WHY RSTUDIO WORKBENCH?

Extend the open source server with a commercial license, support, and more:

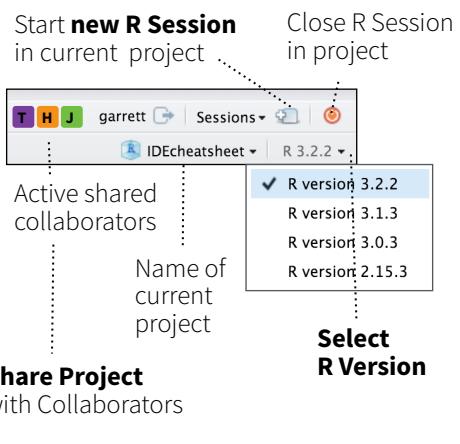
- open and run multiple R sessions at once
- tune your resources to improve performance
- administrative tools for managing user sessions
- collaborate real-time with others in shared projects
- switch easily from one version of R to a different version
- integrate with your authentication, authorization, and audit practices
- work in the RStudio IDE, JupyterLab, Jupyter Notebooks, or VS Code

Download a free 45 day evaluation at www.rstudio.com/products/workbench/evaluation/

Share Projects

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.



Visual Editor

The screenshot shows the RStudio Visual Editor interface with several keyboard shortcuts overlaid. The editor is displaying an R Markdown document with code and text. Key features shown include:

- Check spelling
- Render output
- Choose output format
- Choose output location
- Insert code chunk
- Jump to previous chunk
- Jump to next chunk
- Run selected lines
- Publish to server
- Show file outline
- Back to Source Editor (front page)
- File outline
- Add/Edit attributes
- Set knitr chunk options
- Run this and all previous code chunks
- Run this code chunk

Below the editor, there's a preview pane showing R code and its rendered output.

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents.

```
{r cars}  
summary(cars)
```

Jump to chunk or header

Run Remote Jobs

Run R on remote clusters (Kubernetes/Slurm) via the Job Launcher

The screenshot shows the RStudio Job Launcher interface. It displays a list of jobs and their status:

Job	Status	Time	Location	
fast.R	Running	Local	0:09	
sleepy.R	Succeeded	11:22 AM	Local	0:41
sleepy.R	Idle		KubernetesX	Waiting

Below the list, there's a "Launcher" tab showing the status of launcher jobs:

Job	Status	Time	Location	
Start Launcher Job	Sorted by submission time			
fast.R	Running	Local	0:09	
sleepy.R	Succeeded	11:22 AM	Local	0:41
sleepy.R	Idle		KubernetesX	Waiting

Shiny :: CHEAT SHEET



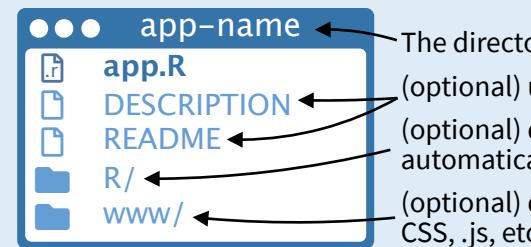
Building an App

A **Shiny** app is a web page (**ui**) connected to a computer running a live R session (**server**).



Users can manipulate the UI, which will cause the server to update the UI's displays (by running R code).

Save your template as **app.R**. Keep your app in a directory along with optional extra files.



Launch apps stored in a directory with **runApp(<path to directory>)**.

Share

Share your app in three ways:

1. **Host it on shinyapps.io**, a cloud based service from RStudio. To deploy Shiny apps:

Create a free or professional account at [shinyapps.io](#)

Click the Publish icon in RStudio IDE, or run: `rsconnect::deployApp("<path to directory>")`

2. **Purchase RStudio Connect**, a publishing platform for R and Python. [rstudio.com/products/connect/](#)

3. **Build your own Shiny Server** [rstudio.com/products/shiny/shiny-server/](#)

To generate the template, type **shinyapp** and press **Tab** in the RStudio IDE or go to **File > New Project > New Directory > Shiny Web Application**

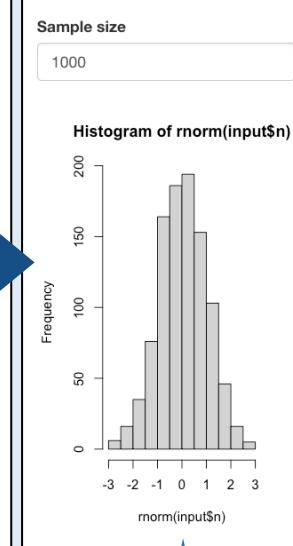
```
# app.R
library(shiny)

ui <- fluidPage(
  numericInput(inputId = "n",
    "Sample size", value = 25),
  plotOutput(outputId = "hist")
)

server <- function(input, output, session) {
  output$hist <- renderPlot({
    hist(rnorm(input$n))
  })
}

shinyApp(ui = ui, server = server)
```

Call **shinyApp()** to combine **ui** and **server** into an interactive app!



See annotated examples of Shiny apps by running **runExample(<example name>)**. Run **runExample()** with no arguments for a list of example names.

Inputs

Collect values from the user.

Access the current value of an input object with **input\$<inputId>**. Input values are **reactive**.

Action

ActionButton(inputId, label, icon, width, ...)

Link

actionLink(inputId, label, icon, ...)

checkbox

checkboxGroupInput(inputId, label, choices, selected, inline, width, choiceNames, choiceValues)

checkbox

checkboxInput(inputId, label, value, width)

date

dateInput(inputId, label, value, min, max, format, startview, weekstart, language, width, autoclose, datesdisabled, daysofweekdisabled)

dateRange

dateRangeInput(inputId, label, start, end, min, max, format, startview, weekstart, language, separator, width, autoclose)

file

fileInput(inputId, label, multiple, accept, width, buttonLabel, placeholder)

number

numericInput(inputId, label, value, min, max, step, width)

password

passwordInput(inputId, label, value, width, placeholder)

radio

radioButtons(inputId, label, choices, selected, inline, width, choiceNames, choiceValues)

select

selectInput(inputId, label, choices, selected, multiple, selectize, width, size)
Also **selectizeInput()**

slider

sliderInput(inputId, label, min, max, value, step, round, format, locale, ticks, animate, width, sep, pre, post, timeFormat, timezone, dragRange)

submit

submitButton(text, icon, width)
(Prevent reactions for entire app)

text

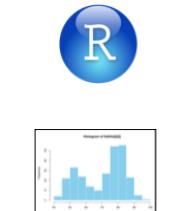
textInput(inputId, label, value, width, placeholder)
Also **textAreaInput()**

Outputs

render*() and ***Output()** functions work together to add R output to the UI.



DT::renderDataTable(expr, options, searchDelay, callback, escape, env, quoted, outputArgs)



renderImage(expr, env, quoted, deleteFile, outputArgs)



renderPrint(expr, env, quoted, width, outputArgs)



renderTable(expr, striped, hover, bordered, spacing, width, align, rownames, colnames, digits, na, ..., env, quoted, outputArgs)



renderText(expr, env, quoted, outputArgs, sep)

renderUI(expr, env, quoted, outputArgs)

dataTableOutput(outputId)

imageOutput(outputId, width, height, click, dblclick, hover, brush, inline)

plotOutput(outputId, width, height, click, dblclick, hover, brush, inline)

verbatimTextOutput(outputId, placeholder)

tableOutput(outputId)

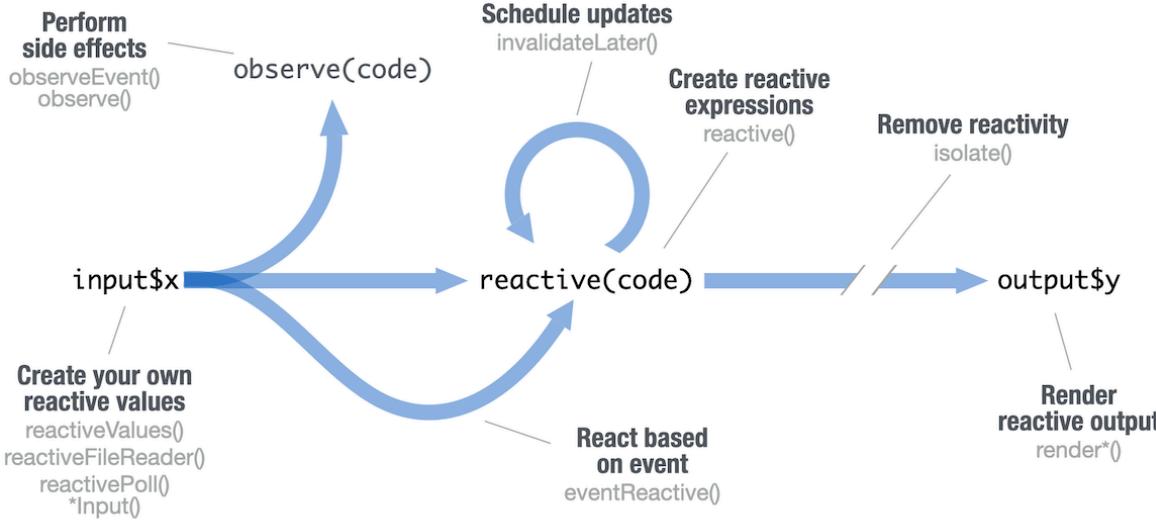
textOutput(outputId, container, inline)

uiOutput(outputId, inline, container, ...)
htmlOutput(outputId, inline, container, ...)

These are the core output types. See [htmlwidgets.org](#) for many more options.

Reactivity

Reactive values work together with reactive functions. Call a reactive value from within the arguments of one of these functions to avoid the error **Operation not allowed without an active reactive context**.



CREATE YOUR OWN REACTIVE VALUES

```
# *Input() example
ui <- fluidPage(
  textInput("a", "", "A")
)
```

```
#reactiveValues example
server <-
function(input, output){
  rv <- reactiveValues()
  rv$number <- 5
}
```

*Input() functions (see front page)

Each input function creates a reactive value stored as **input\$<inputId>**.

reactiveValues(...)

Creates a list of reactive values whose values you can set.

CREATE REACTIVE EXPRESSIONS

```
library(shiny)
ui <- fluidPage(
 textInput("a", "", "A"),
 textInput("z", "", "Z"),
  textOutput("b"))
server <-
function(input, output){
  re <- reactive({
    paste(input$a, input$z)
  })
  output$b <- renderText({
    re()
  })
}
shinyApp(ui, server)
```

reactive(x, env, quoted, label, domain)

Reactive expressions:

- cache their value to reduce computation
 - can be called elsewhere
 - notify dependencies when invalidated
- Call the expression with function syntax, e.g. **re()**.

PERFORM SIDE EFFECTS

```
library(shiny)
ui <- fluidPage(
  textInput("a", "", "A"),
  actionButton("go", "Go"))
server <-
function(input, output){
  observeEvent(input$go, {
    print(input$a)
  })
}
shinyApp(ui, server)
```

observeEvent(eventExpr, handlerExpr, event.env, event.quoted, handler.env, handler.quoted, ..., label, suspended, priority, domain, autoDestroy, ignoreNULL, ignoreInit, once)

Runs code in 2nd argument when reactive values in 1st argument change. See **observe()** for alternative.

REACT BASED ON EVENT

```
library(shiny)
ui <- fluidPage(
  textInput("a", "", "A"),
  actionButton("go", "Go"),
  textOutput("b"))
server <-
function(input, output){
  re <- eventReactive(
    input$go, {input$a})
  output$b <- renderText({
    re()
  })
}
shinyApp(ui, server)
```

eventReactive(eventExpr, valueExpr, event.env, event.quoted, value.env, value.quoted, ..., label, domain, ignoreNULL, ignoreInit)

Creates reactive expression with code in 2nd argument that only invalidates when reactive values in 1st argument change.

REMOVE REACTIVITY

```
library(shiny)
ui <- fluidPage(
  textInput("a", "", "A"),
  textOutput("b"))
server <-
function(input, output){
  output$b <- renderText({
    isolate({input$a})
  })
}
shinyApp(ui, server)
```

isolate(expr)

Runs a code block. Returns a **non-reactive** copy of the results.

UI - An app's UI is an HTML document.

Use Shiny's functions to assemble this HTML with R.

```
fluidPage(
  textInput("a", ""))
## <div class="container-fluid">
##   <div class="form-group shiny-input-container">
##     <label for="a"></label>
##     <input id="a" type="text"
##           class="form-control" value="" />
##   </div>
## </div>
```

Returns HTML

Add static HTML elements with **tags**, a list of functions that parallel common HTML tags, e.g. **tags\$a()**. Unnamed arguments will be passed into the tag; named arguments will become tag attributes.

Run **names(tags)** for a complete list.
tags\$h1("Header") → `<h1>Header</h1>`

The most common tags have wrapper functions. You do not need to prefix their names with **tags\$**

```
ui <- fluidPage(
  h1("Header 1"),
  hr(),
  br(),
  p(strong("bold")),
  p(em("italic")),
  p(code("code")),
  a(href="", "link"),
  HTML("<p>Raw html</p>"))
)
```



To include a CSS file, use **includeCSS()**, or 1. Place the file in the **www** subdirectory 2. Link to it with:

```
tags$head(tags$link(rel = "stylesheet",
  type = "text/css", href = "<file name>"))
```

To include JavaScript, use **includeScript()** or 1. Place the file in the **www** subdirectory 2. Link to it with:

```
tags$head(tags$script(src = "<file name>"))
```

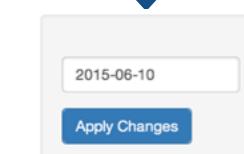
To include an image:

1. Place the file in the **www** subdirectory
2. Link to it with `img(src = "<file name>")`

Layouts

Combine multiple elements into a "single element" that has its own properties with a panel function, e.g.

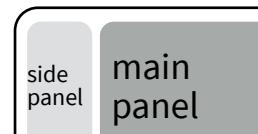
```
wellPanel(
  dateInput("a", ""),
  submitButton()
)
```



```
absolutePanel()
conditionalPanel()
fixedPanel()
headerPanel()
inputPanel()
mainPanel()
navlistPanel()
sidebarPanel()
tabPanel()
tabsetPanel()
titlePanel()
wellPanel()
```

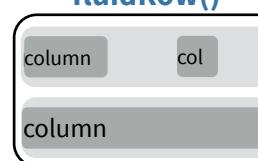
Organize panels and elements into a layout with a layout function. Add elements as arguments of the layout functions.

sidebarLayout()



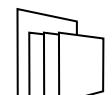
```
ui <- fluidPage(
  sidebarLayout(
    sidebarPanel(),
    mainPanel()
  )
)
```

fluidRow()



```
ui <- fluidPage(
  fluidRow(column(width = 4),
    column(width = 2, offset = 3)),
  fluidRow(column(width = 12))
)
```

Also **flowLayout()**, **splitLayout()**, **verticalLayout()**, **fixedPage()**, and **fixedRow()**.



Layer tabPanels on top of each other, and navigate between them, with:



```
ui <- fluidPage( tabsetPanel(
  tabPanel("tab 1", "contents"),
  tabPanel("tab 2", "contents"),
  tabPanel("tab 3", "contents"))
)
```



```
ui <- fluidPage( navlistPanel(
  tabPanel("tab 1", "contents"),
  tabPanel("tab 2", "contents"),
  tabPanel("tab 3", "contents"))
)
```

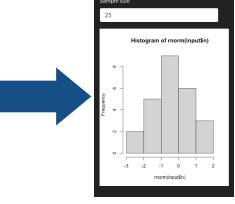


```
ui <- navbarPage(title = "Page",
  tabPanel("tab 1", "contents"),
  tabPanel("tab 2", "contents"),
  tabPanel("tab 3", "contents"))
```

Themes

Use the **bslib** package to add existing themes to your Shiny app ui, or make your own.

```
library(bslib)
ui <- fluidPage(
  theme = bs_theme(
    bootswatch = "darkly",
    ...
  )
)
```

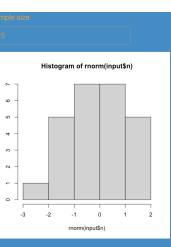


bootswatch_themes() Get a list of themes.

Build your own theme by customizing individual arguments.

```
bs_theme(bg = "#558AC5",
  fg = "#F9B02D",
  ...)
```

?**bs_theme** for a full list of arguments.



bs_themer() Place within the server function to use the interactive theming widget.



rmarkdown :: CHEAT SHEET

What is rmarkdown?



.Rmd files • Develop your code and ideas side-by-side in a single document. Run code as individual chunks or as an entire document.

Dynamic Documents • Knit together plots, tables, and results with narrative text. Render to a variety of formats like HTML, PDF, MS Word, or MS Powerpoint.

Reproducible Research • Upload, link to, or attach your report to share. Anyone can read or run your code to reproduce your work.

Workflow

- 1 Open a **new .Rmd file** in the RStudio IDE by going to *File > New File > R Markdown*.
- 2 **Embed code** in chunks. Run code by line, by chunk, or all at once.
- 3 **Write text** and add tables, figures, images, and citations. Format with Markdown syntax or the RStudio Visual Markdown Editor.
- 4 **Set output format(s) and options** in the YAML header. Customize themes or add parameters to execute or add interactivity with Shiny.
- 5 **Save and render** the whole document. Knit periodically to preview your work as you write.
- 6 **Share your work!**

Embed Code with knitr

CODE CHUNKS

Surround code chunks with `{{r}}` and `{{` or use the Insert Code Chunk button. Add a chunk label and/or chunk options inside the curly braces after **r**.

```
```{r chunk-label, include=FALSE}
summary(mtcars)
```
```

SET GLOBAL OPTIONS

Set options for the entire document in the first chunk.

```
```{r include=FALSE}
knitr::opts_chunk$message = FALSE
```
```

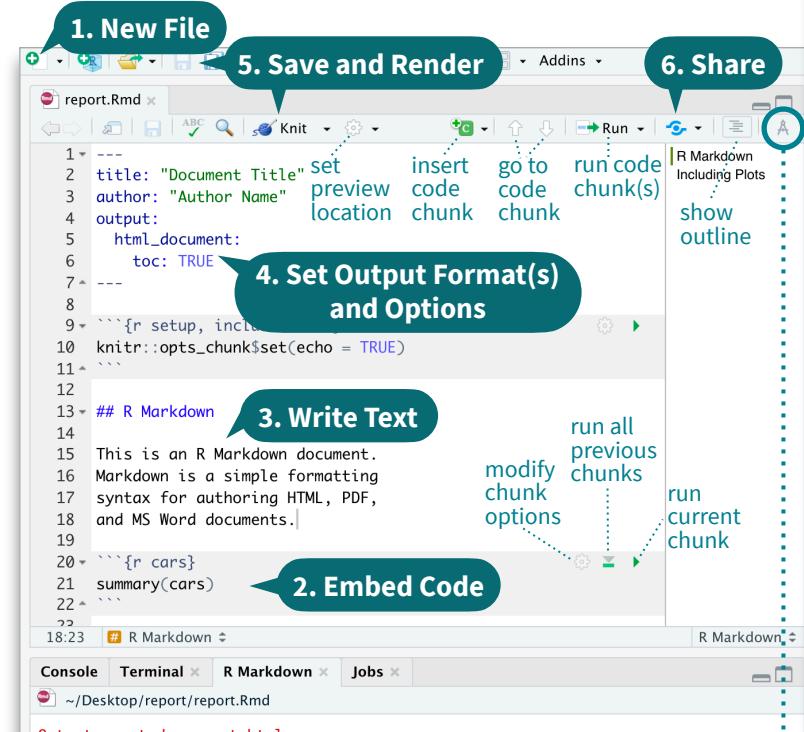
INLINE CODE

Insert `r <code>` into text sections. Code is evaluated at render and results appear as text.

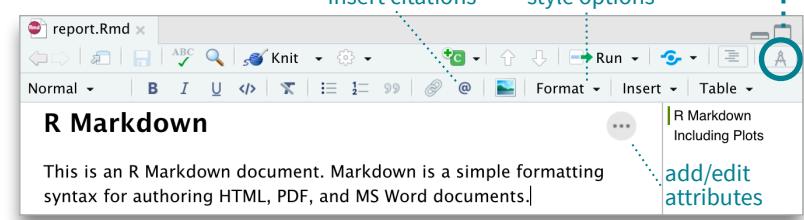
"Built with `r getRversion()`" --> "Built with 4.1.0"



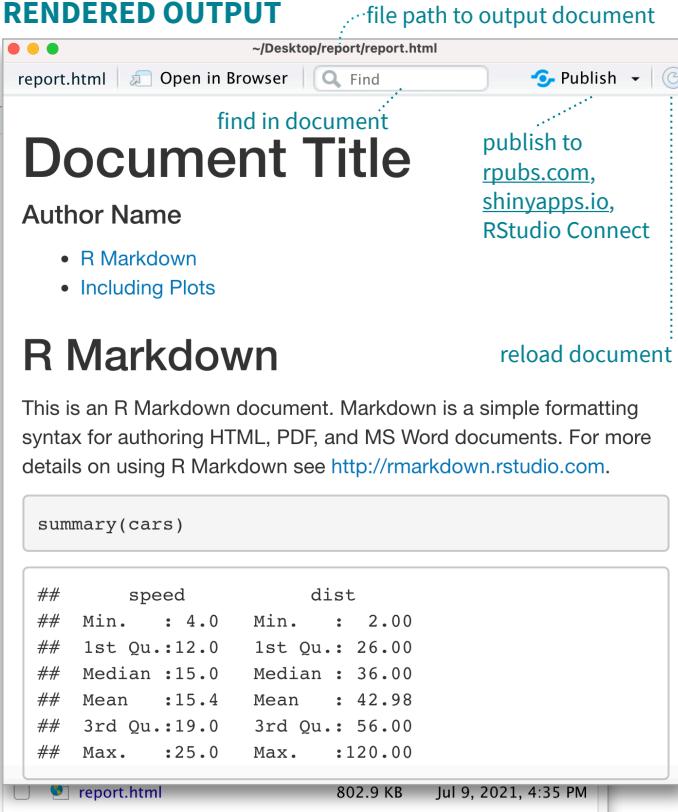
SOURCE EDITOR



VISUAL EDITOR



RENDERED OUTPUT



Write with Markdown

The syntax on the left renders as the output on the right.

Plain text.

Plain text.

End a line with two spaces to start a new paragraph.

End a line with two spaces to start a new paragraph.

Also end with a backslash\ to make a new line.

Also end with a backslash\ to make a new line.

italics* and ****bold****

italics and **bold**

superscript²/subscript₂

superscript²/subscript₂

~~strikethrough~~

strikethrough

escaped: `*` \`

escaped: * \`

endash: --, emdash: ---

endash: -, emdash: -

Header 1 Header 2

...

Header 6

- unordered list

• item 2

- item 2a (indent 1 tab)

• item 2b

1. ordered list

2. item 2

- item 2a (indent 1 tab)

• item 2b

<link url>

[This is a link.](link url)

[This is another link][id].

This is another link.

<http://www.rstudio.com/>

This is a link.

This is another link.



Caption.

verbatim code

multiple lines of verbatim code

block quotes

multiple lines of verbatim code

block quotes

equation: $e^{i\pi} + 1 = 0$

equation block:

$$E = mc^2$$

horizontal rule:

| Right | Left | Default | Center |
|-------|------|---------|--------|
| 12 | 12 | 12 | 12 |
| 123 | 123 | 123 | 123 |
| 1 | 1 | 1 | 1 |

HTML Tabsets

```
# Results {.tabset}
## Plots text
text
```

Tables
more text

Results

| Plots | Tables |
|-------|--------|
| | |

| OPTION | DEFAULT | EFFECTS |
|-----------------------------------|-----------|---|
| echo | TRUE | display code in output document |
| error | FALSE | TRUE (display error messages in doc)
FALSE (stop render when error occurs) |
| eval | TRUE | run code in chunk |
| include | TRUE | include chunk in doc after running |
| message | TRUE | display code messages in document |
| warning | TRUE | display code warnings in document |
| results | "markup" | "asis" (passthrough results)
"hide" (don't display results)
"hold" (put all results below all code) |
| fig.align | "default" | "left", "right", or "center" |
| fig.alt | NULL | alt text for a figure |
| fig.cap | NULL | figure caption as a character string |
| fig.path | "figure/" | prefix for generating figure file paths |
| fig.width & fig.height | 7 | plot dimensions in inches |
| out.width | | rescales output width, e.g. "75%", "300px" |
| collapse | FALSE | collapse all sources & output into a single block |
| comment | "##" | prefix for each line of results |
| child | NULL | files(s) to knit and then include |
| purl | TRUE | include or exclude a code chunk when extracting source code with knitr::purl() |

See more options and defaults by running `str(knitr::opts_chunk$get())`



Set Output Formats and their Options in YAML

Use the document's YAML header to set an **output format** and customize it with **output options**.

```
---
```

```
title: "My Document"
author: "Author Name"
output:
  html_document:
    toc: TRUE
---
```

**Indent format 2 characters,
indent options 4 characters**

| OUTPUT FORMAT | CREATES |
|-------------------------|------------------------------|
| html_document | .html |
| pdf_document* | .pdf |
| word_document | Microsoft Word (.docx) |
| powerpoint_presentation | Microsoft Powerpoint (.pptx) |
| odt_document | OpenDocument Text |
| rtf_document | Rich Text Format |
| md_document | Markdown |
| github_document | Markdown for Github |
| ioslides_presentation | ioslides HTML slides |
| slidy_presentation | Slidy HTML slides |
| beamer_presentation* | Beamer slides |

* Requires LaTeX, use `tinytex::install_tinytex()`
Also see `flexdashboard`, `bookdown`, `distill`, and `blogdown`.

| IMPORTANT OPTIONS | DESCRIPTION | HTML | PDF | MS Word | MS PPT |
|---------------------|--|---------|-----|---------|--------|
| anchor_sections | Show section anchors on mouse hover (TRUE or FALSE) | X | | | |
| citation_package | The LaTeX package to process citations ("default", "natbib", "biblatex") | X | | | |
| code_download | Give readers an option to download the .Rmd source code (TRUE or FALSE) | X | | | |
| code_folding | Let readers to toggle the display of R code ("none", "hide", or "show") | X | | | |
| css | CSS or SCSS file to use to style document (e.g. "style.css") | X | | | |
| dev | Graphics device to use for figure output (e.g. "png", "pdf") | X X | | | |
| df_print | Method for printing data frames ("default", "kable", "tibble", "paged") | X X X X | | | |
| fig_caption | Should figures be rendered with captions (TRUE or FALSE) | X X X X | | | |
| highlight | Syntax highlighting ("tango", "pygments", "kate", "zenburn", "textmate") | X X X | | | |
| includes | File of content to place in doc ("in_header", "before_body", "after_body") | X X | | | |
| keep_md | Keep the Markdown .md file generated by knitting (TRUE or FALSE) | X X X X | | | |
| keep_tex | Keep the intermediate TEX file used to convert to PDF (TRUE or FALSE) | X | | | |
| latex_engine | LaTeX engine for producing PDF output ("pdflatex", "xelatex", or "lualatex") | X | | | |
| reference_docx/_doc | docx/pptx file containing styles to copy in the output (e.g. "file.docx", "file.pptx") | X X | | | |
| theme | Theme options (see Bootswatch and Custom Themes below) | X | | | |
| toc | Add a table of contents at start of document (TRUE or FALSE) | X X X X | | | |
| toc_depth | The lowest level of headings to add to table of contents (e.g. 2, 3) | X X X X | | | |
| toc_float | Float the table of contents to the left of the main document content (TRUE or FALSE) | X | | | |

Use `?<output format>` to see all of a format's options, e.g. `?html_document`

More Header Options

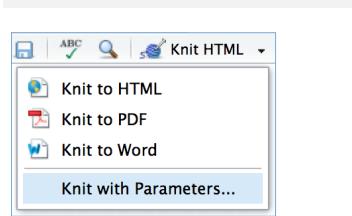
PARAMETERS

Parameterize your documents to reuse with new inputs (e.g., data, values, etc.).

1. **Add parameters** in the header as sub-values of `params`.
2. **Call parameters** in code using `params$<name>`.
3. **Set parameters** with Knit with Parameters or the `params` argument of `render()`.

REUSABLE TEMPLATES

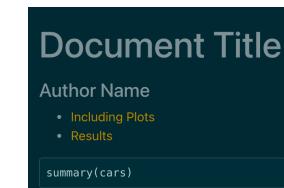
1. **Create a new package** with a `inst/rmarkdown/templates` directory.
2. **Add a folder** containing `template.yaml` (below) and `skeleton.Rmd` (template contents).
3. **Install** the package to access template by going to **File > New R Markdown > From Template**.



BOOTSWATCH THEMES

Customize HTML documents with Bootswatch themes from the `bslib` package using the theme output option.

Use `bslib::bootswatch_themes()` to list available themes.



```
---
```

```
title: "Document Title"
author: "Author Name"
output:
  html_document:
    theme:
      bootswatch: solar
---
```

CUSTOM THEMES

Customize individual HTML elements using `bslib` variables. Use `?bs_theme` to see more variables.

```
---
```

```
output:
  html_document:
    theme:
      bg: "#121212"
      fg: "#E4E4E4"
      base_font:
        google: "Prompt"
---
```

More on `bslib` at pkgs.rstudio.com/bslib/.

STYLING WITH CSS AND SCSS

Add CSS and SCSS to your document by adding a path to a file with the `css` option in the YAML header.

```
---
```

```
title: "My Document"
author: "Author Name"
output:
  html_document:
    css: "style.css"
---
```

Apply CSS styling by writing HTML tags directly or:

- Use markdown to apply style attributes inline.

Bracketed Span
A [green]{.my-color} word.

A green word.

Fenced Div
:::{.my-color}
All of these words
are green.
:::

All of these words
are green.

- Use the Visual Editor. Go to **Format > Div/Span** and add CSS styling directly with Edit Attributes.

.my-css-tag ...
This is a div with some text in it.

Render

When you render a document, rmarkdown:

1. Runs the code and embeds results and text into an .md file with knitr.
2. Converts the .md file into the output format with Pandoc.



Save, then **Knit** to preview the document output. The resulting HTML/PDF/MS Word/etc. document will be created and saved in the same directory as the .Rmd file.

Use `rmarkdown::render()` to render/knit in the R console. See `?render` for available options.

Share

Publish on RStudio Connect

to share R Markdown documents securely, schedule automatic updates, and interact with parameters in real time.

rstudio.com/products/connect/



INTERACTIVITY

Turn your report into an interactive Shiny document in 4 steps:

1. Add `runtime: shiny` to the YAML header.
2. Call Shiny input functions to embed input objects.
3. Call Shiny render functions to embed reactive output.
4. Render with `rmarkdown::run()` or click **Run Document** in RStudio IDE.

```
---
```

```
output: html_document
runtime: shiny
---
```

```
```{r, echo = FALSE}
numericInput("n",
 "How many cars?", 5)
renderTable({
 head(cars, input$n)
})
```

speed	dist
1	4.00
2	4.00
3	7.00
4	7.00
5	8.00
	2.00
	10.00
	4.00
	22.00
	16.00

Also see Shiny Prerendered for better performance.  
[rmarkdown.rstudio.com/authoring\\_shiny\\_prerendered](https://rmarkdown.rstudio.com/authoring_shiny_prerendered)

Embed a complete app into your document with `shiny::shinyAppDir()`. More at [bookdown.org/yihui/rmarkdown/shiny-embedded.html](https://bookdown.org/yihui/rmarkdown/shiny-embedded.html).