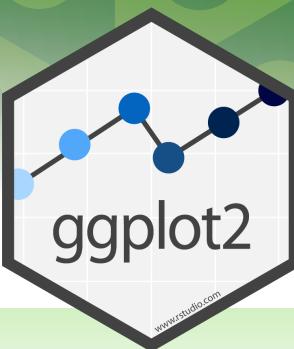


# Data Visualization with ggplot2 :: CHEAT SHEET



## 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, 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))

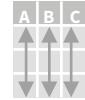
- k + geom\_map(aes(map\_id = state), map = map) + expand\_limits(x = map\$long, y = map\$lat)** - map\_id, alpha, color, fill, linetype, size



# 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)`



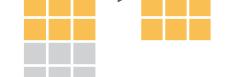
`distinct(.data, ..., .keep_all = FALSE)` Remove rows with duplicate values.  
`distinct(iris, Species)`



`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)`  
`ends_with(match)`  
`matches(match)`

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

### MAKE NEW VARIABLES

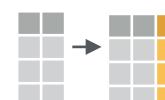
These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).



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



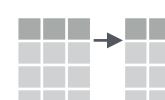
`transmute(.data, ...)`  
Compute new column(s), drop others.  
`transmute(mtcars, gpm = 1/mpg)`



`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

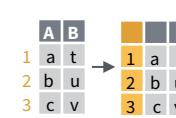
**quantile()** - nth quantile  
**min()** - minimum value  
**max()** - maximum value

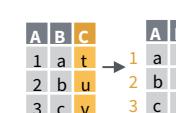
## SPREAD

**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.

 **rownames\_to\_column()**  
Move row names into col.  
a <- rownames\_to\_column(iris, var = "C")

 **column\_to\_rownames()**  
Move col in row names.  
column\_to\_rownames(a, var = "C")

Also **has\_rownames()**, **remove\_rownames()**

# Combine Tables

## COMBINE VARIABLES

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

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

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

A B C D	full_join(x, y, by = NULL,
a t 1 3 b u 2 2 d w NA 1	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	y	=
A B C a t 1 b u 2 c v 3	A B C C v 3 d w 4	

Use **bind\_rows()** to paste tables below each other as they are.

df A B C	bind_rows(..., .id = NULL)
x a t 1 x b u 2 x c v 3 z c v 3 z d w 4	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)

A B C	intersect(x, y, ...)
c v 3	Rows that appear in both x and y.

A B C	setdiff(x, y, ...)
a t 1 b u 2	Rows that appear in x but not y.

A B C	union(x, y, ...)
a t 1 b u 2 c v 3 d w 4	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	y	=
A B C a t 1 b u 2 c v 3	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.

A B C	semi_join(x, y, by = NULL, ...)
a t 1 b u 2	Return rows of x that have a match in y.

A B C	anti_join(x, y, by = NULL, ...)
c v 3	Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.

# R Markdown :: CHEAT SHEET

## What is R Markdown?

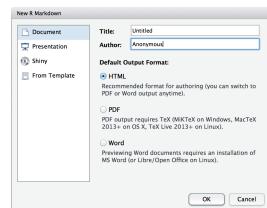


**.Rmd files** • An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.

**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 in a variety of formats, including html, pdf, MS Word, or RTF documents; html or pdf based slides, Notebooks, and more.

## Workflow



① Open a new .Rmd file at File ► New File ► R Markdown. Use the wizard that opens to pre-populate the file with a template

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 server

⑥ Examine build log in R Markdown console

⑦ Use output file that is saved along side .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

## 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

### CODE CHUNKS

One or more lines surrounded with `{{r}}` and `{{ }}`. Place chunk options within curly braces, after r. Insert with {{getRVersion()}}

```
```{r echo=TRUE}
getRVersion()
```

```

### GLOBAL OPTIONS

Set with knitr::opts\_chunk\$set(), e.g.

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

```

### 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 = "#")

**dependson** - chunk dependencies for caching (default = NULL)

**echo** - Display code in output document (default = TRUE)

**engine** - code language used in chunk (default = 'R')

**error** - Display error messages in doc (TRUE) or stop render when errors occur (FALSE) (default = FALSE)

**eval** - Run code in chunk (default = TRUE)

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

**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)

**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)

## .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

## Parameters

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

1. **Add parameters** • Create and set parameters in the header as sub-values of params

```
---
params:
  n: 100
  d: ! 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"))

Knit to HTML  
Knit to PDF  
Knit to Word  
Knit with Parameters...

## 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

```
---
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 |

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

NOTE: 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.





# Pandoc's Markdown

Write with syntax on the left to create effect on right (after render)

```
Plain text
End a line with two spaces
to start a new paragraph.
*italics* and **bold**
`verbatim` code
sub/superscript22
~~strikethrough~~
escaped: `*` \\
endash: --, emdash: ---
equation: $A = \pi * r^2$
```

```
equation block:
```

```
$$E = mc^2$$
```

```
Plain text
End a line with two spaces
to start a new paragraph.
italics and bold
`verbatim` code
sub/superscript22
strikethrough
escaped: `*` \\
endash: --, emdash: ---
equation:  $A = \pi * r^2$ 
```

```
equation block:
```

```
 $E = mc^2$ 
```

```
block quote
```

```
# Header1 {#anchor}
## Header 2 {#css_id}
```

```
### Header 3 {.css_class}
```

```
#### Header 4
```

```
##### Header 5
```

```
##### Header 6
```

```
<!--Text comment-->
```

```
\textbf{Text ignored in HTML}
```

```
<em>HTML ignored in pdfs</em>
```

```
<http://www.rstudio.com>
[link](www.rstudio.com)
Jump to [Header 1]{#anchor}
image:
```

```
![Caption](smallorb.png)
```

```
* unordered list
+ sub-item 1
+ sub-item 2
- sub-sub-item 1
```

```
* item 2
```

```
Continued (indent 4 spaces)
```

```
1. ordered list
2. item 2
  i) sub-item 1
    A. sub-sub-item 1
```

```
(@) A list whose numbering
```

```
continues after
```

```
2. an interruption
```

```
Term 1
```

```
Definition 1
```

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

- slide bullet 1
- slide bullet 2

```
(>- to have bullets appear on click)
```

```
horizontal rule/slide break:
```

```
***
```

```
A footnote [^1]
```

```
[^1]: Here is the footnote.
```

# Set render options with YAML

When you render, R Markdown

1. runs the R code, embeds results and text into .md file with knitr
2. then converts the .md file into the finished format with pandoc



Set a document's default output format in the YAML header:

```
---  
output: html_document  
---  
# Body
```

## output value

## creates

```
html_document
```

```
html
```

```
pdf_document
```

```
pdf (requires Tex)
```

```
word_document
```

```
Microsoft Word (.docx)
```

```
odt_document
```

```
OpenDocument Text
```

```
rtf_document
```

```
Rich Text Format
```

```
md_document
```

```
Markdown
```

```
github_document
```

```
Github compatible markdown
```

```
ioslides_presentation
```

```
ioslides HTML slides
```

```
slidy_presentation
```

```
slidy HTML slides
```

```
beamer_presentation
```

```
Beamer pdf slides (requires Tex)
```

Customize output with sub-options (listed to the right):

```
---  
output: html_document:  
  code_folding: hide  
  toc_float: TRUE  
---  
# Body
```

## html tabs

Use tablet css class to place sub-headers into tabs

```
# Tabset {.tabset .tabset-fade .tabset-pills}  
## Tab 1  
text 1  
## Tab 2  
text 2  
### End tabset
```



# Create a Reusable Template

1. **Create a new package** with a `inst/rmarkdown/templates` directory

2. In the directory, **Place a folder** that contains:

**template.yaml** (see below)

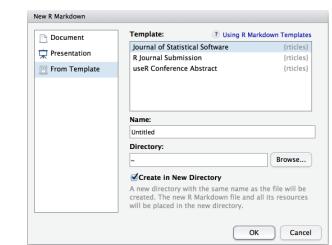
**skeleton.Rmd** (contents of the template)

any supporting files

3. **Install the package**

4. **Access template** in wizard at File ▶ New File ▶ R Markdown template.yaml

```
---  
name: My Template  
---
```



## sub-option

## description

|                       |   | html | pdf | word | odt | rtf | md | gitlab | ioslides | slidy | beamer |
|-----------------------|---|------|-----|------|-----|-----|----|--------|----------|-------|--------|
| citation_package      | The LaTeX package to process citations, natbib, biblatex or none                | X    |     |      |     |     |    |        |          |       |        |
| code_folding          | Let readers to toggle the display of R code, "none", "hide", or "show"          |      | X   |      |     |     |    |        |          |       |        |
| colortheme            | Beamer color theme to use   |      |     |      |     |     |    |        |          |       | X      |
| css                   | CSS file to use to style document   |      | X   |      |     |     |    |        | X        | X     |        |
| dev                   | Graphics device to use for figure output (e.g. "png")                           |      | X   | X    |     |     |    |        | X        | X     | X      |
| duration              | Add a countdown timer (in minutes) to footer of slides                          |      |     |      |     |     |    |        |          |       | X      |
| fig_caption           | Should figures be rendered with captions?                                       | X    | X   | X    | X   |     |    |        | X        | X     | X      |
| fig_height, fig_width | Default figure height and width (in inches) for document                        | X    | X   | X    | X   | X   | X  | X      | X        | X     | X      |
| highlight             | Syntax highlighting: "tango", "pygments", "kate", "zenburn", "textmate"         | X    | X   | X    |     |     |    |        | X        | X     |        |
| includes              | File of content to place in document (in_header, before_body, after_body)       | X    | X   | X    | X   | X   | X  | X      | X        | X     |        |
| incremental           | Should bullets appear one at a time (on presenter mouse clicks)?                |      |     |      |     |     |    |        | X        | X     | X      |
| keep_md               | Save a copy of .md file that contains knitr output                              | X    | X   | X    | X   |     |    |        | X        | X     |        |
| keep_tex              | Save a copy of .tex file that contains knitr output                             |      |     |      |     |     |    |        |          |       | X      |
| latex_engine          | Engine to render latex, "pdflatex", "xelatex", or "lualatex"                    |      |     |      |     |     |    |        |          |       | X      |
| lib_dir               | Directory of dependency files to use (Bootstrap, MathJax, etc.)                 |      | X   |      |     |     |    |        | X        | X     |        |
| mathjax               | Set to local or a URL to use a local/URL version of MathJax to render equations | X    |     |      |     |     |    |        | X        | X     |        |
| md_extensions         | Markdown extensions to add to default definition or R Markdown                  | X    | X   | X    | X   | X   | X  | X      | X        | X     |        |
| number_sections       | Add section numbering to headers  |      | X   | X    |     |     |    |        |          |       |        |
| pandoc_args           | Additional arguments to pass to Pandoc  | X    | X   | X    | X   | X   | X  | X      | X        | X     |        |
| preserve_yaml         | Preserve YAML front matter in final document?                                   |      |     |      |     |     |    |        |          |       | X      |
| reference_docx        | docx file whose styles should be copied when producing docx output              |      |     |      |     |     |    |        |          |       | X      |
| self_contained        | Embed dependencies into the doc   |      |     |      |     |     |    |        |          |       | X      |
| slide_level           | The lowest heading level that defines individual slides                         |      |     |      |     |     |    |        |          |       | X      |
| smaller               | Use the smaller font size in the presentation?                                  |      |     |      |     |     |    |        |          |       | X      |
| smart                 | Convert straight quotes to curly, dashes to em-dashes, ... to ellipses, etc.    | X    |     |      |     |     |    |        |          | X     | X      |
| template              | Pandoc template to use when rendering file quarterly_report.html                | X    | X   | X    |     |     |    |        | X        | X     |        |
| theme                 | Bootswatch or Beamer theme to use for page                                      | X    |     |      |     |     |    |        |          |       | X      |
| toc                   | Add a table of contents at start of document                                    | X    | X   | X    | X   | X   | X  | X      | X        | X     |        |
| toc_depth             | The lowest level of headings to add to table of contents                        | X    | X   | X    | X   | X   | X  | X      | X        | X     |        |
| toc_float             | Float the table of contents to the left of the main content                     | X    |     |      |     |     |    |        |          |       |        |

# Table Suggestions

Several functions format R data into tables

| Table with kable |         |
|------------------|---------|
| eruptions        | waiting |
| 3.600            | 79      |
| 1.800            | 54      |
| 3.333            | 74      |
| 2.283            | 62      |

| eruptionswaiting |      |
|------------------|------|
| 1                | 3.60 |
| 2                | 1.80 |
| 3                | 3.33 |
| 4                | 2.28 |

| Table with xtable |         |
|-------------------|---------|
| eruptions         | waiting |
| 3.600             | 79      |
| 1.800             | 54      |
| 3.333             | 74      |
| 2.283             | 62      |

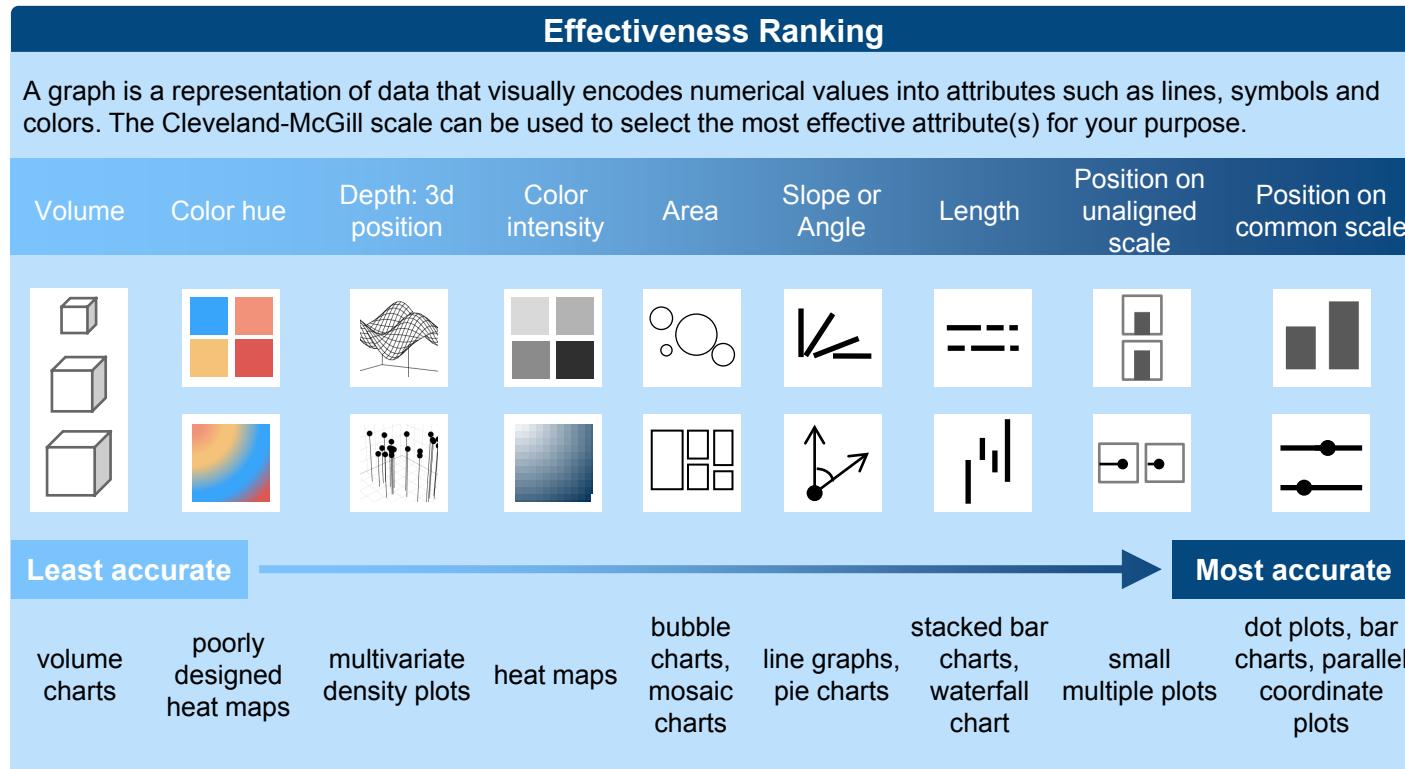
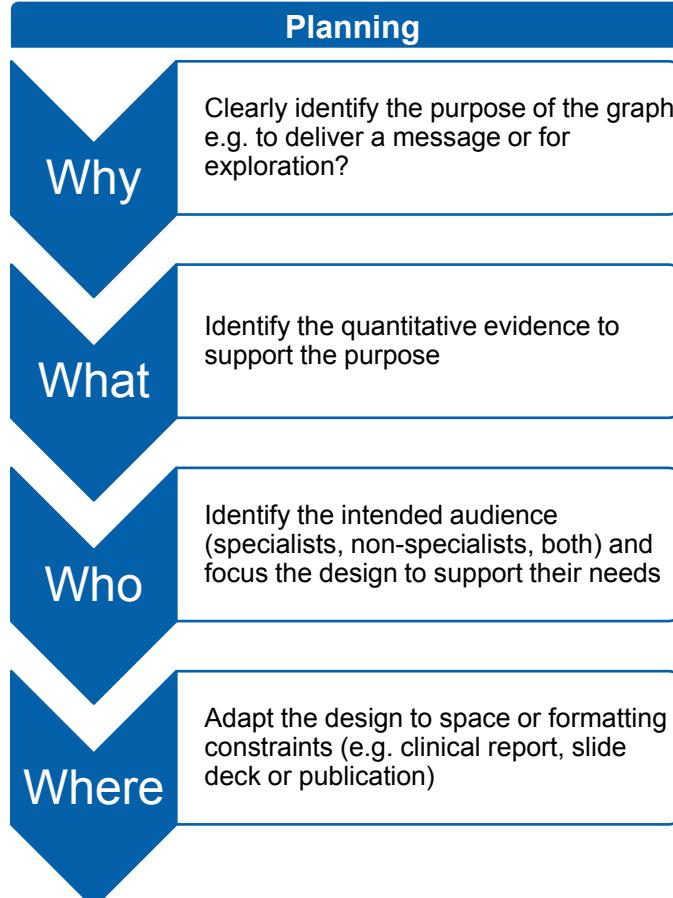
# Graphics Principles Cheat Sheet v1.1

**Communication**

Effective visualizations communicate complex statistical and quantitative information facilitating insight, understanding, and decision making.

But what is an effective graph?

This cheat sheet provides general guidance and points to consider.



**Principles of Effective Graphic Design**

**Proximity** – group related elements together

**Alignment** – elements on the same vertical or horizontal plane are perceived as having similar properties

**Simplicity** – cut anything superfluous, only include elements that add value, limit to 2-3 colors or fonts

**White space** (empty space) – use white space to minimize distraction & provide clarity

**Legibility** – sans serif fonts are easier to read, use color for emphasis instead of a new typeface

**Color** – select colors that present enough contrast to make the graph legible. Choose monochromatic color schemes to prevent clashing. Use dark colors and accent colors to emphasize important information

**Visual Hierarchy** – use color, font, image size, typeface, alignment & placement to create a viewing order

**Focal Points** – primary area of interest that immediately attracts the eye, emphasize the most important concept and make it your focal point. Use contrasting colors to draw attention

**Repetition** – repeating elements can be visually appealing, repeated shapes, labels, colors

**Familiarity** – using familiar styles, icons, navigation structure makes viewers feel confident

**Consistency** – be consistent with heading sizes, font choices, color scheme, and spacing. Use images with similar styles

**Selecting the right base graph**

Consider if a standard graph can be used by identifying suitable designs based on the:

- (i) **purpose** (i.e. message to be conveyed or question to answer)
- (ii) **data** (i.e. variables to display)

**Example plots categorized by purpose**

| Deviation | Correlation | Ranking | Distribution | Evolution | Part-to-whole | Magnitude   |
|-----------|-------------|---------|--------------|-----------|---------------|-------------|
|           |             |         |              |           |               |             |
| Waterfall | Heat map    | Dotplot | Histogram    | Line plot | Tree map      | Forest plot |

**Facilitating Comparisons**

**Proximity improves association**

Place labels next to data instead of using legends

Group together elements to be compared directly

**Ease visual inspection**

Order values to help compare across many categories

Judgments are easier to make on a common vertical scale

**Reduce mental arithmetic**

Plot the final comparison e.g. mean difference not two means  
Exception: if comparator is of interest in itself

Use reference lines and other visual anchors.

**Color for emphasis or distinction**

Restrained use of color is highly effective in organizing a narrative and calling attention to certain elements.

Think carefully before introducing additional color. Do you really need it?

Do not use color to differentiate between categories of the same variable

Use colors or shades to represent meaningful differences such as positive/negative values, treatments or doses

Be consistent, use the same color to mean the same thing in a series of graphs (e.g. treatment, dose)

Use a bold, saturated or contrasting color to emphasize important details.

Emphasize the data by minimizing unnecessary ink, e.g. soften gridlines with a light color

Utilize existing resources for selection of appropriate palettes such as Color brewer or Munsell

## Implementation Considerations

Plot cause on the x-axis and effect on the y-axis. Use this standard convention in order to avoid misinterpretation.

Aspect ratio can influence interpretation. Aim for a 45 degree angle of change to avoid over-interpretation of slope.

Use position for comparisons rather than length (i.e. dots instead of bars), especially for non-linear scales (e.g. log scale or % change).

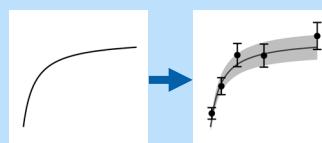
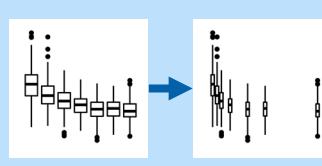
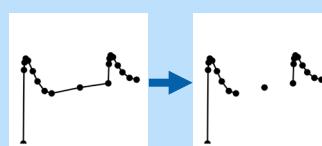
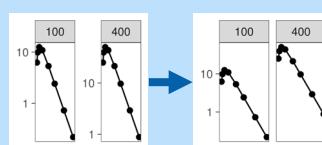
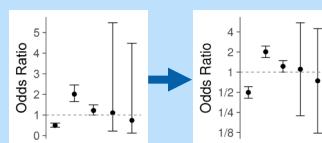
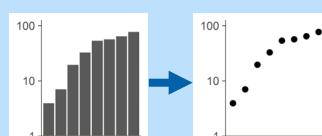
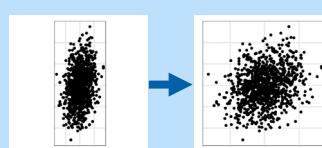
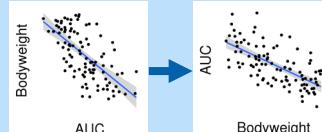
Do not plot log-normally distributed variables on a linear scale (e.g. hazard ratio, AUC, CL)

When displaying data measured on the same scale, also plot them on the same scale for easy comparison.

Connected data imply continuity. Do not connect data across a disconnected or uneven time scale.

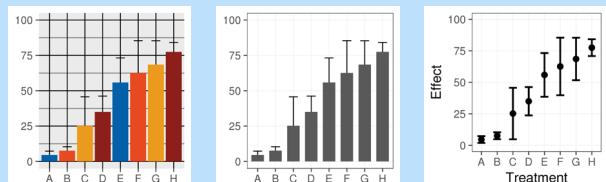
Visits displayed close together are perceived to be closer in time. Space the visits proportional to the time between each in order to avoid confusion. Exception: baseline or pre-dose

Plot data and inferences to support stories about models.



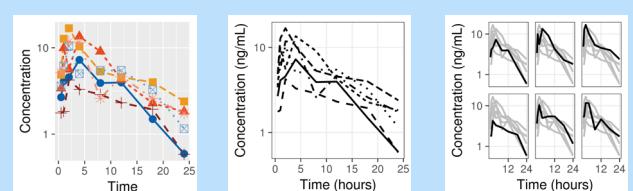
## Putting it all together – Remove the clutter & emphasize the message

Creating a graph is an iterative process: produce, review and refine.



Colors, backgrounds, and borders can be removed and gridlines reduced.

It is easier to see differences in position over a difference in length, i.e. a dot over a bar.



Using too many colors can be distracting. Use white background and try using other methods to distinguish different curves.

One solution could be repeating the data in different panels, highlighting individual curves in a darker color.

## Good graph checklist

### Clear Communication

- Is the message of the graph as clear as possible?
- Is it easy for someone unfamiliar with the data to interpret the graph?
- Are the patterns/relationships easily identified?
- Is the graph tailored to its primary purpose and audience?
- Is the correct graph type used?

### Facilitating Comparisons

- Are elements to be compared grouped together?
- Are labels placed next to data instead of in legends?
- Have categories been ordered for easy comparison?
- Can the plot be read without doing mental calculations?
- Are the estimates of interest plotted (e.g. mean differences with confidence intervals)?

### Color for emphasis or distinction

- Are graphical elements displayed in a dark color on a light background?
- Are grid lines drawn with a thin line and a light color such as grey?
- Are colors used sparingly (e.g. max 3)?
- Do all elements in the graph have a purpose (e.g. colors, textures, grid lines)?
- Are the same colors used to mean the same thing in a series of graphs?

### Implementation Considerations

- Are multiple panels plotted on the same scale?
- Are lognormally distributed variables plotted on a log scale?
- Are common baselines used wherever possible?
- Does the orientation of the axes aid interpretation?
- Does the aspect ratio allow the reader to see variations in the data?
- Are data across a disconnected time scale kept disconnected?
- Are data spaced proportionally to the actual time interval (instead of according to visit number)?
- Are data and inferences plotted to support stories about models?
- Are number of patients by group reported if this adds context?

### Legibility and Clarity

- Can all graphical elements be seen?
- Does the graph have a clear title, axis labels, annotations and data units?
- Can the font be read without eye strain or effort?
- Are sans-serif fonts used?
- Do text sizes have correct hierarchy (big to small, main text to subtext)?
- Are the elements of the graph clearly labeled (e.g. points, error bars, lines, shaded regions)?
- Are labels oriented horizontally where possible?

## Resources

### Books:

- E. R. Tufte, The visual display of quantitative information, Connecticut, Graphics Press, 2001.
- Cleveland, W.S. and McGill, Robert, Graphical perception: theory, experimentation and application to the development of graphical methods, JASA, Vol. 79, No. 387, pp. 531 – 554, 1984.
- S. Few, Show Me The Numbers - Designing Tables and Graphs to Enlighten (2nd. Edition), Burlingame, CA: Analytics Press, 2012.
- D. M. Wong, The Wall Street Journal Guide to Information Graphics: The Dos and Don'ts of Presenting Data, Facts, and Figures. December 16, 2013.
- J. Doumont, Trees, maps, and theorems: Effective communication for rational minds. PRINCIPIAE.
- N. B. Robbins, Creating More Effective Graphs. Chart House.

- Online resources:
- <https://www.perceptualedge.com/> (S. Few)
  - <https://www.edwardtufte.com/tufte/> (E. Tufte)
  - <http://flowingdata.com/> (N. Yau)
  - <http://www.principiae.be/> (J. Doumont)
  - <http://andrewgelman.com/> (A. Gelman)
  - <http://www.thefunctionalart.com/> (A. Cairo)
  - <http://www.nbr-graphs.com/> (N. Robbins)

## Authors

Alison Margolskee, Mark Baillie, Baldur Magnusson, Julie Jones, Marc Vandemeulebroecke



# 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())
```

| a,b,c | 1,2,3 | 4,5,NA |
|-------|-------|--------|
| 1     | 2     | 3      |
| 4     | 5     | NA     |

| a;b;c | 1;2;3 | 4;5;NA |
|-------|-------|--------|
| 1     | 2     | 3      |
| 4     | 5     | NA     |

| a b c | 1 2 3 | 4 5 NA |
|-------|-------|--------|
| 1     | 2     | 3      |
| 4     | 5     | NA     |

| a b c | 1 2 3 | 4 5 NA |
|-------|-------|--------|
| 1     | 2     | 3      |
| 4     | 5     | NA     |

### 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")
```

### Tab Delimited Files

```
read_tsv("file.tsv") Also read_table().
```

```
write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")
```

## USEFUL ARGUMENTS

| a,b,c | 1,2,3 | 4,5,NA |
|-------|-------|--------|
| 1     | 2     | 3      |
| 4     | 5     | NA     |

### Example file

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

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

### Skip lines

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

| a,b,c | 1,2,3 | 4,5,NA |
|-------|-------|--------|
| 1     | 2     | 3      |
| 4     | 5     | NA     |

### No header

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

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

### Read in a subset

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

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

### Provide header

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

| A  | B | C  |
|----|---|----|
| NA | 2 | 3  |
| 4  | 5 | NA |

### Missing Values

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

## Read Non-Tabular Data

### Read a file into a single string

```
read_file(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) | 4   | manual(m5) |
| 161   | 1999         | 4          | auto(l4)   | 4   | auto(l4)   |
| 162   | 2008         | 4          | manual(m5) | 4   | manual(m5) |
| 163   | 2008         | 4          | manual(m5) | 4   | 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

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

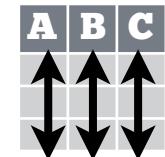


R Studio

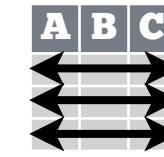
# Tidy Data with tidyverse

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

A table is tidy if:



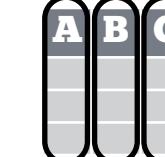
&



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

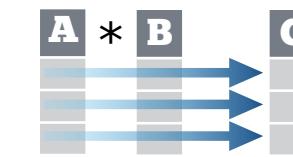
Each **observation**, or **case**, is in its own **row**

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

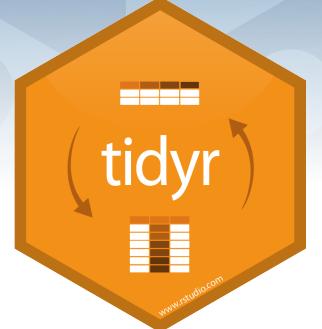
| country | year | type  | count |
|---------|------|-------|-------|
| A       | 1999 | cases | 0.7K  |
| A       | 1999 | pop   | 19M   |
| A       | 2000 | cases | 2K    |
| A       | 2000 | pop   | 20M   |
| B       | 1999 | cases | 37K   |
| B       | 1999 | pop   | 172M  |
| B       | 2000 | cases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

separate\_rows(table3, rate)

| 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 |
| C       | 1999 | 212K |
| C       | 1999 | 1T   |
| C       | 2000 | 213K |
| C       | 2000 | 1T   |

**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

**drop\_na(data, ...)**

Drop rows containing NA's in ... columns.

| x1 | x2 |
|----|----|
| A  | 1  |
| B  | NA |
| C  | NA |
| D  | 3  |
| E  | NA |

`drop_na(x, x2)`

**fill(data, ..., .direction = c("down", "up"))**

Fill in NA's in ... columns with most recent non-NA values.

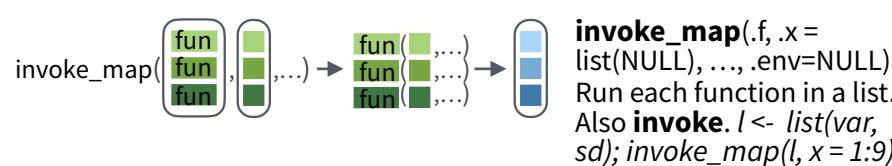
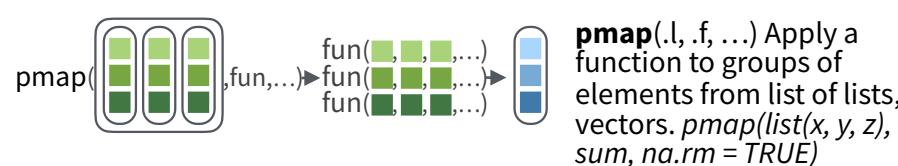
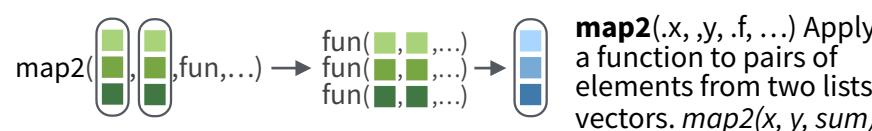
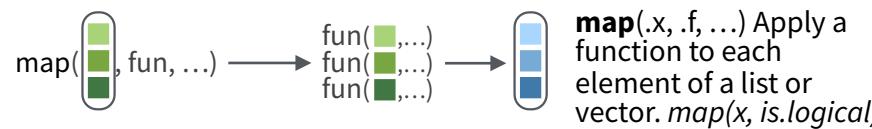
| x1 | x2 |
|----|----|
| A  | 1  |
| B  | NA |
| C  | NA |
| D  | 3  |

# Apply functions with purrr :: CHEAT SHEET



## Apply Functions

Map functions apply a function iteratively to each element of a list or vector.



**lmap(.x, .f, ...)** Apply function to each list-element of a list or vector.  
**imap(.x, .f, ...)** Apply .f to each element of a list or vector and its index.

### OUTPUT

**map()**, **map2()**, **pmap()**, **imap** and **invoke\_map** each return a list. Use a suffixed version to return the results as a specific type of flat vector, e.g. **map2\_chr**, **pmap\_lgl**, etc.

Use **walk**, **walk2**, and **pwalk** to trigger side effects. Each return its input invisibly.

### SHORTCUTS - within a purrr function:

"**name**" becomes `function(x) x[["name"]]`, e.g. `map(l, "a")` extracts `a` from each element of `l`

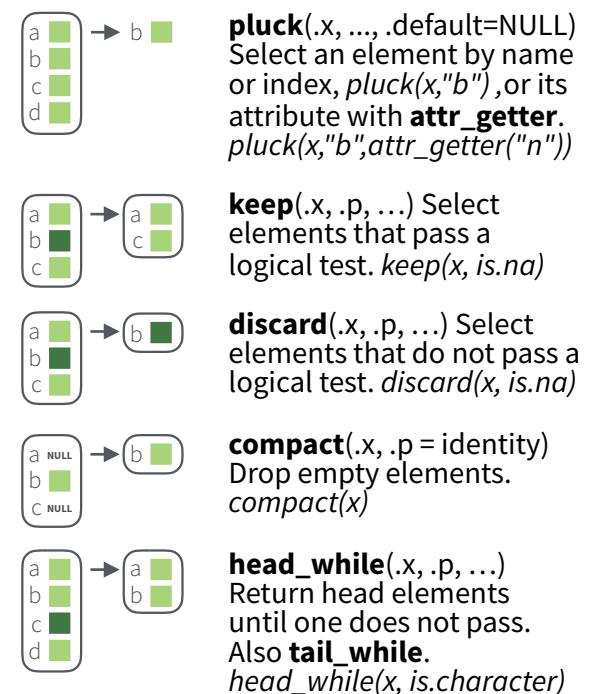
`~.x` becomes **function(x) x**, e.g. `map(l, ~2+x)` becomes `map(l, function(x) 2+x)`

### function returns

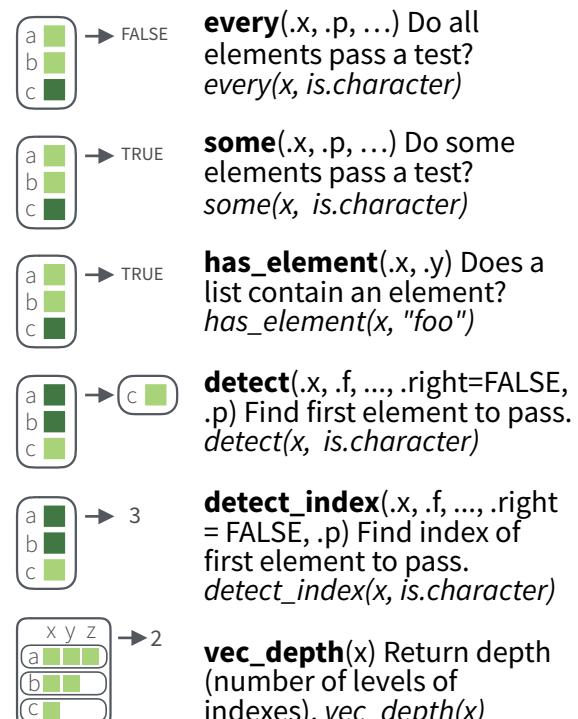
|                 |  |
|-----------------|--|
| <b>function</b> | <b>returns</b>                                     |
| <b>map</b>      | list   |
| <b>map_chr</b>  | character vector                                   |
| <b>map_dbl</b>  | double (numeric) vector                            |
| <b>map_dfc</b>  | data frame (column bind)                           |
| <b>map_dfr</b>  | data frame (row bind)                              |
| <b>map_int</b>  | integer vector                                     |
| <b>map_lgl</b>  | logical vector                                     |
| <b>walk</b>     | triggers side effects, returns the input invisibly |

## Work with Lists

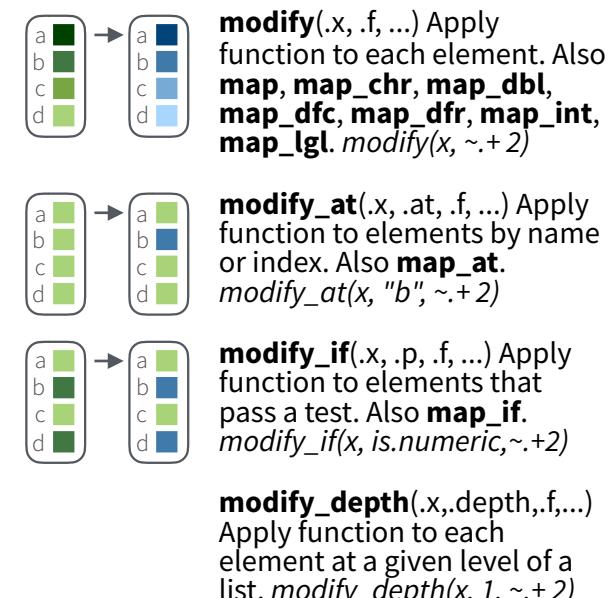
### FILTER LISTS



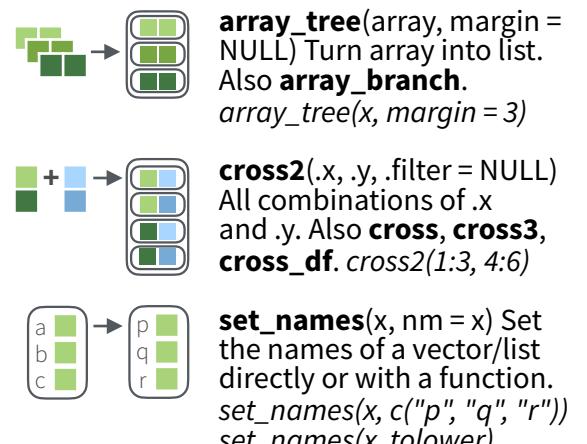
### SUMMARISE LISTS



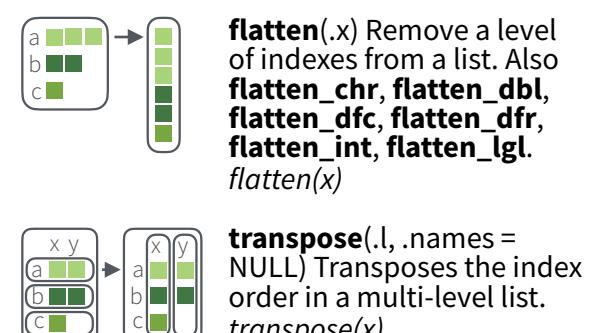
### TRANSFORM LISTS



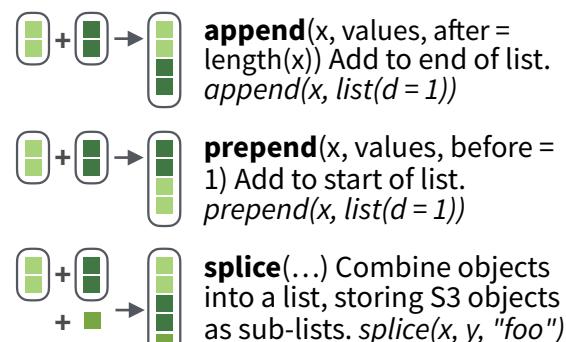
### WORK WITH LISTS



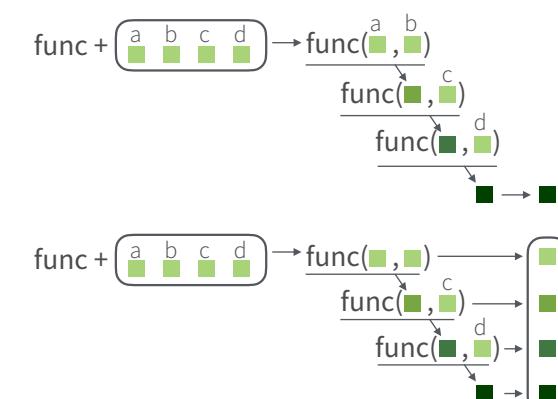
### RESHAPE LISTS



### JOIN (TO) LISTS



## Reduce Lists



**reduce(.x, .f, ..., .init)** Apply function recursively to each element of a list or vector. Also **reduce\_right**, **reduce2**, **reduce2\_right**. `reduce(x, sum)`

**accumulate(.x, .f, ..., .init)** Reduce, but also return intermediate results. Also **accumulate\_right**. `accumulate(x, sum)`

**compose()** Compose multiple functions.

**lift()** Change the type of input a function takes. Also **lift\_dl**, **lift\_lv**, **lift\_id**, **lift\_iv**, **lift\_vd**, **lift\_vl**.

**rerun()** Rerun expression n times.

**negate()** Negate a predicate function (a pipe friendly !)

**partial()** Create a version of a function that has some args preset to values.

**safely()** Modify func to return list of results whenever an error occurs (instead of error).

**quietly()** Modify function to return list of results, output, messages, warnings.

**possibly()** Modify function to return default value whenever an error occurs (instead of error).



# Nested Data

A **nested data frame** stores individual tables within the cells of a larger, organizing table.

| "cell" contents |         |         |         |
|-----------------|---------|---------|---------|
| Sepal.L         | Sepal.W | Petal.L | Petal.W |
| 5.1             | 3.5     | 1.4     | 0.2     |
| 4.9             | 3.0     | 1.4     | 0.2     |
| 4.7             | 3.2     | 1.3     | 0.2     |
| 4.6             | 3.1     | 1.5     | 0.2     |
| 5.0             | 3.6     | 1.4     | 0.2     |

n\_iris\$data[[1]]

| nested data frame |            | Species    | data              |
|-------------------|------------|------------|-------------------|
| setosa            | setosa     | setosa     | <tibble [50 x 4]> |
| versicolor        | versicolor | versicolor | <tibble [50 x 4]> |
| virginica         | virginica  | virginica  | <tibble [50 x 4]> |

n\_iris

| Sepal.L | Sepal.W | Petal.L | Petal.W |
|---------|---------|---------|---------|
| 7.0     | 3.2     | 4.7     | 1.4     |
| 6.4     | 3.2     | 4.5     | 1.5     |
| 6.9     | 3.1     | 4.9     | 1.5     |
| 5.5     | 2.3     | 4.0     | 1.3     |
| 6.5     | 2.8     | 4.6     | 1.5     |

n\_iris\$data[[2]]

| Sepal.L | Sepal.W | Petal.L | Petal.W |
|---------|---------|---------|---------|
| 6.3     | 3.3     | 6.0     | 2.5     |
| 5.8     | 2.7     | 5.1     | 1.9     |
| 7.1     | 3.0     | 5.9     | 2.1     |
| 6.3     | 2.9     | 5.6     | 1.8     |
| 6.5     | 3.0     | 5.8     | 2.2     |

n\_iris\$data[[3]]

Use a nested data frame to:

- preserve relationships between observations and subsets of data
- manipulate many sub-tables at once with the **purrr** functions **map()**, **map2()**, or **pmap()**.

Use a two step process to create a nested data frame:

1. Group the data frame into groups with **dplyr::group\_by()**
2. Use **nest()** to create a nested data frame with one row per group

|                                 |                                 |
|---------------------------------|---------------------------------|
| Species   S.L   S.W   P.L   P.W | Species   S.L   S.W   P.L   P.W |
| setosa 5.1 3.5 1.4 0.2          | setosa 5.1 3.5 1.4 0.2          |
| setosa 4.9 3.0 1.4 0.2          | setosa 4.9 3.0 1.4 0.2          |
| setosa 4.7 3.2 1.3 0.2          | setosa 4.7 3.2 1.3 0.2          |
| setosa 4.6 3.1 1.5 0.2          | setosa 4.6 3.1 1.5 0.2          |
| setosa 5.0 3.6 1.4 0.2          | setosa 5.0 3.6 1.4 0.2          |
| versi 7.0 3.2 4.7 1.4           | versi 7.0 3.2 4.7 1.4           |
| versi 6.4 3.2 4.5 1.5           | versi 6.4 3.2 4.5 1.5           |
| versi 6.9 3.1 4.9 1.5           | versi 6.9 3.1 4.9 1.5           |
| versi 5.5 2.3 4.0 1.3           | versi 5.5 2.3 4.0 1.3           |
| versi 6.5 2.8 4.6 1.5           | versi 6.5 2.8 4.6 1.5           |
| virgini 6.3 3.3 6.0 2.5         | virgini 6.3 3.3 6.0 2.5         |
| virgini 5.8 2.7 5.1 1.9         | virgini 5.8 2.7 5.1 1.9         |
| virgini 7.1 3.0 5.9 2.1         | virgini 7.1 3.0 5.9 2.1         |
| virgini 6.3 2.9 5.6 1.8         | virgini 6.3 2.9 5.6 1.8         |
| virgini 6.5 3.0 5.8 2.2         | virgini 6.5 3.0 5.8 2.2         |

n\_iris <- iris %>% group\_by(Species) %>% nest()

**tidy::nest(data, ..., .key = data)**

For grouped data, moves groups into cells as data frames.

Unnest a nested data frame with **unnest()**:

n\_iris %>% unnest()

**tidy::unnest(data, ..., .drop = NA, .id=NULL, .sep=NULL)**

Unnests a nested data frame.

# List Column Workflow

Nested data frames use a **list column**, a list that is stored as a column vector of a data frame. A typical **workflow** for list columns:

## 1 Make a list column

| Species | S.L | S.W | P.L | P.W |
|---------|-----|-----|-----|-----|
| setosa  | 5.1 | 3.5 | 1.4 | 0.2 |
| setosa  | 4.9 | 3.0 | 1.4 | 0.2 |
| setosa  | 4.7 | 3.2 | 1.3 | 0.2 |
| setosa  | 4.6 | 3.1 | 1.5 | 0.2 |
| setosa  | 5.0 | 3.6 | 1.4 | 0.2 |
| versi   | 7.0 | 3.2 | 4.7 | 1.4 |
| versi   | 6.4 | 3.2 | 4.5 | 1.5 |
| versi   | 6.9 | 3.1 | 4.9 | 1.5 |
| versi   | 5.5 | 2.3 | 4.0 | 1.3 |
| versi   | 6.5 | 2.8 | 4.6 | 1.5 |
| virgini | 6.3 | 3.3 | 6.0 | 2.5 |
| virgini | 5.8 | 2.7 | 5.1 | 1.9 |
| virgini | 7.1 | 3.0 | 5.9 | 2.1 |
| virgini | 6.3 | 2.9 | 5.6 | 1.8 |
| virgini | 6.5 | 3.0 | 5.8 | 2.2 |

```
n_iris <- iris %>%  
  group_by(Species) %>%  
  nest()
```

## 2 Work with list columns

| Species | data            | model    |
|---------|-----------------|----------|
| setosa  | <tibble [50x4]> | <S3: lm> |
| versi   | <tibble [50x4]> | <S3: lm> |
| virgini | <tibble [50x4]> | <S3: lm> |

```
mod_fun <- function(df)  
  lm(Sepal.Length ~ ., data = df)
```

```
m_iris <- n_iris %>%  
  mutate(model = map(data, mod_fun))
```

## 3 Simplify the list column

| Species | beta |
|---------|------|
| setos   | 2.35 |
| versi   | 1.89 |
| virgini | 0.69 |

```
b_fun <- function(mod)  
  coefficients(mod)[[1]]
```

```
m_iris %>% transmute(Species,  
  beta = map_dbl(model, b_fun))
```

### 1. MAKE A LIST COLUMN

You can create list columns with functions in the **tibble** and **dplyr** packages, as well as **tidyR**'s **nest()**

**tibble::tribble(...)**

Makes list column when needed

| max | seq                  |
|-----|----------------------|
| 3   | <code>int [3]</code> |
| 4   | <code>int [4]</code> |
| 5   | <code>int [5]</code> |

**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 tibble with list cols

`enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')`

**dplyr::mutate(.data, ...)** Also **transmute()**

Returns list col when result returns list.

`mtcars %>% mutate(seq = map(cyl, seq))`

### 2. WORK WITH LIST COLUMNS

Use the purrr functions **map()**, **map2()**, and **pmap()** to apply a function that returns a result element-wise to the cells of a list column. **walk()**, **walk2()**, and **pwalk()** work the same way, but return a side effect.

**purrr::map(.x, .f, ...)**

Apply .f element-wise to .x as .f(x)

`n_iris %>% mutate(n = map(data, dim))`

**purrr::map2(.x, .y, .f, ...)**

Apply .f element-wise to .x and .y as .f(x, .y)

`m_iris %>% mutate(n = map2(data, model, list))`

**purrr::pmap(.l, .f, ...)**

Apply .f element-wise to vectors saved in .l

`m_iris %>%  
 mutate(n = pmap(list(data, model, data), list))`

**map(data, fun, ...)**

`fun(<tibble [50x4]>, ...)`

`fun(<tibble [50x4]>, ...)`

`fun(<tibble [50x4]>, ...)`

**map2(data, model, fun, ...)**

`fun(<tibble [50x4]>, ...)`

`fun(<tibble [50x4]>, ...)`

`fun(<tibble [50x4]>, ...)`

# 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.<br><code>str_detect(fruit, "a")</code>                                     |
|  | <code>str_which(string, pattern)</code> Find the indexes of strings that contain a pattern match.<br><code>str_which(fruit, "a")</code>                                 |
|  | <code>str_count(string, pattern)</code> Count the number of matches in a string.<br><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> .<br><code>str_locate(fruit, "a")</code> |

## Subset Strings

|  |  |
|--|--|
|  | <code>str_sub(string, start = 1L, end = -1L)</code> Extract substrings from a character vector.<br><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.<br><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.<br><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> .<br><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).<br><code>str_length(fruit)</code>                            |
|  | <code>str_pad(string, width, side = c("left", "right", "both"), pad = " ")</code> Pad strings to constant width.<br><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.<br><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.<br><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.<br><code>str_sub(fruit, 1, 3) &lt;- "str"</code> |
|  | <code>str_replace(string, pattern, replacement)</code> Replace the first matched pattern in each string.<br><code>str_replace(fruit, "a", "-")</code>  |
|  | <code>str_replace_all(string, pattern, replacement)</code> Replace all matched patterns in each string.<br><code>str_replace_all(fruit, "a", "-")</code>                                       |
|  | <code>str_to_lower(string, locale = "en")<sup>1</sup></code> Convert strings to lower case.<br><code>str_to_lower(sentences)</code>  |
|  | <code>str_to_upper(string, locale = "en")<sup>1</sup></code> Convert strings to upper case.<br><code>str_to_upper(sentences)</code>  |
|  | <code>str_to_title(string, locale = "en")<sup>1</sup></code> Convert strings to title case.<br><code>str_to_title(sentences)</code>  |

## Join and Split

|  |  |
|--|--|
|  | <code>str_c(..., sep = "", collapse = NULL)</code> Join multiple strings into a single string.<br><code>str_c(letters, LETTERS)</code>   |
|  | <code>str_c(..., sep = "", collapse = NULL)</code> Collapse a vector of strings into a single string.<br><code>str_c(letters, collapse = "")</code>  |
|  | <code>str_dup(string, times)</code> Repeat strings times times.<br><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.<br><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.<br><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.<br><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, ...)<sup>1</sup></code> Return the vector of indexes that sorts a character vector.<br><code>x[str_order(x)]</code> |
|  | <code>str_sort(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)<sup>1</sup></code> Sort a character vector.<br><code>str_sort(x)</code>   |

## Helpers

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

<sup>1</sup> See [bit.ly/ISO639-1](http://bit.ly/ISO639-1) for a complete list of locales.



# Scoped verbs

## Suffixes

| suffix | use when   |
|--------|--|
| _all   | you want to apply the verb to all columns                        |
| _at    | you want to apply the verb to specified columns                  |
| _if    | you want to apply the verb to all the columns with some property |

## Examples

`mutate()`, `summarize()`, `select()`, and `rename()`

## Named functions

| Verb                       | Example                          | Example explanation             |
|----------------------------|----------------------------------|---------------------------------|
| <code>summarize_all</code> | <code>summarize_all(mean)</code> | finds the mean of all variables |

| <b>Verb</b>  | <b>Example</b>                         | <b>Example explanation</b>                       |
|--------------|--|--|
| summarize_at | summarize_at(vars(x, y), mean)         | finds the mean of variables x and y              |
| summarize_if | summarize_if(is.double, mean)          | finds the mean of all double variables           |
| mutate_all   | mutate_all(as.character)               | converts all variables to characters             |
| mutate_at    | mutate_at(vars(x, y),<br>as.character) | converts variables x and y to characters         |
| mutate_if    | mutate_if(is.factor, as.character)     | converts all factor variables to characters      |
| rename_all   | rename_all(str_to_lower)               | changes all column names to lowercase            |
| rename_at    | rename_at(vars(X, Y),<br>str_to_lower) | changes the names of columns X and Y to x and y  |
| rename_if    | rename_if(is.double,<br>str_to_lower)  | changes the names of double columns to lowercase |

| Verb       | Example                             | Example explanation   |
|------------|-------------------------------------|---|
| select_all | select_all(str_to_lower)            | selects all columns and changes their names to lowercase (better to use rename_all()) |
| select_at  | select_at(vars(X, Y), str_to_lower) | selects just columns X and Y and changes their names to x and y                       |
| select_if  | select_if(is.double, str_to_lower)  | selects just double columns and changes their names to lowercase                      |

## Extra arguments

| verb          | example                                       | example_explanation   |
|---------------|---|---|
| summarize_if  | summarize_if(is.double, mean, na.rm = TRUE)   | finds the mean, excluding NAs, of all double variables  |
| summarize_all | summarize_all(mean, trim = 0.1, na.rm = TRUE) | finds the mean of all variables, excluding NAs. Removes the bottom and top 10% of values of each variable before computing mean |

## Anonymous functions

| <b>verb</b>   | <b>example</b>                 | <b>example_explanation</b>                                 |
|---------------|--------------------------------|--|
| summarize_all | summarize_all(~ sum(is.na(.))) | determines the number of NAs in each column                |
| select_if     | select_if(~ n_distinct(.) > 1) | selects only the columns with more than one distinct value |

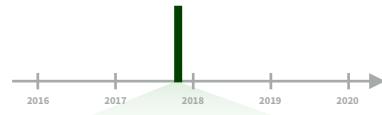
#### **filter()**

| <b>verb</b> | <b>example</b>                             | <b>example_explanation</b>                              |
|-------------|--|---|
| filter_all  | filter_all(all_vars(!is.na(.)))            | finds rows without any NAs                              |
| filter_all  | filter_all(any_vars(!is.na(.)))            | finds rows with at least one non-NA value               |
| filter_at   | filter_at(vars(x, y), all_vars(!is.na(.))) | finds rows where both x and y are non-NA                |
| filter_at   | filter_at(vars(x, y), any_vars(!is.na(.))) | finds rows where at least one of x and y is non-NA      |
| filter_if   | filter_if(is.double, all_vars(!is.na(.)))  | finds rows where all double variables are non-NA        |
| filter_if   | filter_if(is.double, any_vars(!is.na(.)))  | finds rows where at least one double variable is non-NA |

# 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)

1. Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data.
2. 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()**.

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

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

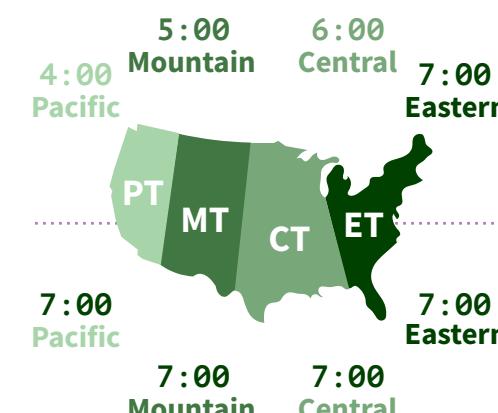
**Tip:** use a date with day > 12

## 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()**



**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:

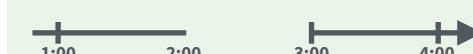
A normal day

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



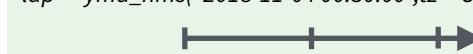
The start of daylight savings (spring forward)

```
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")
```



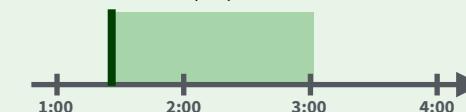
Leap years and leap seconds

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



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

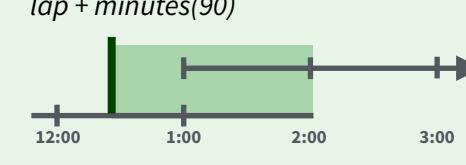
```
nor + minutes(90)
```



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



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

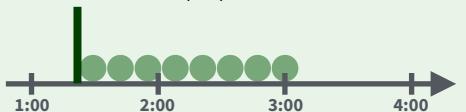


```
leap + years(1)
```

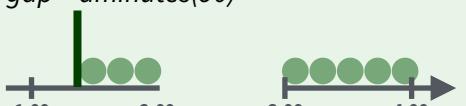


**Durations** track the passage of physical time, which deviates from clock time when irregularities occur.

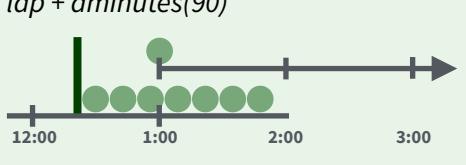
```
nor + dminutes(90)
```



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



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



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



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

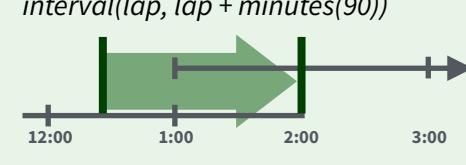
```
interval(nor, nor + minutes(90))
```



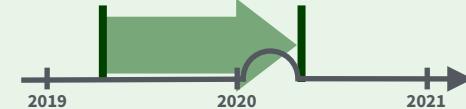
```
interval(gap, gap + minutes(90))
```



```
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.

Start Date End Date

```
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
```

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