

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.



`m[2,]` - Select a row



`m[, 1]` - Select a column



`m[2, 3]` - Select an element

`t(m)`

Transpose

`m %*% n`

Matrix Multiplication

`solve(m, n)`

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.

`l[[2]]`

Second element of l.

`l[1]`

New list with only the first element.

`l$x`

Element named x.

`l['y']`

New list with only element named y.

Also see the `dplyr` package.

Data Frames

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

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

x	y
1	a
2	b
3	c

Matrix subsetting

`df[, 2]`



`df[2,]`



`df[2, 2]`



List subsetting



Understanding a data frame

`View(df)`

See the full data frame.

`head(df)`

See the first 6 rows.

`nrow(df)`

Number of rows.

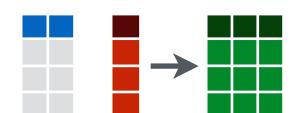
`ncol(df)`

Number of columns.

`dim(df)`

Number of columns and rows.

`cbind` - Bind columns.



`rbind` - Bind rows.



Strings

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

Join multiple vectors together.

`paste(x, collapse = ' ')`

Join elements of a vector together.

`grep(pattern, x)`

Find regular expression matches in x.

`gsub(pattern, replace, x)`

Replace matches in x with a string.

`toupper(x)`

Convert to uppercase.

`tolower(x)`

Convert to lowercase.

`nchar(x)`

Number of characters in a string.

Factors

`factor(x)`

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

`cut(x, breaks = 4)`

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

Statistics

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

Linear model.

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

Generalised linear model.

`summary`

Get more detailed information out a model.

`t.test(x, y)`

Perform a t-test for difference between means.

`pairwise.t.test`

Perform a t-test for paired data.

`aov`

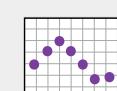
Analysis of variance.

Distributions

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

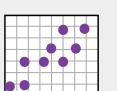
Plotting

Also see the `ggplot2` package.



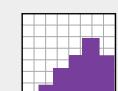
`plot(x)`

Values of x in order.



`plot(x, y)`

Values of x against y.



`hist(x)`

Histogram of x.

Dates

See the `lubridate` package.

RStudio IDE :: CHEAT SHEET

Documents and Apps



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

Access markdown guide at **Help > Markdown Quick Reference**

Jump to chunk Set knitr chunk options Run this and all previous code chunks Run this code chunk

RStudio recognizes that files named **app.R**, **server.R**, **ui.R**, and **global.R** belong to a shiny app

Run app Choose location to view app Publish to shinyapps.io or server Manage publish accounts

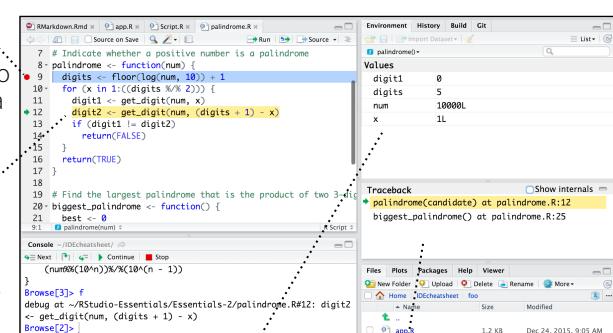
Debug Mode

Open with **debug()**, **browser()**, or a breakpoint. RStudio will open the debugger mode when it encounters a breakpoint while executing code.

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

Highlighted line shows where execution has paused

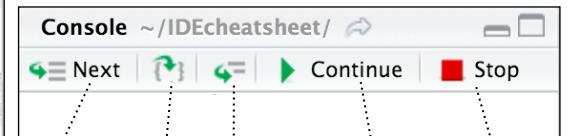
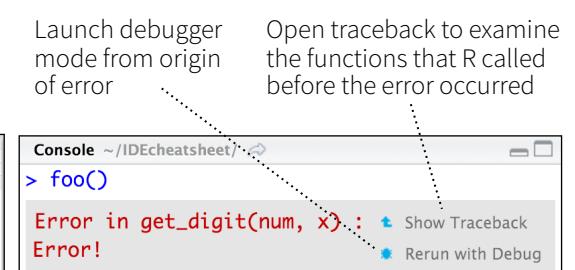
Run commands in environment where execution has paused



Examine variables in executing environment

Select function in traceback to debug

Launch debugger mode from origin of error



Step through code one line at a time

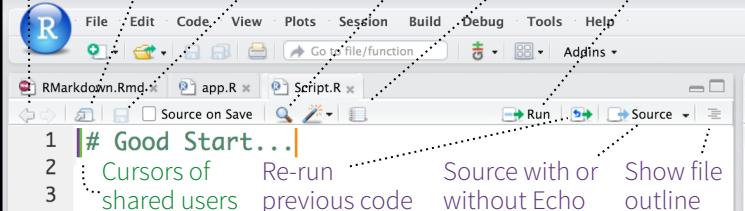
Step into and out of functions to run

Resume execution mode

Quit debug

Write Code

Navigate tabs Open in new window Save Find and replace Compile as notebook Run selected code



Cursors of shared users Re-run previous code Source with or without Echo Show file outline

Multiple cursors/column selection with **Alt + mouse drag**.

Code diagnostics that appear in the margin. Hover over diagnostic symbols for details.

Syntax highlighting based on your file's extension

Tab completion to finish function names, file paths, arguments, and more.

Multi-language code snippets to quickly use common blocks of code.

Jump to function in file

Change file type

Working Directory

Press **↑** to see command history

Maximize, minimize panes

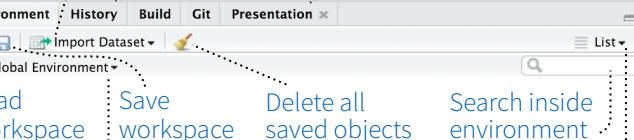
Drag pane boundaries

R Support

Import data with wizard

History of past commands to run/copy

Display .RPres slideshows **File > New File > R Presentation**



Choose environment to display from list of parent environments

Display objects as list or grid

Data 150 obs. of 5 variables
Values a 1
Functions foo function (x)

Displays saved objects by type with short description

View in data viewer View function source code



Path to displayed directory

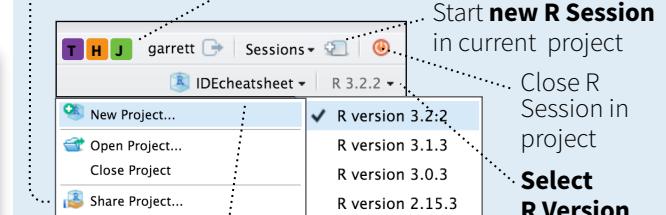
Maximize, minimize panes

Drag pane boundaries

A File browser keyed to your working directory. Click on file or directory name to open.

Pro Features

Share Project Active shared with Collaborators



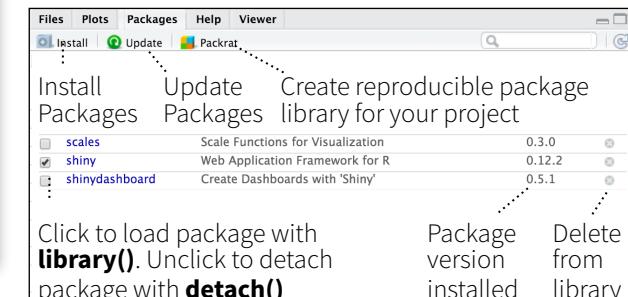
PROJECT SYSTEM **File > New Project**

RStudio saves the call history, workspace, and working directory associated with a project. It reloads each when you re-open a project.

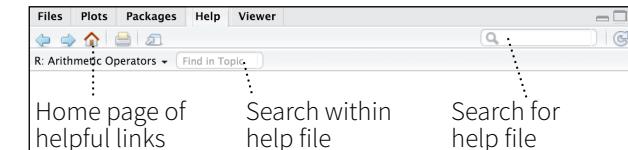
RStudio opens plots in a dedicated Plots pane



GUI Package manager lists every installed package



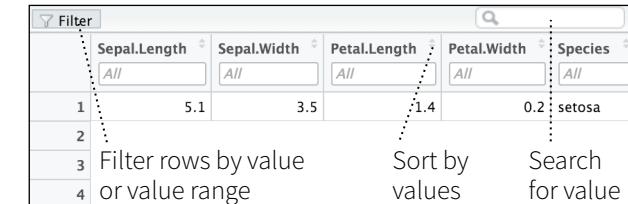
RStudio opens documentation in a dedicated Help pane



Viewer Pane displays HTML content, such as Shiny apps, RMarkdown reports, and interactive visualizations



View(<data>) opens spreadsheet like view of data set





1 LAYOUT

Move focus to Source Editor
Move focus to Console
Move focus to Help
Show History
Show Files
Show Plots
Show Packages
Show Environment
Show Git/SVN
Show Build

Windows/Linux **Mac**
Ctrl+1 Ctrl+1
Ctrl+2 Ctrl+2
Ctrl+3 Ctrl+3
Ctrl+4 Ctrl+4
Ctrl+5 Ctrl+5
Ctrl+6 Ctrl+6
Ctrl+7 Ctrl+7
Ctrl+8 Ctrl+8
Ctrl+9 Ctrl+9
Ctrl+0 Ctrl+0

2 RUN CODE

Search command history

Navigate command history
Move cursor to start of line
Move cursor to end of line
Change working directory

Interrupt current command

Clear console

Quit Session (desktop only)

Restart R Session

Run current (retain cursor)
Run from current to end
Run the current function
Source a file

Source the current file

Source with echo

Windows/Linux **Mac**

Ctrl+↑ **Cmd+↑**
↑/↓ **↑/↓**
Home **Cmd+←**
End **Cmd+→**
Ctrl+Shift+H **Ctrl+Shift+H**

Esc **Esc**
Ctrl+L **Ctrl+L**
Ctrl+Q **Cmd+Q**

Ctrl+Shift+F10 **Cmd+Shift+F10**

Ctrl+Enter **Cmd+Enter**

Alt+Enter Option+Enter

Ctrl+Alt+E Cmd+Option+E

Ctrl+Alt+F Cmd+Option+F

Ctrl+Alt+G Cmd+Option+G

Ctrl+Shift+S **Cmd+Shift+S**

Ctrl+Shift+Enter Cmd+Shift+Enter

4 WRITE CODE

Attempt completion

Navigate candidates
Accept candidate
Dismiss candidates
Undo
Redo
Cut
Copy
Paste
Select All
Delete Line

Select
Select Word

Select to Line Start

Select to Line End

Select Page Up/Down

Select to Start/End

Delete Word Left

Delete Word Right

Delete to Line End

Delete to Line Start

Indent

Outdent

Yank line up to cursor

Yank line after cursor

Insert yanked text

Insert <->

Insert %>%

Show help for function

Show source code

New document

New document (Chrome)

Open document

Save document

Close document

Close document (Chrome)

Close all documents

Extract function

Extract variable

Reindent lines

(Un)Comment lines

Reflow Comment

Reformat Selection

Select within braces

Show Diagnostics

Transpose Letters

Move Lines Up/Down

Copy Lines Up/Down

Add New Cursor Above

Add New Cursor Below

Move Active Cursor Up

Move Active Cursor Down

Find and Replace

Use Selection for Find

Replace and Find

Windows /Linux

Tab or Ctrl+Space

↑/↓
Enter, Tab, or →
Esc
Ctrl+Z
Ctrl+Shift+Z
Ctrl+X
Ctrl+C
Ctrl+V
Ctrl+A
Ctrl+D
Shift+[Arrow]
Ctrl+Shift+←→

Alt+Shift+←

Alt+Shift+→

Shift+PageUp/Down

Shift+Alt+↑↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Ctrl+Shift+M

F1

F2

Cmd+Shift+N

Ctrl+Alt+Shift+N

Ctrl+O

Ctrl+S

Cmd+W

Cmd+Option+W

Ctrl+Shift+W

Ctrl+Alt+X

Ctrl+Option+X

Ctrl+Option+V

Cmd+I

Ctrl+Shift+C

Ctrl+Shift+/

Ctrl+Shift+A

Ctrl+Shift+E

Ctrl+Shift+E

Cmd+Shift+Opt+P

Ctrl+T

Alt+↑↓

Shift+Alt+↑↓

Cmd+Option+↑↓

Ctrl+Option+Up

Ctrl+Option+Down

Ctrl+Option+Shift+Up

Ctrl+Opt+Shift+Down

Ctrl+Shift+J

Mac

Tab or Cmd+Space

↑/↓

Enter, Tab, or →

Esc

Cmd+Z

Cmd+Shift+Z

Cmd+X

Cmd+C

Cmd+V

Cmd+A

Cmd+D

Shift+[Arrow]

Option+Shift+←→

Cmd+Shift+←

Cmd+Shift+→

Shift+PageUp/Down

Shift+PageUp/Down

Ctrl+Shift+↑↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Ctrl+Shift+M

F1

F2

Cmd+Shift+N

Cmd+Shift+Opt+N

Cmd+O

Ctrl+S

Cmd+W

Cmd+Option+W

Ctrl+Shift+W

Ctrl+Alt+X

Ctrl+Option+X

Ctrl+Option+V

Cmd+I

Cmd+Shift+C

Cmd+Shift+/

Cmd+Shift+A

Cmd+Shift+E

Cmd+Shift+E

Cmd+Shift+Opt+P

Ctrl+T

Option+↑↓

Cmd+Option+↑↓

Ctrl+Option+Up

Ctrl+Option+Down

Ctrl+Option+Shift+Up

Ctrl+Opt+Shift+Down

Ctrl+Shift+J

WHY RSTUDIO SERVER PRO?

RSP extends 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
 - edit the same project at the same time as others
 - see what you and others are doing on your server
 - switch easily from one version of R to a different version
 - integrate with your authentication, authorization, and audit practices
- Download a free 45 day evaluation at www.rstudio.com/products/rstudio-server-pro/

5 DEBUG CODE

Toggle Breakpoint
Execute Next Line
Step Into Function
Finish Function/Loop
Continue
Stop Debugging

Windows/Linux Mac

Shift+F9	Shift+F9
F10	F10
Shift+F4	Shift+F4
Shift+F6	Shift+F6
Shift+F5	Shift+F5
Shift+F8	Shift+F8

6 VERSION CONTROL

Show diff
Commit changes
Scroll diff view
Stage/Unstage (Git)
Stage/Unstage and move to next

R Markdown Cheat Sheet

learn more at rmarkdown.rstudio.com



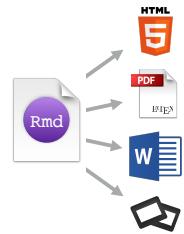
.Rmd files

An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.



Reproducible Research

At the click of a button, or the type of a command, you can rerun the code in an R Markdown file to reproduce your work and export the results as a finished report.



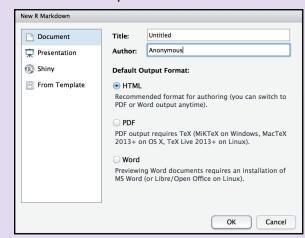
Dynamic Documents

You can choose to export the finished report as a html, pdf, MS Word, ODT, RTF, or markdown document; or as a html or pdf based slide show.

Workflow

1 Open a new .Rmd file

Use the wizard that opens to pre-populate the file with a template



.Rmd structure

YAML Header

Optional section of render (e.g. pandoc) options written as key:value pairs (YAML).

- At start of file
- Between lines of ---

Text

Narration formatted with markdown, mixed with:

Code chunks

Chunks of embedded code. Each chunk:

- Begins with `r`
- ends with ``

R Markdown will run the code and append the results to the doc.

It will use the location of the .Rmd file as the **working directory**

2 Write document

by editing template

3 Knit document to create report

Use knit button or `render()` to knit

4 Preview Output

in IDE window

5 Publish (optional)

to web or server
Sync publish button to accounts at

- rpubs.com
- shinyapps.io
- RStudio Connect

Reload document

Find in document

File path to output document

6 Examine build log

in R Markdown console

7 Use output file

that is saved alongside .Rmd

render()

Use `rmarkdown::render()` to render/knit at cmd line. Important args:

input - file to render

output_format

output_options - List of render options (as in YAML)

output_file

output_dir

params - list of params to use

envir - environment to evaluate code chunks in

encoding - of input file

Interactive Documents

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

Embed a complete app into your document with `shiny::shinyAppDir()`

* Your report will be rendered as a Shiny app, which means you must choose an html output format, like `html_document`, and serve it with an active R Session.

Embed code with knitr syntax

Inline code

Insert with `r<code>`. Results appear as text without code.

Built with
`r getRVersion()` → Built with 3.2.3

Important chunk options

cache - cache results for future knits (default = FALSE)

cache.path - directory to save cached results in (default = "cache/")

child - file(s) to knit and then include (default = NULL)

collapse - collapse all output into single block (default = FALSE)

comment - prefix for each line of results (default = '##')

Options not listed above: R.options, aniopts, autodep, background, cache.comments, cache.lazy, cache.rebuild, cache.vars, dev, dev.args, dpi, engine.opts, engine.path, fig.asp, fig.env, fig.ext, fig.keep, fig.lp, fig.path, fig.pos, fig.process, fig.retina, fig.scap, fig.show, fig.showtext, fig.subcap, interval, out.extra, out.height, out.width, prompt, purl, ref.label, render, size, split, tidy.opts

Code chunks

One or more lines surrounded with `r` and ``. Place chunk options within curly braces, after r. Insert with ↗

```{r echo=TRUE}  
getRVersion()  
...  
getRVersion()  
## [1] '3.2.3'

**fig.align** - 'left', 'right', or 'center' (default = 'default')

**fig.cap** - figure caption as character string (default = NULL)

**fig.height, fig.width** - Dimensions of plots in inches

**highlight** - highlight source code (default = TRUE)

**include** - Include chunk in doc after running (default = TRUE)

### Global options

Set with `knitr::opts_chunk$set()`, e.g.

```{r include=FALSE}  
knitr::opts_chunk\$set(echo = TRUE)
...
knitr::opts_chunk\$set(echo = TRUE)

message - display code messages in document (default = TRUE)

results (default = 'markup') 'asis' - passthrough results

'hide' - do not display results
'hold' - put all results below all code

tidy - tidy code for display (default = FALSE)

warning - display code warnings in document (default = TRUE)

Parameters

Parameterize your documents to reuse with different inputs (e.g., data sets, values, etc.)

1 Add parameters

Create and set parameters in the header as sub-values of **params**

params:
n: 100
d: !r Sys.Date()

2 Call parameters

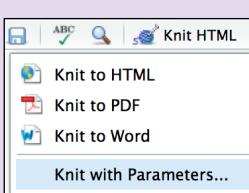
Call parameter values in code as `params$<name>`

Today's date is `r params\$d`

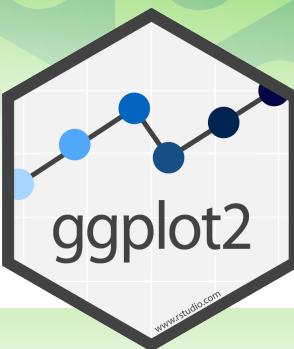
3 Set parameters

Set values with **Knit with parameters** or the `params` argument of `render()`:

`render("doc.Rmd", params = list(n = 1, d = as.Date("2015-01-01")))`



Data Visualization with ggplot2 :: CHEAT SHEET



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

Not required, sensible defaults supplied

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

aesthetic mappings **data** **geom**

qplot(x = cty, y = hwy, data = mpg, geom = "point") Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

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.

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()**
(Useful for expanding limits)
- b + geom_curve(aes(yend = lat + 1, xend = long + 1, curvature = z))** - x, xend, y, 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(group = group))**
x, y, alpha, color, fill, group, 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(f1))
- d + geom_bar()**
x, alpha, color, fill, linetype, size, weight

TWO VARIABLES

continuous x , continuous y

- e <- ggplot(mpg, aes(cty, hwy))
- e + geom_label(aes(label = cty), nudge_x = 1, nudge_y = 1, check_overlap = TRUE)** x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

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

- 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, check_overlap = TRUE)**, x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

discrete x , continuous y

- 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

discrete x , discrete y

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

- g + geom_count()**, x, y, alpha, color, fill, shape, size, stroke

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, colour, group, linetype, size, weight

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_density2d()**
x, y, alpha, colour, group, linetype, size

- h + geom_hex()**
x, y, alpha, colour, 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

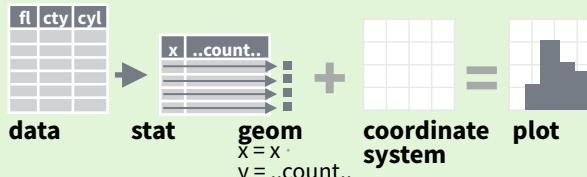
- data <- data.frame(murder = USArrests\$Murder, state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))

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

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 `..name..` syntax to map stat variables to aesthetics.



`c + stat_bin(binwidth = 1, origin = 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() + stat_function(aes(x = -3:3), n = 99, fun = dnorm, args = list(sd = 0.5)) x | ..x.., ..y..`

`e + stat_identity(na.rm = TRUE)`

`ggplot() + stat_qq(aes(sample = 1:100), dist = qt, dparam = list(df = 5)) sample, x, y | ..sample.., ..theoretical..`

`e + stat_sum(x, y, size | ..n.., ..prop..)`

`e + stat_summary(fun.data = "mean_cl_boot")`

`h + stat_summary_bin(fun.y = "mean", geom = "bar")`

`e + stat_unique()`

Scales

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_*_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 x values as date times. Use same arguments as `scale_x_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 direction of x axis

`scale_x_sqrt()` - Plot x on square root scale

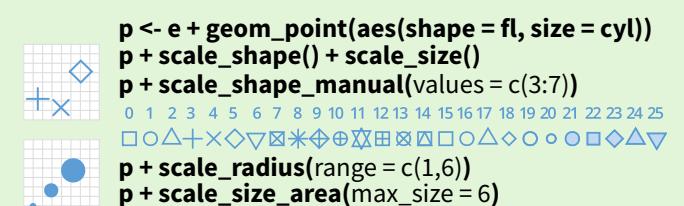
COLOR AND FILL SCALES (DISCRETE)



COLOR AND FILL SCALES (CONTINUOUS)



SHAPE AND SIZE SCALES



Coordinate Systems

`r <- d + geom_bar()`

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

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

`r + coord_flip()`
Flipped Cartesian coordinates

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

`r + coord_trans(xtrans = "sqrt")`
xtrans, ytrans, limx, limy
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

`π + coord_quickmap()`

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

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 labels

`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)))`

`αc αd αe αp αr`

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

`c d e p r`

Labels

`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", <AES> = "New <AES> legend title")`

`t + annotate(geom = "text", x = 8, y = 9, label = "A")`

`geom to place manual values for geom's aesthetics`

Legends

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

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

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

Zooming

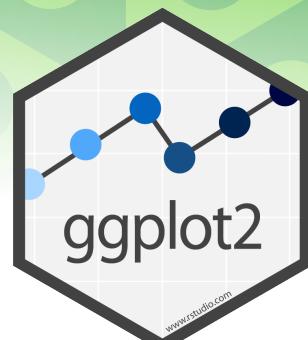
Without clipping (preferred)

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

With clipping (removes unseen data points)

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

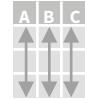
`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`



Data Transformation with dplyr :: CHEAT SHEET



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



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



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



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

Summarise Cases

These 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).



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

`count(x, ..., wt = NULL, sort = FALSE)`
Count number of rows in each group defined by the variables in ... Also **tally()**.
`count(iris, Species)`

VARIATIONS

`summarise_all()` - Apply funs to every column.

`summarise_at()` - Apply funs to specific columns.

`summarise_if()` - Apply funs to all cols of one type.

Group Cases

Use **group_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



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

`group_by(.data, ..., add = FALSE)`
Returns copy of table grouped by ...
`g_iris <- group_by(iris, Species)`

`ungroup(x, ...)`
Returns ungrouped copy of table.
`ungroup(g_iris)`

Manipulate Cases

EXTRACT CASES

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



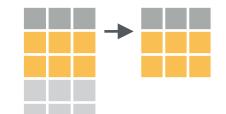
`filter(.data, ...)` Extract rows that meet logical criteria. `filter(iris, Sepal.Length > 7)`



`distinct(.data, ..., .keep_all = FALSE)` Remove rows with duplicate values.
`distinct(iris, Species)`



`sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame())` Randomly select fraction of rows.
`sample_frac(iris, 0.5, replace = TRUE)`



`slice(.data, ...)` Select rows by position.
`slice(iris, 10:15)`



`top_n(x, n, wt)` Select and order top n entries (by group if grouped data). `top_n(iris, 5, Sepal.Width)`

Manipulate Variables

EXTRACT VARIABLES

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



`pull(.data, var = -1)` Extract column values as a vector. Choose by name or index.
`pull(iris, Sepal.Length)`



`select(.data, ...)` Extract columns as a table. Also `select_if()`.
`select(iris, Sepal.Length, Species)`

Use these helpers with `select()`,
e.g. `select(iris, starts_with("Sepal"))`

`contains(match)` `num_range(prefix, range)` : e.g. `mpg:cyl`
`ends_with(match)` `one_of(...)` -, e.g. `-Species`
`matches(match)` `starts_with(match)`

MAKE NEW VARIABLES

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



`mutate(.data, ...)`
Compute new column(s).
`mutate(mtcars, gpm = 1/mpg)`

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

`mutate_all(.tbl, .funs, ...)` Apply funs to every column. Use with `funs()`. Also `mutate_if()`.
`mutate_all(faithful, funs(log(.), log2(.)))`
`mutate_if(iris, is.numeric, funs(log(.)))`

`mutate_at(.tbl, .cols, .funs, ...)` Apply funs to specific columns. Use with `funs()`, `vars()` and the helper functions for `select()`.
`mutate_at(iris, vars(-Species), funs(log(.)))`

`add_column(.data, ..., .before = NULL, .after = NULL)` Add new column(s). Also `add_count()`, `add_tally()`.
`add_column(mtcars, new = 1:32)`

`rename(.data, ...)` Rename columns.
`rename(iris, Length = Sepal.Length)`



Vector Functions

TO USE WITH MUTATE ()

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

vectorized function →

OFFSETS

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

CUMULATIVE AGGREGATES

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

RANKINGS

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

MATH

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

MISC

dplyr::case_when() - multi-case if_else()
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()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

Summary Functions

TO USE WITH SUMMARISE ()

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

summary function →

COUNTS

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

LOCATION

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

LOGICALS

mean() - Proportion of TRUE's
sum() - # of TRUE's

POSITION/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 |
|---|---|
| 1 | a |
| 2 | b |
| 3 | c |

rownames_to_column()

Move row names into col.
a <- rownames_to_column(iris, var = "C")

| A | B | C |
|---|---|---|
| 1 | a | t |
| 2 | b | u |
| 3 | c | v |

column_to_rownames()

Move col in row names.
column_to_rownames(a, var = "C")

Also **has_rownames()**, **remove_rownames()**

Combine Tables

COMBINE VARIABLES

| | | | | | | |
|---|----------------------------------|---|---|----------------------------------|---|--|
| X | A B C
a t 1
b u 2
c v 3 | + | y | A B D
a t 3
b u 2
d w 1 | = | A B C A B D
a t 1 a t 3
b u 2 b u 2
c v 3 d w 1 |
|---|----------------------------------|---|---|----------------------------------|---|--|

Use **bind_cols()** to paste tables beside each other as they are.

bind_cols(...) Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

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
a t 1 3
b u 2 2
c v 3 NA | left_join(x, y, by = NULL,
copy=FALSE, suffix=c("x","y"),...)
Join matching values from y to x. |
|---|--|

| | |
|---|---|
| A B C D
a t 1 3
b u 2 2
d w NA 1 | right_join(x, y, by = NULL, copy = FALSE,
suffix=c("x","y"),...)
Join matching values from x to y. |
|---|---|

| | |
|-------------------------------|---|
| A B C D
a t 1 3
b u 2 2 | inner_join(x, y, by = NULL, copy = FALSE,
suffix=c("x","y"),...)
Join data. Retain only rows with matches. |
|-------------------------------|---|

| | |
|---|--|
| A B C D
a t 1 3
b u 2 2
d w NA 1 | full_join(x, y, by = NULL,
copy=FALSE, suffix=c("x","y"),...)
Join data. Retain all values, all rows. |
|---|--|

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

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

Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

COMBINE CASES

| | | | | |
|---|----------------------------------|---|---|-------------------------|
| X | A B C
a t 1
b u 2
c v 3 | + | y | A B C
C v 3
d w 4 |
|---|----------------------------------|---|---|-------------------------|

Use **bind_rows()** to paste tables below each other as they are.

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)

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).

EXTRACT ROWS

| | | | | | |
|---|----------------------------------|---|---|----------------------------------|---|
| X | A B C
a t 1
b u 2
c v 3 | + | y | A B D
a t 3
b u 2
d w 1 | = |
|---|----------------------------------|---|---|----------------------------------|---|

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

semi_join(x, y, by = NULL, ...)
Return rows of x that have a match in y.
USEFUL TO SEE WHAT WILL BE JOINED.

anti_join(x, y, by = NULL, ...)
Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.

Data Import :: CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.

The front side of this sheet shows how to read text files into R with **readr**.

The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

OTHER TYPES OF DATA

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

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

Save Data

Save **x**, an R object, to **path**, a file path, as:

Comma delimited file

```
write_csv(x, path, na = "NA", append = FALSE,  
          col_names = !append)
```

File with arbitrary delimiter

```
write_delim(x, path, delim = " ", na = "NA",  
            append = FALSE, col_names = !append)
```

CSV for excel

```
write_excel_csv(x, path, na = "NA", append =  
                FALSE, col_names = !append)
```

String to file

```
write_file(x, path, append = FALSE)
```

String vector to file, one element per line

```
write_lines(x, path, na = "NA", append = FALSE)
```

Object to RDS file

```
write_rds(x, path, compress = c("none", "gz",  
                                "bz2", "xz"), ...)
```

Tab delimited files

```
write_tsv(x, path, na = "NA", append = FALSE,  
          col_names = !append)
```



Read Tabular Data

- These functions share the common arguments:

```
read_*(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),  
       quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000,  
       n_max), progress = interactive())
```

Comma Delimited Files

```
read_csv("file.csv")
```

To make file.csv run:

```
write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")
```

Semi-colon Delimited Files

```
read_csv2("file2.csv")
```

```
write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")
```

Files with Any Delimiter

```
read_delim("file.txt", delim = "|")
```

```
write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")
```

Fixed Width Files

```
read_fwf("file.fwf", col_positions = c(1, 3, 5))
```

```
write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")
```

USEFUL ARGUMENTS

Example file

```
write_file("a,b,c\n1,2,3\n4,5,NA","file.csv")  
f <- "file.csv"
```

Skip lines

```
read_csv(f, skip = 1)
```

No header

```
read_csv(f, col_names = FALSE)
```

Read in a subset

```
read_csv(f, n_max = 1)
```

Provide header

```
read_csv(f, col_names = c("x", "y", "z"))
```

Missing Values

```
read_csv(f, na = c("1", "!"))
```

Read Non-Tabular Data

Read a file into a single string

```
read_file(locale = default_locale())
```

Read each line into its own string

```
read_lines(file, skip = 0, n_max = -1L, na = character(),  
          locale = default_locale(), progress = interactive())
```

Read Apache style log files

```
read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())
```

Read a file into a raw vector

```
read_file_raw(file)
```

Read each line into a raw vector

```
read_lines_raw(file, skip = 0, n_max = -1L,  
               progress = interactive())
```



Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:  
## cols(  
##   age = col_integer(),  
##   sex = col_character(),  
##   earn = col_double()  
## )
```

age is an integer

sex is a character

1. Use **problems()** to diagnose problems.

```
x <- read_csv("file.csv"); problems(x)
```

2. Use a **col_** function to guide parsing.

- **col_guess()** - the default
 - **col_character()**
 - **col_double()**, **col_euro_double()**
 - **col_datetime(format = "")** Also **col_date(format = "")**, **col_time(format = "")**
 - **col_factor(levels, ordered = FALSE)**
 - **col_integer()**
 - **col_logical()**
 - **col_number()**, **col_numeric()**
 - **col_skip()**
- x <- read_csv("file.csv", col_types = cols(
A = col_double(),
B = col_logical(),
C = col_factor()))**

3. Else, read in as character vectors then parse with a **parse_** function.

- **parse_guess()**
 - **parse_character()**
 - **parse_datetime()** Also **parse_date()** and **parse_time()**
 - **parse_double()**
 - **parse_factor()**
 - **parse_integer()**
 - **parse_logical()**
 - **parse_number()**
- x\$A <- parse_number(x\$A)**

Tibbles - an enhanced data frame

The **tibble** package provides a new S3 class for storing tabular data, the **tibble**. Tibbles inherit the data frame class, but improve three behaviors:

- **Subsetting** - [always returns a new tibble, [[and \$ always return a vector.
- **No partial matching** - You must use full column names when subsetting
- **Display** - When you print a tibble, R provides a concise view of the data that fits on one screen

| # A tibble: 234 × 6 | manufacturer | model | displ | cyl | trans |
|---|--------------|------------|------------|-----|------------|
| 1 | audi | a4 | 1.80 | 4 | manual |
| 2 | audi | a4 | 1.80 | 4 | manual |
| 3 | audi | a4 | 2.00 | 4 | manual |
| 4 | audi | a4 | 2.00 | 4 | manual |
| 5 | audi | a4 | 2.00 | 4 | manual |
| 6 | audi | a4 | 2.00 | 4 | manual |
| 7 | audi | a4 | 3.10 | 6 | manual |
| 8 | audi | a4 quattro | 1.80 | 4 | quattro |
| 9 | audi | a4 quattro | 1.80 | 4 | quattro |
| 10 | audi | a4 quattro | 1.80 | 4 | quattro |
| ... | ... | ... | ... | ... | ... |
| 156 | 1999 | 6 | auto(l4) | 4 | auto(l4) |
| 157 | 1999 | 6 | auto(l4) | 4 | auto(l4) |
| 158 | 2008 | 8 | auto(l4) | 4 | auto(l4) |
| 159 | 2008 | 8 | auto(s4) | 4 | auto(s4) |
| 160 | 1999 | 4 | manual(m5) | 5 | manual(m5) |
| 161 | 1999 | 4 | auto(l4) | 4 | auto(l4) |
| 162 | 2008 | 4 | manual(m5) | 5 | manual(m5) |
| 163 | 2008 | 4 | manual(m5) | 5 | manual(m5) |
| 164 | 2008 | 4 | auto(l4) | 4 | auto(l4) |
| 165 | 2008 | 4 | auto(l4) | 4 | auto(l4) |
| 166 | 1999 | 4 | auto(l4) | 4 | auto(l4) |
| [reached getOption("max.print") -- omitted 68 rows] | | | | | |

tibble display

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

A large table to display

data frame display

- Control the default appearance with options:
`options(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf)`
- View full data set with **View()** or **glimpse()**
- Revert to data frame with **as.data.frame()**

CONSTRUCT A TIBBLE IN TWO WAYS

| tibble(...) | Construct by columns.
<code>tibble(x = 1:3, y = c("a", "b", "c"))</code> | Both make this tibble |
|---------------------|---|---|
| tribble(...) | Construct by rows.
<code>tribble(~x, ~y, 1, "a", 2, "b", 3, "c")</code> | <code>A tibble: 3 × 2</code>
<code>x <int> y <chr></code>
1 1 a
2 2 b
3 3 c |

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

enframe(x, name = "name", value = "value")
Convert named vector to a tibble

is_tibble(x) Test whether x is a tibble.

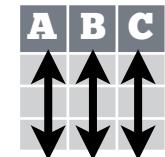


R Studio

Tidy Data with tidyverse

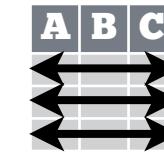
Tidy data is a way to organize tabular data. It provides a consistent data structure across packages.

A table is tidy if:



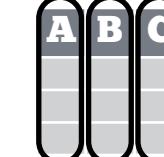
Each **variable** is in its own **column**

&



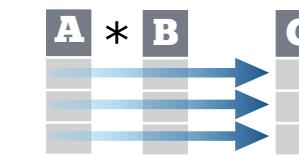
Each **observation**, or **case**, is in its own **row**

Tidy data:



Makes variables easy to access as vectors

$A * B \rightarrow C$



Preserves cases during vectorized operations

Reshape Data - change the layout of values in a table

Use **gather()** and **spread()** to reorganize the values of a table into a new layout.

gather(data, key, value, ..., na.rm = FALSE, convert = FALSE, factor_key = FALSE)

gather() moves column names into a **key** column, gathering the column values into a single **value** column.

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 |

key value

`gather(table4a, `1999`, `2000`, key = "year", value = "cases")`

table2

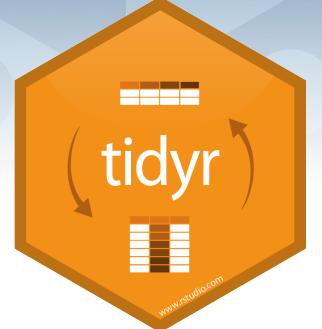
| country | year | type | count |
|---------|------|-------|-------|
| A | 1999 | cases | 0.7K |
| A | 1999 | pop | 19M |
| A | 2000 | cases | 2K |
| A | 2000 | pop | 20M |
| B | 1999 | cases | 37K |
| B | 1999 | pop | 172M |
| B | 2000 | cases | 80K |
| B | 2000 | pop | 174M |
| C | 1999 | cases | 212K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

key value

`spread(table2, type, count)`

Split Cells

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



separate(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", ...)

Separate each cell in a column to make several columns.

table3

| country | year | rate |
|---------|------|----------|
| 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 |

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

| 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 |

separate_rows(data, ..., sep = "[^[:alnum:]].+", convert = FALSE)

Separate each cell in a column to make several rows. Also **separate_rows_()**.

table3

| country | year | rate |
|---------|------|----------|
| A | 1999 | 0.7K/19M |
| A | 1999 | 19M |
| A | 2000 | 2K |
| A | 2000 | 20M |
| B | 1999 | 37K/172M |
| B | 1999 | 37K |
| B | 2000 | 80K |
| B | 2000 | 174M |
| C | 1999 | 212K/1T |
| C | 1999 | 212K |
| C | 2000 | 213K |
| C | 2000 | 1T |

`separate_rows(table3, rate)`

unite(data, col, ..., sep = "_", remove = TRUE)

Collapse cells across several columns to make a single column.

table5

| country | century | year |
|---------|---------|------|
| Afghan | 19 | 99 |
| Afghan | 20 | 0 |
| Brazil | 19 | 99 |
| Brazil | 20 | 0 |
| China | 19 | 99 |
| China | 20 | 0 |

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

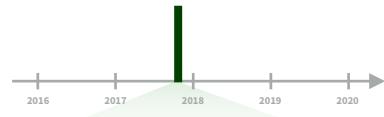
| country | year |
|---------|------|
| Afghan | 1999 |
| Afghan | 2000 |
| Brazil | 1999 |
| Brazil | 2000 |
| China | 1999 |
| China | 2000 |

Handle Missing Values

Dates and times with lubridate :: CHEAT SHEET



Date-times



2017-11-28 12:00:00

2017-11-28 12:00:00

A **date-time** is a point on the timeline, stored as the number of seconds since 1970-01-01 00:00:00 UTC

```
dt <- as_datetime(1511870400)
## "2017-11-28 12:00:00 UTC"
```

PARSE DATE-TIMES (Convert strings or numbers to date-times)

- Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data.
- Use the function below whose name replicates the order. Each accepts a wide variety of input formats.

2017-11-28T14:02:00

ymd_hms(), **ymd_hm()**, **ymd_h()**.
ymd_hms("2017-11-28T14:02:00")

2017-22-12 10:00:00

ydm_hms(), **ydm_hm()**, **ydm_h()**.
ydm_hms("2017-22-12 10:00:00")

11/28/2017 1:02:03

mdy_hms(), **mdy_hm()**, **mdy_h()**.
mdy_hms("11/28/2017 1:02:03")

1 Jan 2017 23:59:59

dmy_hms(), **dmy_hm()**, **dmy_h()**.
dmy_hms("1 Jan 2017 23:59:59")

20170131

ymd(), **ydm()**. **ymd(20170131)**

July 4th, 2000

mdy(), **myd()**. **mdy("July 4th, 2000")**

4th of July '99

dmy(), **dym()**. **dmy("4th of July '99")**

2001: Q3

yq() Q for quarter. **yq("2001: Q3")**

2:01

hms::hms() Also lubridate::hms(), **hm()** and **ms()**, which return periods.* **hms::hms(sec = 0, min = 1, hours = 2)**

2017.5

date_decimal(decimal, tz = "UTC")
date_decimal(2017.5)



now(tzone = "") Current time in tz (defaults to system tz). **now()**

today(tzone = "") Current date in a tz (defaults to system tz). **today()**

fast.strptime() Faster strftime.
fast.strptime('9/1/01', '%y/%m/%d')

parse_date_time() Easier strftime.
parse_date_time("9/1/01", "ymd")

2017-11-28

A **date** is a day stored as the number of days since 1970-01-01

```
d <- as_date(17498)
## "2017-11-28"
```

12:00:00

An hms is a **time** stored as the number of seconds since 00:00:00

```
t <- hms::as.hms(85)
## 00:01:25
```

GET AND SET COMPONENTS

Use an accessor function to get a component. Assign into an accessor function to change a component in place.

```
d ## "2017-11-28"
day(d) ## 28
day(d) <- 1
d ## "2017-11-01"
```

2018-01-31 11:59:59

date(x) Date component. **date(dt)**

2018-01-31 11:59:59

year(x) Year. **year(dt)**
isoyear(x) The ISO 8601 year.
epiyear(x) Epidemiological year.

2018-01-31 11:59:59

month(x, label, abbr) Month.
month(dt)

2018-01-31 11:59:59

day(x) Day of month. **day(dt)**
wday(x, label, abbr) Day of week.
qday(x) Day of quarter.

2018-01-31 11:59:59

hour(x) Hour. **hour(dt)**

2018-01-31 11:59:59

minute(x) Minutes. **minute(dt)**

2018-01-31 11:59:59

second(x) Seconds. **second(dt)**

2018-01-31 11:59:59

week(x) Week of the year. **week(dt)**
isoweek() ISO 8601 week.
epiweek() Epidemiological week.

2018-01-31 11:59:59

quarter(x, with_year = FALSE) Quarter. **quarter(dt)**

2018-01-31 11:59:59

semester(x, with_year = FALSE) Semester. **semester(dt)**

2018-01-31 11:59:59

am(x) Is it in the am? **am(dt)**
pm(x) Is it in the pm? **pm(dt)**

2018-01-31 11:59:59

dst(x) Is it daylight savings? **dst(dt)**

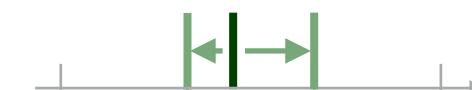
2018-01-31 11:59:59

leap_year(x) Is it a leap year?
leap_year(dt)

2018-01-31 11:59:59

update(object, ..., simple = FALSE)
update(dt, mday = 2, hour = 1)

Round Date-times



floor_date(x, unit = "second")
Round down to nearest unit.
floor_date(dt, unit = "month")

round_date(x, unit = "second")
Round to nearest unit.
round_date(dt, unit = "month")

ceiling_date(x, unit = "second", change_on_boundary = NULL)
Round up to nearest unit.
ceiling_date(dt, unit = "month")

rollback(dates, roll_to_first = FALSE, preserve_hms = TRUE)
Roll back to last day of previous month. **rollback(dt)**

Stamp Date-times

stamp() Derive a template from an example string and return a new function that will apply the template to date-times. Also **stamp_date()** and **stamp_time()**.

- Derive a template, create a function
`sf <- stamp("Created Sunday, Jan 17, 1999 3:34")`

- Apply the template to dates
`sf(ymd("2010-04-05"))`
`## [1] "Created Monday, Apr 05, 2010 00:00"`

Tip: use a date with day > 12

Time Zones

R recognizes ~600 time zones. Each encodes the time zone, Daylight Savings Time, and historical calendar variations for an area. R assigns one time zone per vector.

Use the **UTC** time zone to avoid Daylight Savings.

OlsonNames() Returns a list of valid time zone names. **OlsonNames()**

5:00 Mountain **6:00 Central**
4:00 Pacific **7:00 Eastern**

PT **MT** **CT** **ET**

7:00 Pacific **7:00 Mountain** **7:00 Central**
7:00 Eastern

with_tz(time, tzone = "") Get the same date-time in a new time zone (a new clock time).
with_tz(dt, "US/Pacific")

force_tz(time, tzone = "") Get the same clock time in a new time zone (a new date-time).
force_tz(dt, "US/Pacific")





Math with Date-times

Math with date-times relies on the **timeline**, which behaves inconsistently. Consider how the timeline behaves during:

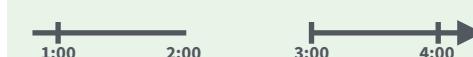
A normal day

```
nor <- ymd_hms("2018-01-01 01:30:00", tz = "US/Eastern")
```



The start of daylight savings (spring forward)

```
gap <- ymd_hms("2018-03-11 01:30:00", tz = "US/Eastern")
```



The end of daylight savings (fall back)

```
lap <- ymd_hms("2018-11-04 00:30:00", tz = "US/Eastern")
```



Leap years and leap seconds

```
leap <- ymd("2019-03-01")
```

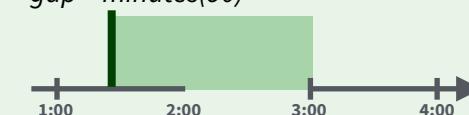


Periods track changes in clock times, which ignore time line irregularities.

```
nor + minutes(90)
```



```
gap + minutes(90)
```



```
lap + minutes(90)
```

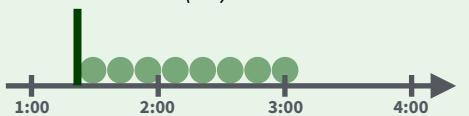


```
leap + years(1)
```



Durations track the passage of physical time, which deviates from clock time when irregularities occur.

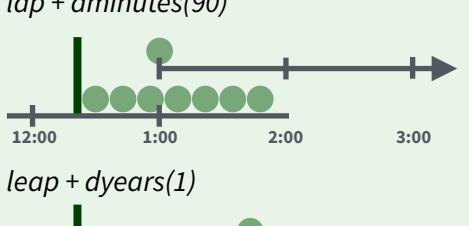
```
nor + dminutes(90)
```



```
gap + dminutes(90)
```



```
lap + dminutes(90)
```



```
leap + dyears(1)
```



Intervals represent specific intervals of the timeline, bounded by start and end date-times.

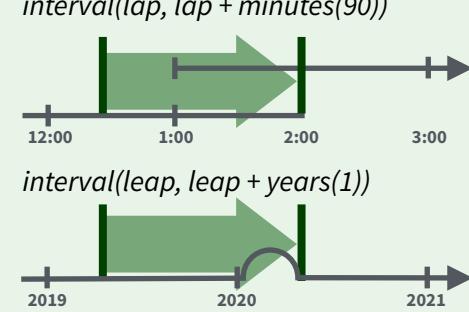
```
interval(nor, nor + minutes(90))
```



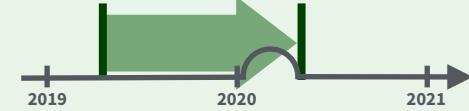
```
interval(gap, gap + minutes(90))
```



```
interval(lap, lap + minutes(90))
```



```
interval(leap, leap + years(1))
```



Not all years are 365 days due to **leap days**.

Not all minutes are 60 seconds due to **leap seconds**.

It is possible to create an imaginary date by adding **months**, e.g. February 31st

```
jan31 <- ymd(20180131)
jan31 + months(1)
## NA
```

%m+% and %m-% will roll imaginary dates to the last day of the previous month.

```
jan31 %m+% months(1)
## "2018-02-28"
```

add_with_rollback(e1, e2, roll_to_first = TRUE) will roll imaginary dates to the first day of the new month.

```
add_with_rollback(jan31, months(1),
roll_to_first = TRUE)
## "2018-03-01"
```

PERIODS

Add or subtract periods to model events that happen at specific clock times, like the NYSE opening bell.

Make a period with the name of a time unit **pluralized**, e.g.

```
p <- months(3) + days(12)
```

```
years(x = 1) x years.
```

```
months(x) x months.
```

```
weeks(x = 1) x weeks.
```

```
days(x = 1) x days.
```

```
hours(x = 1) x hours.
```

```
minutes(x = 1) x minutes.
```

```
seconds(x = 1) x seconds.
```

```
milliseconds(x = 1) x milliseconds.
```

```
microseconds(x = 1) x microseconds
```

```
nanoseconds(x = 1) x nanoseconds.
```

```
picoseconds(x = 1) x picoseconds.
```

```
period(num = NULL, units = "second", ...)
```

An automation friendly period constructor.

```
period(5, unit = "years")
```

as.period(x, unit) Coerce a timespan to a period, optionally in the specified units.

Also **is.period**(). **as.period(i)**

period_to_seconds(x) Convert a period to the "standard" number of seconds implied by the period. Also **seconds_to_period**().

```
period_to_seconds(p)
```

DURATIONS

Add or subtract durations to model physical processes, like battery life. Durations are stored as seconds, the only time unit with a consistent length.

Diftimes are a class of durations found in base R.

Make a duration with the name of a period prefixed with a **d**, e.g.

```
dd <- ddays(14)
```

```
dd
"1209600s (~2 weeks)"
```

Exact length in seconds

Equivalent in common units

```
dyears(x = 1) 31536000x seconds.
```

```
dweeks(x = 1) 604800x seconds.
```

```
ddays(x = 1) 86400x seconds.
```

```
dhours(x = 1) 3600x seconds.
```

```
dminutes(x = 1) 60x seconds.
```

```
dseconds(x = 1) x seconds.
```

```
dmilliseconds(x = 1) x × 10-3 seconds.
```

```
dmicroseconds(x = 1) x × 10-6 seconds.
```

```
dnanoseconds(x = 1) x × 10-9 seconds.
```

```
dpicoseconds(x = 1) x × 10-12 seconds.
```

```
duration(num = NULL, units = "second", ...)
```

An automation friendly duration constructor. **duration(5, unit = "years")**

as.duration(x, ...) Coerce a timespan to a duration. Also **is.duration**(), **is.difftime**(). **as.duration(i)**

make_difftime(x) Make difftime with the specified number of units.

```
make_difftime(99999)
```

INTERVALS

Divide an interval by a duration to determine its physical length, divide an interval by a period to determine its implied length in clock time.

Make an interval with **interval()** or %--%, e.g.

```
i <- interval(ymd("2017-01-01"), d)
```

```
## 2017-01-01 UTC--2017-11-28 UTC
```

```
j <- d %--% ymd("2017-12-31")
```

```
## 2017-11-28 UTC--2017-12-31 UTC
```

Start Date **End Date**

a %within% b Does interval or date-time a fall within interval b? **now()** %within% i

int_start(int) Access/set the start date-time of an interval. Also **int_end**(). **int_start(i) < now(); int_start(i)**

int_aligns(int1, int2) Do two intervals share a boundary? Also **int_overlaps**(). **int_aligns(i, j)**

int_diff(times) Make the intervals that occur between the date-times in a vector.

```
v <- c(dt, dt + 100, dt + 1000); int_diff(v)
```

int_flip(int) Reverse the direction of an interval. Also **int_standardize**(). **int_flip(i)**

int_length(int) Length in seconds. **int_length(i)**

int_shift(int, by) Shifts an interval up or down the timeline by a timespan. **int_shift(i, days(-1))**

as.interval(x, start, ...) Coerce a timespan to an interval with the start date-time. Also **is.interval**(). **as.interval(days(1), start = now())**

String manipulation with stringr :: CHEAT SHEET



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

| | |
|--|---|
| | <code>str_detect(string, pattern)</code> Detect the presence of a pattern match in a string.
<code>str_detect(fruit, "a")</code> |
| | <code>str_which(string, pattern)</code> Find the indexes of strings that contain a pattern match.
<code>str_which(fruit, "a")</code> |
| | <code>str_count(string, pattern)</code> Count the number of matches in a string.
<code>str_count(fruit, "a")</code> |
| | <code>str_locate(string, pattern)</code> Locate the positions of pattern matches in a string. Also <code>str_locate_all</code> .
<code>str_locate(fruit, "a")</code> |

Subset Strings

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

Manage Lengths

| | |
|--|--|
| | <code>str_length(string)</code> The width of strings (i.e. number of code points, which generally equals the number of characters).
<code>str_length(fruit)</code> |
| | <code>str_pad(string, width, side = c("left", "right", "both"), pad = " ")</code> Pad strings to constant width.
<code>str_pad(fruit, 17)</code> |
| | <code>str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...")</code> Truncate the width of strings, replacing content with ellipsis.
<code>str_trunc(fruit, 3)</code> |
| | <code>str_trim(string, side = c("both", "left", "right"))</code> Trim whitespace from the start and/or end of a string.
<code>str_trim(fruit)</code> |

Mutate Strings

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

Join and Split

| | |
|--|--|
| | <code>str_c(..., sep = "", collapse = NULL)</code> Join multiple strings into a single string.
<code>str_c(letters, LETTERS)</code> |
| | <code>str_c(..., sep = "", collapse = NULL)</code> Collapse a vector of strings into a single string.
<code>str_c(letters, collapse = "")</code> |
| | <code>str_dup(string, times)</code> Repeat strings times times. <code>str_dup(fruit, times = 2)</code> |
| | <code>str_split_fixed(string, pattern, n)</code> Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also <code>str_split</code> to return a list of substrings.
<code>str_split_fixed(fruit, " ", n=2)</code> |
| | <code>str_glue(..., .sep = "", .envir = parent.frame())</code> Create a string from strings and {expressions} to evaluate. <code>str_glue("Pi is {pi}")</code> |
| | <code>str_glue_data(.x, ..., .sep = "", .envir = parent.frame(), .na = "NA")</code> Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate.
<code>str_glue_data(mtcars, "{rownames(mtcars)} has {hp} hp")</code> |

Order Strings

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

Helpers

| | |
|--|---|
| | <code>str_conv(string, encoding)</code> Override the encoding of a string. <code>str_conv(fruit, "ISO-8859-1")</code> |
| | <code>str_view(string, pattern, match = NA)</code> View HTML rendering of first regex match in each string. <code>str_view(fruit, "[aeiou]")</code> |
| | <code>str_view_all(string, pattern, match = NA)</code> View HTML rendering of all regex matches. <code>str_view_all(fruit, "[aeiou]")</code> |
| | <code>str_wrap(string, width = 80, indent = 0, exdent = 0)</code> Wrap strings into nicely formatted paragraphs. <code>str_wrap(sentences, 20)</code> |

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

