

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

a;b;c 1;2;3 4;5;NA	→	A B C 1 2 3 4 5 NA
--------------------------	---	--------------------------------------

a b c 1 2 3 4 5 NA	→	A B C 1 2 3 4 5 NA
--------------------------	---	--------------------------------------

a b c 1 2 3 4 5 NA	→	A B C 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

Example file

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

1	2	3
4	5	NA

Skip lines

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

A B C 1 2 3 4 5 NA

No header

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

A B C 1 2 3

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(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
earn is a double (numeric)

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 large table to display

tibble display

A tibble: 234 x 6 manufacture... <chr> <dbl> 1 audi a4 1.8 2 audi a4 2.0 3 audi a4 2.0 4 audi a4 2.0 5 audi a4 2.0 6 audi a4 2.0 7 audi a4 2.0 8 audi a4 quattro 1.8 9 audi a4 quattro 1.8 10 audi a4 quattro 1.8 ... with 234 more rows: audi a4 quattro and 3 more variables: year <int>, cyl <int>, trans <chr>
--

data frame display

156 1999 6 auto(14) 157 2008 6 auto(14) 158 2008 6 auto(14) 159 2008 6 manual(5) 160 2008 6 manual(5) 161 1999 4 auto(14) 162 2008 4 manual(5) 163 2008 4 manual(5) 164 2008 4 auto(14) 165 2008 4 auto(14) 166 2008 4 auto(14) 167 2008 4 auto(14) 168 2008 4 auto(14) 169 2008 4 auto(14) 170 2008 4 auto(14) [reached getOption("max.print") -- omitted 68 rows]
--

- 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.

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

Both make this tibble

tribble(...)

Construct by rows.

`tribble(~x, ~y,
1, "a",
2, "b",
3, "c")`

A tibble: 3 x 2
x y
1 a
2 b
3 c

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

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

Convert named vector to a tibble

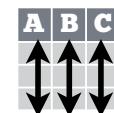
`is_tibble(x)` Test whether x is a tibble.



Tidy Data with tidyr

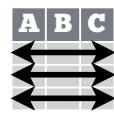
Tidy data is a way to organize tabular data. It provides a consistent data structure across packages.

A table is tidy if:



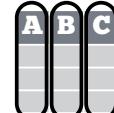
Each **variable** is in its own **column**

&



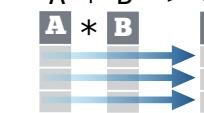
Each **observation**, or **case**, is in its own **row**

Tidy data:



Makes variables easy to access as vectors

$A * B \rightarrow C$



Preserves cases during vectorized operations

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

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

`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

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

`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

country	year
Afghan	1999
Afghan	2000
Brazil	1999
Brazil	2000
China	1999
China	2000

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

Handle Missing Values

drop_na(data, ...)

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

x1 x2
A 1 → A 1
B NA → D 3
C NA → E 3

`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 → A 1
B NA → B 1
C NA → C 1
D 3 → D 3
E NA → E 3

`fill(x, x2)`

replace_na(data, replace = list(), ...)

Replace NA's by column.

x1 x2
A 1 → A 1
B NA → B 2
C NA → C 2
D 3 → D 3
E NA → E 2

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

Expand Tables - quickly create tables with combinations of values

complete(data, ..., fill = list())

Adds to the data missing combinations of the values of the variables listed in ...

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

expand(data, ...)

Create new tibble with all possible combinations of the values of the variables listed in ...

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

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)`



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)`
- `sample_n(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame())` Randomly select size rows. `sample_n(iris, 10, 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)`

Logical and boolean operators to use with filter()

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

See `?base::logic` and `?Comparison` for help.

ARRANGE CASES

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

ADD CASES

`add_row(.data, ..., .before = NULL, .after = NULL)` Add one or more rows to a table.
`add_row(faithful, eruptions = 1, waiting = 1)`

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.

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	C	A	B	D
a	t	1	a	t	3
b	u	2	b	u	2
c	v	3	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.

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

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

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

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 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.

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 y

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

+ =

A	B	C
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 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>
```

`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.

R Studio

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",
  linemiter=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()
  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

l + geom_raster(aes(fill = z), hjust=0.5, vjust=0.5,
  interpolate=FALSE)
  x, y, alpha, fill

l + geom_tile(aes(fill = z))
  x, y, alpha, color, fill, linetype, size, width
```

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
```


Factors withforcats :: CHEAT SHEET

The **forcats** package provides tools for working with factors, which are R's data structure for categorical data.

Factors

R represents categorical data with factors. A **factor** is an integer vector with a **levels** attribute that stores a set of mappings between integers and categorical values. When you view a factor, R displays not the integers, but the values associated with them.

	stored	displayed
integer vector	<code>1 = a 2 = b 3 = c 1 = a</code>	<code>a b c a</code>
levels	<code>1 2 3 = c 1</code>	

`a
c
b
a` → `a
c
b
a` *Create a factor with factor()*
`factor(x = character(), levels, labels = levels, exclude = NA, ordered = is.ordered(x), nmax = NA)` Convert a vector to a factor. Also `as_factor()`.
`f <- factor(c("a", "c", "b", "a"), levels = c("a", "b", "c"))`

`a
c
b
a` → `a
b
c` *Return its levels with levels()*
`levels(x)` Return/set the levels of a factor. `levels(f); levels(f) <- c("x", "y", "z")`

Use unclass() to see its structure

Inspect Factors

`a
c
b
a` → `f
n` *fct_count(f, sort = FALSE)*
Count the number of values with each level. `fct_count(f)`

`a
c
b
a` → `a
c
b` *fct_unique(f)* Return the unique values, removing duplicates. `fct_unique(f)`

Combine Factors

`a
c
2=c + b
1=a
2=b
3=b` → `a
c
b
a` *fct_c(...)* Combine factors with different levels.
`f1 <- factor(c("a", "c"))
f2 <- factor(c("b", "a"))
fct_c(f1, f2)`

`a
b
1=a
2=b
a
c
1=a
2=c
3=c` → `a
b
c
1=a
2=b
3=c` *fct_unify(fs, levels = lvs_union(fs))* Standardize levels across a list of factors. `fct_unify(list(f2, f1))`



Change the order of levels

`a
c
b
a` → `a
c
b
a` *fct_relevel(f, ..., after = 0L)*
Manually reorder factor levels.
`fct_relevel(f, c("b", "c", "a"))`

`c
c
a` → `c
c
a` *fct_infreq(f, ordered = NA)*
Reorder levels by the frequency in which they appear in the data (highest frequency first).
`f3 <- factor(c("c", "c", "a"))
fct_infreq(f3)`

`b
a` → `b
a` *fct_inorder(f, ordered = NA)*
Reorder levels by order in which they appear in the data.
`fct_inorder(f2)`

`a
b
c` → `a
b
c` *fct_rev(f)* Reverse level order.
`f4 <- factor(c("a", "b", "c"))
fct_rev(f4)`

`a
b
c` → `a
b
c` *fct_shift(f)* Shift levels to left or right, wrapping around end.
`fct_shift(f4)`

`a
b
c` → `a
b
c` *fct_shuffle(f, n = 1L)* Randomly permute order of factor levels.
`fct_shuffle(f4)`

`a
b
c` → `b
c
a` *fct_reorder(f, .x, .fun = median, ..., .desc = FALSE)* Reorder levels by their relationship with another variable.
`boxplot(data = iris, Sepal.Width ~ fct_reorder(Species, Sepal.Width))`

`a
b
1=a
2=b
1=a
2=c
3=a` → `a
b
1=b
2=c
3=a` *fct_reorder2(f, .x, .y, .fun = last2, ..., .desc = TRUE)* Reorder levels by their final values when plotted with two other variables.
`ggplot(data = iris, aes(Sepal.Width, Sepal.Length, color = fct_reorder2(Species, Sepal.Width, Sepal.Length))) + geom_smooth()`

Change the value of levels

`a
c
b
a` → `v
z
x
v` *fct_recode(f, ...)* Manually change levels. Also `fct_relabel` which obeys purrr::map syntax to apply a function or expression to each level.
`fct_recode(f, v = "a", x = "b", z = "c")
fct_relabel(f, ~ paste0("x", .x))`

`a
c
b
a` → `2
1
3
2` *fct_anon(f, prefix = "")*
Anonymize levels with random integers. `fct_anon(f)`

`a
c
b
a` → `x
c
x
x` *fctCollapse(f, ...)* Collapse levels into manually defined groups.
`fctCollapse(f, x = c("a", "b"))`

`a
c
b
a` → `a
2=Other
Other
a` *fct_lump(f, n, prop, w = NULL, other_level = "Other", ties.method = c("min", "average", "first", "last", "random", "max"))* Lump together least/most common levels into a single level. Also `fct_lump_min`.
`fct_lump(f, n = 1)`

`a
c
b
a` → `a
2=Other
Other
b
a` *fct_other(f, keep, drop, other_level = "Other")* Replace levels with "other."
`fct_other(f, keep = c("a", "b"))`

Add or drop levels

`a
b
1=a
2=b
3=x` → `a
b
1=a
2=b` *fct_drop(f, only)* Drop unused levels.
`f5 <- factor(c("a", "b", "c", "a", "b", "x"))
f6 <- fct_drop(f5)`

`a
b
1=a
2=b` → `a
b
1=a
2=b
3=x` *fct_expand(f, ...)* Add levels to a factor.
`fct_expand(f6, "x")`

`a
b
NA` → `a
b
3=x` *fct_explicit_na(f, na_level = "Missing")*
Assigns a level to NAs to ensure they appear in plots, etc.
`fct_explicit_na(factor(c("a", "b", NA)))`

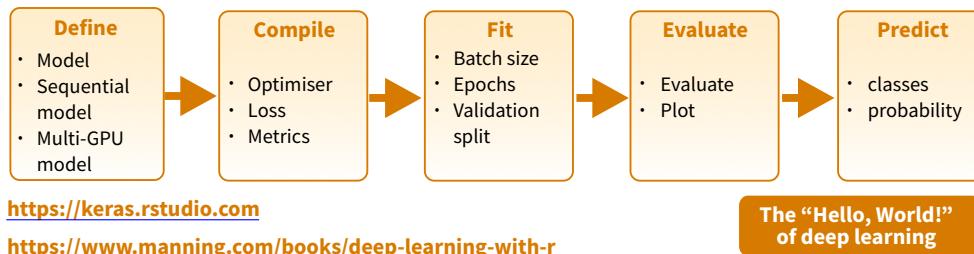
Deep Learning with Keras :: CHEAT SHEET



Intro

Keras is a high-level neural networks API developed with a focus on enabling fast experimentation. It supports multiple backends, including TensorFlow, CNTK and Theano.

TensorFlow is a lower level mathematical library for building deep neural network architectures. The keras R package makes it easy to use Keras and TensorFlow in R.



<https://keras.rstudio.com>

<https://www.manning.com/books/deep-learning-with-r>

INSTALLATION

The keras R package uses the Python keras library. You can install all the prerequisites directly from R.

https://keras.rstudio.com/reference/install_keras.html

```
library(keras)  
install_keras()
```

See ?install_keras
for GPU instructions

This installs the required libraries in an Anaconda environment or virtual environment 'r-tensorflow'.

Working with keras models

DEFINE A MODEL

`keras_model()` Keras Model

`keras_model_sequential()` Keras Model composed of a linear stack of layers

`multi_gpu_model()` Replicates a model on different GPUs

COMPILE A MODEL

`compile(object, optimizer, loss, metrics = NULL)`

Configure a Keras model for training

FIT A MODEL

`fit(object, x = NULL, y = NULL, batch_size = NULL, epochs = 10, verbose = 1, callbacks = NULL, ...)`
Train a Keras model for a fixed number of epochs (iterations)

`fit_generator()` Fits the model on data yielded batch-by-batch by a generator

`train_on_batch() test_on_batch()` Single gradient update or model evaluation over one batch of samples

EVALUATE A MODEL

`evaluate(object, x = NULL, y = NULL, batch_size = NULL)` Evaluate a Keras model

`evaluate_generator()` Evaluates the model on a data generator

PREDICT

`predict()` Generate predictions from a Keras model

`predict_proba() and predict_classes()`
Generates probability or class probability predictions for the input samples

`predict_on_batch()` Returns predictions for a single batch of samples

`predict_generator()` Generates predictions for the input samples from a data generator

OTHER MODEL OPERATIONS

`summary()` Print a summary of a Keras model

`export_savedmodel()` Export a saved model

`get_layer()` Retrieves a layer based on either its name (unique) or index

`pop_layer()` Remove the last layer in a model

`save_model_hdf5(); load_model_hdf5()` Save/Load models using HDF5 files

`serialize_model(); unserialize_model()`
Serialize a model to an R object

`clone_model()` Clone a model instance

`freeze_weights(); unfreeze_weights()`
Freeze and unfreeze weights

CORE LAYERS

 `layer_input()` Input layer

 `layer_dense()` Add a densely-connected NN layer to an output

 `layer_activation()` Apply an activation function to an output

 `layer_dropout()` Applies Dropout to the input

 `layer_reshape()` Reshapes an output to a certain shape

 `layer_permute()` Permute the dimensions of an input according to a given pattern

 `layer_repeat_vector()` Repeats the input n times

 `layer_lambda(object, f)` Wraps arbitrary expression as a layer

 `layer_activity_regularization()` Layer that applies an update to the cost function based input activity

 `layer_masking()` Masks a sequence by using a mask value to skip timesteps

 `layer_flatten()` Flattens an input

TRAINING AN IMAGE RECOGNIZER ON MNIST DATA

```
# input layer: use MNIST images  
mnist <- dataset_mnist()  
x_train <- mnist$train$x; y_train <- mnist$train$y  
x_test <- mnist$test$x; y_test <- mnist$test$y
```

5 0 4 1

```
# reshape and rescale  
x_train <- array_reshape(x_train, c(nrow(x_train), 784))  
x_test <- array_reshape(x_test, c(nrow(x_test), 784))  
x_train <- x_train / 255; x_test <- x_test / 255
```

```
y_train <- to_categorical(y_train, 10)  
y_test <- to_categorical(y_test, 10)
```

```
# defining the model and layers  
model <- keras_model_sequential()  
model %>%  
  layer_dense(units = 256, activation = 'relu',  
             input_shape = c(784)) %>%  
  layer_dropout(rate = 0.4) %>%  
  layer_dense(units = 128, activation = 'relu') %>%  
  layer_dense(units = 10, activation = 'softmax')
```

```
# compile (define loss and optimizer)  
model %>% compile(  
  loss = 'categorical_crossentropy',  
  optimizer = optimizer_rmsprop(),  
  metrics = c('accuracy'))
```

```
# train (fit)  
model %>% fit(  
  x_train, y_train,  
  epochs = 30, batch_size = 128,  
  validation_split = 0.2)  
model %>% evaluate(x_test, y_test)  
model %>% predict_classes(x_test)
```

More layers

CONVOLUTIONAL LAYERS



`layer_conv_1d()` 1D, e.g. temporal convolution



`layer_conv_2d_transpose()` Transposed 2D (deconvolution)



`layer_conv_2d()` 2D, e.g. spatial convolution over images



`layer_conv_3d_transpose()` Transposed 3D (deconvolution)
`layer_conv_3d()` 3D, e.g. spatial convolution over volumes



`layer_conv_lstm_2d()` Convolutional LSTM



`layer_separable_conv_2d()` Depthwise separable 2D



`layer_upsampling_1d()`
`layer_upsampling_2d()`
`layer_upsampling_3d()` Upsampling layer



`layer_zero_padding_1d()`
`layer_zero_padding_2d()`
`layer_zero_padding_3d()` Zero-padding layer



`layer_cropping_1d()`
`layer_cropping_2d()`
`layer_cropping_3d()` Cropping layer

POOLING LAYERS



`layer_max_pooling_1d()`
`layer_max_pooling_2d()`
`layer_max_pooling_3d()` Maximum pooling for 1D to 3D



`layer_average_pooling_1d()`
`layer_average_pooling_2d()`
`layer_average_pooling_3d()` Average pooling for 1D to 3D



`layer_global_max_pooling_1d()`
`layer_global_max_pooling_2d()`
`layer_global_max_pooling_3d()` Global maximum pooling



`layer_global_average_pooling_1d()`
`layer_global_average_pooling_2d()`
`layer_global_average_pooling_3d()` Global average pooling



R Studio

ACTIVATION LAYERS



`layer_activation(object, activation)`
Apply an activation function to an output



`layer_activation_leaky_relu()`
Leaky version of a rectified linear unit



`layer_activation_parametric_relu()`
Parametric rectified linear unit



`layer_activation_thresholded_relu()`
Thresholded rectified linear unit



`layer_activation_elu()`
Exponential linear unit

DROPOUT LAYERS



`layer_dropout()`
Applies dropout to the input



`layer_spatial_dropout_1d()`
`layer_spatial_dropout_2d()`
`layer_spatial_dropout_3d()` Spatial 1D to 3D version of dropout

RECURRENT LAYERS



`layer_simple_rnn()`
Fully-connected RNN where the output is to be fed back to input

`layer_gru()`
Gated recurrent unit - Cho et al

`layer_cudnn_gru()`
Fast GRU implementation backed by CuDNN

`layer_lstm()`
Long-Short Term Memory unit - Hochreiter 1997

`layer_cudnn_lstm()`
Fast LSTM implementation backed by CuDNN

LOCALLY CONNECTED LAYERS



`layer_locally_connected_1d()`
`layer_locally_connected_2d()`
Similar to convolution, but weights are not shared, i.e. different filters for each patch

Preprocessing

SEQUENCE PREPROCESSING

`pad_sequences()`
Pads each sequence to the same length (length of the longest sequence)

`skipgrams()`
Generates skipgram word pairs

`make_sampling_table()`
Generates word rank-based probabilistic sampling table

TEXT PREPROCESSING

`text_tokenizer()` Text tokenization utility

`fit_text_tokenizer()` Update tokenizer internal vocabulary

`save_text_tokenizer(); load_text_tokenizer()`
Save a text tokenizer to an external file

`texts_to_sequences(); texts_to_sequences_generator()`
Transforms each text in texts to sequence of integers

`texts_to_matrix(); sequences_to_matrix()`
Convert a list of sequences into a matrix

`text_one_hot()` One-hot encode text to word indices

`text_hashing_trick()`
Converts a text to a sequence of indexes in a fixed-size hashing space

`text_to_word_sequence()`
Convert text to a sequence of words (or tokens)

IMAGE PREPROCESSING

`image_load()` Loads an image into PIL format.

`flow_images_from_data()`
`flow_images_from_directory()`
Generates batches of augmented/normalized data from images and labels, or a directory

`image_data_generator()` Generate minibatches of image data with real-time data augmentation.

`fit_image_data_generator()` Fit image data generator internal statistics to some sample data

`generator_next()` Retrieve the next item

`image_to_array(); image_array_resize()`
`image_array_save()` 3D array representation

Pre-trained models

Keras applications are deep learning models that are made available alongside pre-trained weights. These models can be used for prediction, feature extraction, and fine-tuning.

`application_xception()`
`xception_preprocess_input()`
Xception v1 model

`application_inception_v3()`
`inception_v3_preprocess_input()`
Inception v3 model, with weights pre-trained on ImageNet

`application_inception_resnet_v2()`
`inception_resnet_v2_preprocess_input()`
Inception-ResNet v2 model, with weights trained on ImageNet

`application_vgg16(); application_vgg19()`
VGG16 and VGG19 models

`application_resnet50()` ResNet50 model

`application_mobilenet()`
`mobilenet_preprocess_input()`
`mobilenet_decode_predictions()`
`mobilenet_load_model_hdf5()`
MobileNet model architecture

IMAGENET

`ImageNet` is a large database of images with labels, extensively used for deep learning

`imagenet_preprocess_input()`
`imagenet_decode_predictions()`
Preprocesses a tensor encoding a batch of images for ImageNet, and decodes predictions

Callbacks

A callback is a set of functions to be applied at given stages of the training procedure. You can use callbacks to get a view on internal states and statistics of the model during training.

`callback_early_stopping()` Stop training when a monitored quantity has stopped improving
`callback_learning_rate_scheduler()` Learning rate scheduler
`callback_tensorboard()` TensorBoard basic visualizations

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()**.
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(), hms() 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(zone = "") Current time in tz (defaults to system tz). now()

today(zone = "") 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. **wday(dt)**
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)**

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

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

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

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

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

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

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")

Tip: use a date with day > 12

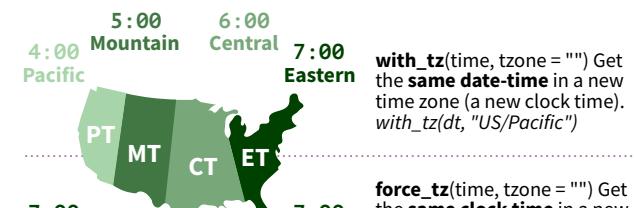
2. Apply the template to dates
sf(ymd("2010-04-05"))
[1] "Created Monday, Apr 05, 2010 00:00"

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()**



force_tz(time, tzzone = "") 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)

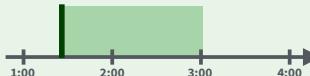
```
lap <- ymd_hms("2018-11-04 00:30:00", tz = "US/Eastern")
```

Leap years and leap seconds

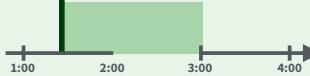
```
leap <- ymd("2019-03-01")
```

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

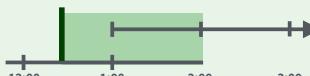
`nor + minutes(90)`



`gap + minutes(90)`



`lap + minutes(90)`

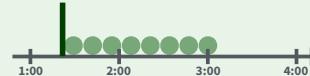


`leap + years(1)`



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

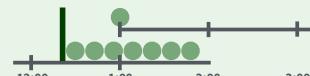
`nor + dminutes(90)`



`gap + dminutes(90)`



`lap + dminutes(90)`



`leap + dyears(1)`



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

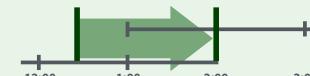
`interval(nor, nor + minutes(90))`



`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)
```

```
p
```

```
"3m 12d 0H 0M 0S"
```

Number of months Number of days etc.

```
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. **Difftimes** 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)`

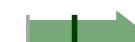
INTERVALS

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())`

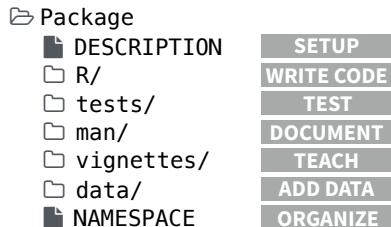
Package Development: : CHEAT SHEET



Package Structure

A package is a convention for organizing files into directories.

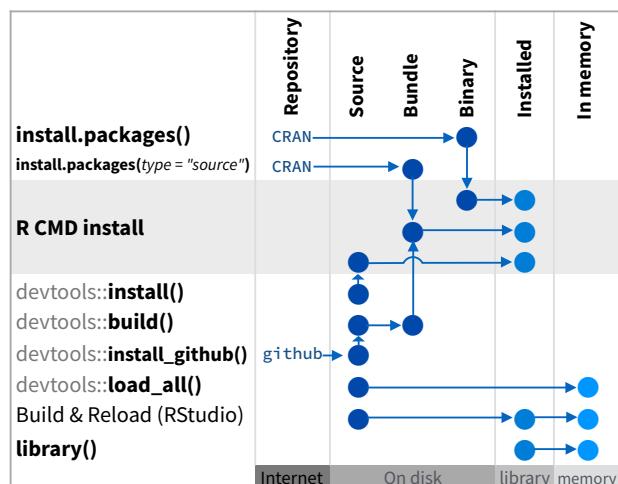
This sheet shows how to work with the 7 most common parts of an R package:



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

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

Or installed into an R library (loaded into memory during an R session) or archived online in a repository. Use the functions below to move between these states.



devtools::use_build_ignore("file")

Adds file to .Rbuildignore, a list of files that will not be included when package is built.

Setup (DESCRIPTION)

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

- You must have a `DESCRIPTION` file
- Add the packages that yours relies on with `devtools::use_package()`
Adds a package to the Imports or Suggests field

CC0

No strings attached.

MIT

MIT license applies to your code if re-shared.

GPL-2

GPL-2 license applies to your code, and all code anyone bundles with it, if re-shared.

Package: mypackage

Title: Title of Package

Version: 0.1.0

Authors@R: person("Hadley", "Wickham", email = "hadley@r-project.org", role = c("aut", "cre"))

Description: What the package does (one paragraph)

Depends: R (>= 3.1.0)

License: GPL-2

LazyData: true

Imports:

dplyr (>= 0.4.0),
ggvis (>= 0.2)

Suggests:
knitr (>= 0.1.0)

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

Suggest packages that are not very essential to yours. Users can install them manually, or not, as they like.

Write Code (R)

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

- Create a new package project with `devtools::create("path/to/name")`
Create a template to develop into a package.
- Save your code in `R/` as scripts (extension `.R`)

WORKFLOW

1. Edit your code.
2. Load your code with one of
`devtools::load_all()`
Re-loads all saved files in `R/` into memory.
3. Experiment in the console.
4. Repeat.
 - Use consistent style with r-pkgs.had.co.nz/r.html#style
 - Click on a function and press **F2** to open its definition
 - Search for a function with **Ctrl + .**



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

Test (tests/)

Use `tests/` to store tests that will alert you if your code breaks.

- Add a `tests/` directory
- Import `testthat` with `devtools::use_testthat()`, which sets up package to use automated tests with `testthat`
- Write tests with `context()`, `test()`, and `expect` statements
- Save your tests as `.R` files in `tests/testthat/`

WORKFLOW

1. Modify your code or tests.
2. Test your code with one of
`devtools::test()`
Runs all tests in `tests/`
3. Repeat until all tests pass

Example Test

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

Expect statement Tests

<code>expect_equal()</code>	is equal within small numerical tolerance?
<code>expect_identical()</code>	is exactly equal?
<code>expect_match()</code>	matches specified string or regular
<code>expect_output()</code>	prints specified output?
<code>expect_message()</code>	displays specified message?
<code>expect_warning()</code>	displays specified warning?
<code>expect_error()</code>	throws specified error?
<code>expect_is()</code>	output inherits from certain class?
<code>expect_false()</code>	returns FALSE?
<code>expect_true()</code>	returns TRUE?

Document (□ man/)

□ man/ contains the documentation for your functions, the help pages in your package.

- Use roxygen comments to document each function beside its definition
- Document the name of each exported data set
- Include helpful examples for each function

WORKFLOW

1. Add roxygen comments in your .R files
2. Convert roxygen comments into documentation with one of:

devtools::document()

Converts roxygen comments to .Rd files and places them in □ man/. Builds NAMESPACE.

Ctrl/Cmd + Shift + D (Keyboard Shortcut)

3. Open help pages with ? to preview documentation
4. Repeat

.Rd FORMATTING TAGS

\emph{italic text}	\email{name@foo.com}
\strong{bold text}	\href{url}{display}
\code{function(args)}	\url{url}
\pkg{package}	
\dontrun{code}	\link[=dest]{display}
\dontshow{code}	\linkS4class{class}
\donttest{code}	\code{\link{function}}
\deqn{a + b (block)}	\code{\link[package]{function}}
\eqn{a + b (inline)}	\tabular{lcr}{ left \tab centered \tab right \cr cell \tab cell \tab cell \cr}

Teach (□ vignettes/)

□ vignettes/ holds documents that teach your users how to solve real problems with your tools.

- Create a □ vignettes/ directory and a template vignette with devtools::use_vignette()
Adds template vignette as vignettes/my-vignette.Rmd.
- Append YAML headers to your vignettes (like right)
- Write the body of your vignettes in R Markdown (rmarkdown.rstudio.com)

ROXYGEN2

The **roxygen2** package lets you write documentation inline in your .R files with a shorthand syntax. devtools implements roxygen2 to make documentation.



- Add roxygen documentation as comment lines that begin with #’.
- Place comment lines directly above the code that defines the object documented.
- Place a roxygen @ tag (right) after #’ to supply a specific section of documentation.
- Untagged lines will be used to generate a title, description, and details section (in that order)

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

COMMON ROXYGEN TAGS

@aliases	@inheritParams	@seealso	
@concepts	@keywords	@format	
@describeln	@param	@source	data
@examples	@rdname	@include	
@export	@return	@slot	S4
	@family	@section	RC

Add Data (□ data/)

The □ data/ directory allows you to include data with your package.



- Save data as .Rdata files (suggested)
- Store data in one of **data/**, **R/Sysdata.rda**, **inst/extdata**
- Always use **LazyData: true** in your DESCRIPTION file.

devtools::use_data()

Adds a data object to data/ (R/Sysdata.rda if **internal = TRUE**)

devtools::use_data_raw()

Adds an R Script used to clean a data set to data-raw/. Includes data-raw/ on .Rbuildignore.

Store data in

- **data/** to make data available to package users
- **R/sysdata.rda** to keep data internal for use by your functions.
- **inst/extdata** to make raw data available for loading and parsing examples. Access this data with **system.file()**

Organize (□ NAMESPACE)

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

- Export functions for users by placing **@export** in their roxygen comments
- Import objects from other packages with **package::object** (recommended) or **@import**, **@importFrom**, **@importClassesFrom**, **@importMethodsFrom** (not always recommended)

WORKFLOW

1. Modify your code or tests.
2. Document your package (devtools::document())
3. Check NAMESPACE
4. Repeat until NAMESPACE is correct

SUBMIT YOUR PACKAGE

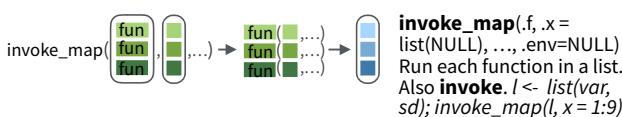
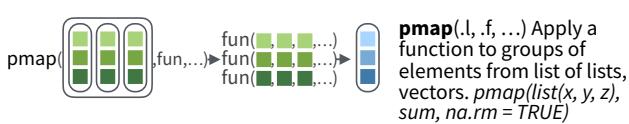
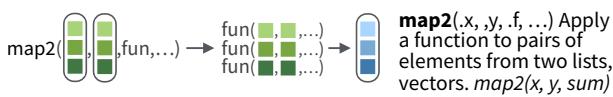
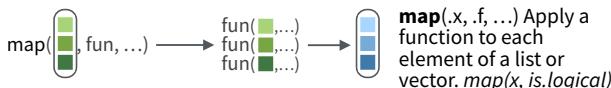
r-pkgs.had.co.nz/release.html

Apply functions with purrr :: CHEAT SHEET



Apply Functions

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



imap(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:

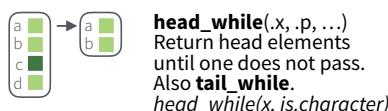
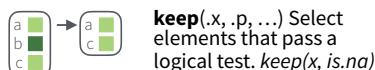
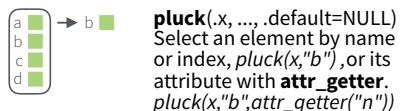
"name" becomes `function(x) x[["name"]]`, e.g. `map(l, "a")` extracts `a` from each element of `l`

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

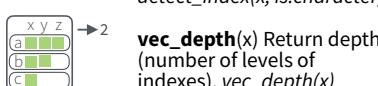
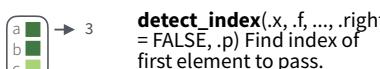
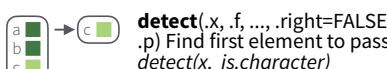
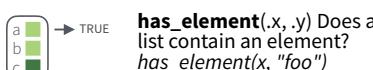
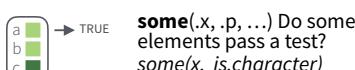
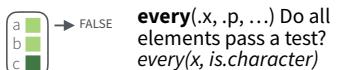
`~ .x ..2` etc becomes `function(.1, .2, etc) ..1 ..2` etc, e.g. `pmap(list(a, b, c), ~ ..3 + ..1 - ..2)` becomes `pmap(list(a, b, c), function(a, b, c) c + a - b)`

Work with Lists

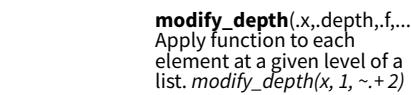
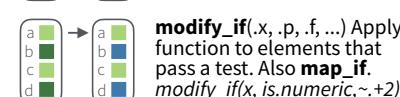
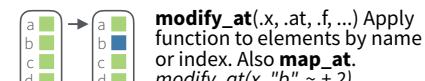
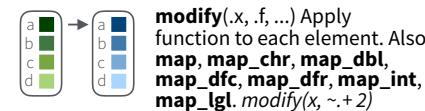
FILTER LISTS



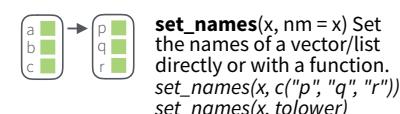
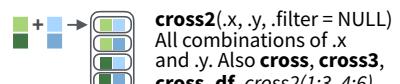
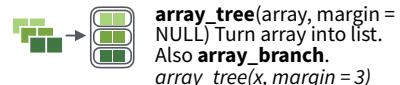
SUMMARISE LISTS



TRANSFORM LISTS



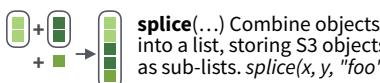
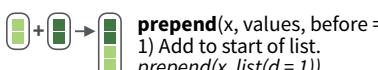
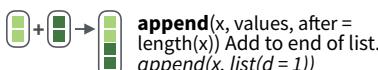
WORK WITH LISTS



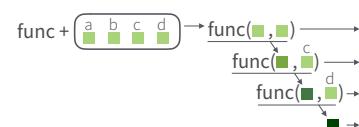
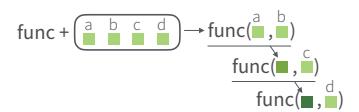
RESHAPE LISTS



JOIN (TO) LISTS



Reduce Lists



compose() Compose multiple functions.

lift() Change the type of input a function takes. Also `lift_dl`, `lift_kv`, `lift_ld`, `lift_lv`, `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 and errors.

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).

Modify function behavior

Nested Data

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

nested data frame

Species	data
setosa	<tibble [50 x 4]>
versicolor	<tibble [50 x 4]>
virginica	<tibble [50 x 4]>

n_iris

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
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	5.8	2.7	5.1	1.9
versi	7.1	3.0	5.9	2.1
versi	6.3	2.9	5.6	1.8
versi	6.5	3.0	5.8	2.2
virgini	5.1	3.5	1.4	0.2
virgini	4.9	3.0	1.4	0.2
virgini	4.7	3.2	1.3	0.2
virgini	4.6	3.1	1.5	0.2
virgini	5.0	3.6	1.4	0.2
virgini	5.5	2.3	4.0	1.3
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()

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.

"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]]

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 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
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
versi	6.5	3.3	6.0	2.5
versi	5.8	2.7	5.1	1.9
versi	7.1	3.0	5.9	2.1
versi	6.3	2.9	5.6	1.8
versi	6.5	3.0	5.8	2.2
virgini	5.1	3.5	1.4	0.2
virgini	4.9	3.0	1.4	0.2
virgini	4.7	3.2	1.3	0.2
virgini	4.6	3.1	1.5	0.2
virgini	5.0	3.6	1.4	0.2
virgini	5.5	2.3	4.0	1.3
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()

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	5.8	2.7	5.1	1.9
versi	7.1	3.0	5.9	2.1
versi	6.3	2.9	5.6	1.8
versi	6.5	3.0	5.8	2.2
virgini	5.1	3.5	1.4	0.2
virgini	4.9	3.0	1.4	0.2
virgini	4.7	3.2	1.3	0.2
virgini	4.6	3.1	1.5	0.2
virgini	5.0	3.6	1.4	0.2
virgini	5.5	2.3	4.0	1.3
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 [50 x 4]>	lm(S.L ~ ., df)
versi	<tibble [50 x 4]>	lm(S.W ~ ., df)
virgini	<tibble [50 x 4]>	lm(P.L ~ ., df)

```
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))
```

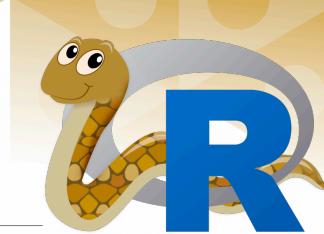
3. SIMPLIFY THE LIST COLUMN (into a regular column)

Use the purrr functions **map_lgl()**, **map_int()**, **map_dbl()**, **map_chr()**, as well as tidy's **unnest()** to reduce a list column into a regular column.

```
purrr::map_lgl(.x, f, ...)  
Apply .f element-wise to .x, return a logical vector  
n_iris %>% transmute(n = map_lgl(data, is.matrix))  
  
purrr::map_int(.x, f, ...)  
Apply .f element-wise to .x, return an integer vector  
n_iris %>% transmute(n = map_int(data, nrow))
```

```
purrr::map_dbl(.x, f, ...)  
Apply .f element-wise to .x, return a double vector  
n_iris %>% transmute(n = map_dbl(data, nrow))  
  
purrr::map_chr(.x, f, ...)  
Apply .f element-wise to .x, return a character vector  
n_iris %>% transmute(n = map_chr(data, nrow))
```

Use Python with R with reticulate :: CHEAT SHEET



The `reticulate` package lets you use Python and R together seamlessly in R code, in R Markdown documents, and in the RStudio IDE.

Python in R Markdown

(Optional) Build Python env to use.

Add `knitr::knit_engines$set(python = reticulate::eng_python)` to the setup chunk to set up the reticulate Python engine (not required for `knitr >= 1.18`).

Suggest the Python environment to use, in your setup chunk.

Begin Python chunks with ````{python}`. Chunk options like `echo`, `include`, etc. all work as expected.

Use the `py` object to access objects created in Python chunks from R chunks.

Python chunks all execute within a **single** Python session so you have access to all objects created in previous chunks.

Use the `r` object to access objects created in R chunks from Python chunks.

Output displays below chunk, including matplotlib plots.

A screenshot of the RStudio interface showing an R Markdown file. It contains several code blocks: an R chunk for setting up a Python environment, followed by two Python chunks (one for fmri data and one for plotting). Then there are R chunks for subset selection and plotting. The RStudio interface shows the R Console and R Markdown tabs at the bottom.

```
1 ````{r setup, include = FALSE}
2 library(reticulate)
3 virtualenv_create("fmri-proj")
4 py_install("seaborn", envname = "fmri-proj")
5 use_virtualenv("fmri-proj")
6 ````

8 ````{python, echo = FALSE}
9 import seaborn as sns
10 fmri = sns.load_dataset("fmri")
11 ````

13 ````{r}
14 f1 <- subset(py$fmri, region == "parietal")
15 ````

17 ````{python}
18 import matplotlib as mpl
19 sns.lmplot("timepoint", "signal", data = r.f1)
20 mpl.pyplot.show()
21 ````
```

A screenshot of the RStudio interface showing R code. It imports reticulate and seaborn, uses virtualenv, and then imports sns. It loads the fmri dataset and prints its dimensions. It then runs a Python script named python.py, which creates a tips dataset and prints its shape. The RStudio interface shows the R Script tab at the bottom.

```
1 library(reticulate)
2 py_install("seaborn")
3 use_virtualenv("r-reticulate")
4
5 sns <- import("seaborn")
6
7 fmri <- sns$load_dataset("fmri")
8 dim(fmri)
9
10 # creates tips
11 source_python("python.py")
12 dim(tips)
13
14 # creates tips in main
15 py_run_file("python.py")
16 dim(py$tips)
17
18 py_run_string("print(tips.shape)")
19
```

Object Conversion

Tip: To index Python objects begin at 0, use integers, e.g. `0L`

Reticulate provides **automatic** built-in conversion between Python and R for many Python types.

R	↔	Python
Single-element vector		Scalar
Multi-element vector		List
List of multiple types		Tuple
Named list		Dict
Matrix/Array		NumPy ndarray
Data Frame		Pandas DataFrame
Function		Python function
NULL, TRUE, FALSE		None, True, False

Or, if you like, you can convert manually with

`py_to_r(x)` Convert a Python object to an R object. Also `r_to_py`. `py_to_r(x)`

`tuple(..., convert = FALSE)` Create a Python tuple. `tuple("a", "b", "c")`

`dict(..., convert = FALSE)` Create a Python dictionary object. Also `py_dict` to make a dictionary that uses Python objects as keys. `dict(foo = "bar", index = 42L)`

`np_array(data, dtype = NULL, order = "C")` Create NumPy arrays. `np_array(c(1:8), dtype = "float16")`

`array_reshape(x, dim, order = c("C", "F"))` Reshape a Python array. `x <- 1:4; array_reshape(x, c(2, 2))`

`py_func(object)` Wrap an R function in a Python function with the same signature. `py_func(xor)`

`py_main_thread_func(object)` Create a function that will always be called on the main thread.

`iterate(..., convert = FALSE)` Apply an R function to each value of a Python iterator or return the values as an R vector, draining the iterator as you go. Also `iter_next` and `as_iterator`. `iterate(iter, print)`

`py_iterator(fn, completed = NULL)` Create a Python iterator from an R function. `seq_gen <- function(x){n <- x; function(){n <- n + 1; n}}; py_iterator(seq_gen(9))`

Helpers

`py_capture_output(expr, type = c("stdout", "stderr"))` Capture and return Python output. Also `py_suppress_warnings`. `py_capture_output("x")`

`py_get_attr(x, name, silent = FALSE)` Get an attribute of a Python object. Also `py_set_attr`, `py_has_attr`, and `py_list_attributes`. `py_get_attr(x)`

`py_help(object)` Open the documentation page for a Python object. `py_help(sns)`

`py_last_error()` Get the last Python error encountered. Also `py_clear_last_error` to clear the last error. `py_last_error()`

`py_save_object(object, filename, pickle = "pickle")` Save and load Python objects with pickle. Also `py_load_object`. `py_save_object(x, "x.pickle")`

`with(data, expr, as = NULL, ...)` Evaluate an expression within a Python context manager. `py <- import_builtin(); with(py$open("output.txt", "w"))%as%file, {file$write("Hello, there!")})`

Python in R code

Call Python from R in three ways:

IMPORT PYTHON MODULES

Use `import()` to import any Python module. Access the attributes of a module with `$`.

- `import(module, as = NULL, convert = TRUE, delay_load = FALSE)` Import a Python module. If `convert = TRUE`, Python objects are converted to their equivalent R types. Also `import_from_path`. `import("pandas")`
- `import_main(convert = TRUE)` Import the main module, where Python executes code by default. `import_main()`
- `import_builtins(convert = TRUE)` Import Python's built-in functions. `import_builtins()`

SOURCE PYTHON FILES

Use `source_python()` to source a Python script and make the Python functions and objects it creates available in the calling R environment.

- `source_python(file, envir = parent.frame(), convert = TRUE)` Run a Python script, assigning objects to a specified R environment. `source_python("file.py")`

RUN PYTHON CODE

Execute Python code into the **main** Python module with `py_run_file()` or `py_run_string()`.

- `py_run_string(code, local = FALSE, convert = TRUE)` Run Python code (passed as a string) in the main module. `py_run_string("x = 10"); py$x`
- `py_run_file(file, local = FALSE, convert = TRUE)` Run Python file in the main module. `py_run_file("script.py")`
- `py_eval(code, convert = TRUE)` Run a Python expression, return the result. Also `py_call`. `py_eval("1 + 1")`

Access the results, and anything else in Python's **main** module, with `py`.

- `py` An R object that contains the Python main module and the results stored there. `py$x`

Python in the IDE

Requires reticulate plus RStudio v1.2 or higher.

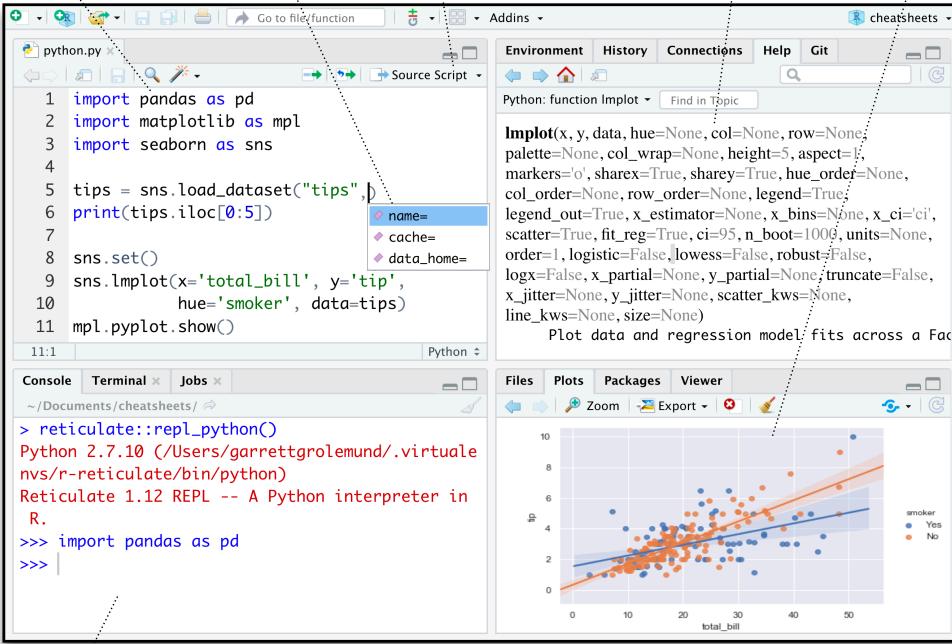
Syntax highlighting for Python scripts and chunks

Tab completion for Python functions and objects (and Python modules imported in R scripts)

Source Python scripts.

Execute Python code line by line with Cmd + Enter (Ctrl + Enter)

Press F1 over a Python symbol to display the help topic for that symbol.

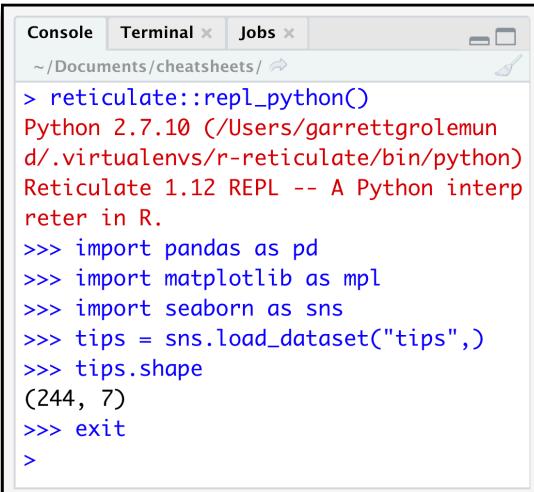


A Python REPL opens in the console when you run Python code with a keyboard shortcut. Type **exit** to close.

Python REPL

A REPL (Read, Eval, Print Loop) is a command line where you can run Python code and view the results.

1. Open in the console with **repl_python()**, or by running code in a Python script with **Cmd + Enter (Ctrl + Enter)**.
 - **repl_python(module = NULL, quiet =getOption("reticulate.repl.quiet", default = FALSE))** Launch a Python REPL. Run **exit** to close. *repl_python()*
2. Type commands at **>>>** prompt
3. Press **Enter** to run code
4. Type **exit** to close and return to R console



Configure Python

Reticulate binds to a local instance of Python when you first call **import()** directly or implicitly from an R session. To control the process, find or build your desired Python instance. Then suggest your instance to reticulate. **Restart R to unbind.**



Find Python

- **py_discover_config()** Return all detected versions of Python. Use **py_config** to check which version has been loaded. *py_config()*
- **py_available(initialize = FALSE)** Check if Python is available on your system. Also **py_module_available**, **py_numpy_module**, **py_available()**

Create a Python env

- **virtualenv_create(envname)** Create a new virtualenv. *virtualenv_create("r-pandas")*
- **conda_create(envname, packages = NULL, conda = "auto")** Create a new Conda env. *conda_create("r-pandas", packages = "pandas")*

Install Packages

Install Python packages with R (below) or the shell:
pip install SciPy
conda install SciPy

- **py_install(packages, envname = "r-reticulate", method = c("auto", "virtualenv", "conda"), conda = "auto", ...)** Installs Python packages into a Python env named "r-reticulate". *py_install("pandas")*
- **virtualenv_install(envname, packages, ignore_installed = FALSE)** Install a package within a virtualenv. *virtualenv_install("r-pandas", packages = "pandas")*
- **virtualenv_remove(envname, packages = NULL, confirm = interactive())** Remove individual packages or an entire virtualenv. *virtualenv_remove("r-pandas", packages = "pandas")*
- **virtualenv_remove(envname, packages = NULL, confirm = interactive())** Remove individual packages or an entire virtualenv. *virtualenv_remove("r-pandas", packages = "pandas")*
- **conda_install(envname, packages, forge = TRUE, pip = FALSE, pip_ignore_installed = TRUE, conda = "auto")** Install a package within a Conda env. *conda_install("r-pandas", packages = "plotly")*
- **conda_remove(envname, packages = NULL, conda = "auto")** Remove individual packages or an entire Conda env. *conda_remove("r-pandas", packages = "plotly")*

Suggest an env to use

To choose an instance of Python to bind to, reticulate scans the instances on your computer in the following order, stopping at the first instance that contains the module called by **import()**.

1. The instance referenced by the environment variable **RETICULATE_PYTHON** (if specified). **Tip:** set in *.Renviron* file.

• **sys.setenv(RETICULATE_PYTHON = PATH)**
Set default Python binary. Persists across sessions! Undo with **sys.unsetenv**.
sys.setenv(RETICULATE_PYTHON = "/usr/local/bin/python")

2. The instances referenced by **use_** functions if called before **import()**. Will fail silently if called after **import** unless **required = TRUE**.

• **use_python(python, required = FALSE)**
Suggest a Python binary to use by path.
use_python("/usr/local/bin/python")

• **use_virtualenv(virtualenv = NULL, required = FALSE)**
Suggest a Python virtualenv.
use_virtualenv("~/myenv")

• **use_condaenv(condaenv = NULL, conda = "auto", required = FALSE)**
Suggest a Conda env to use.
use_condaenv(condaenv = "r-nlp", conda = "/opt/anaconda3/bin/conda")

3. Within virtualenvs and conda envs that carry the same name as the imported module. e.g. *~/anaconda/envs/nltk* for **import("nltk")**

4. At the location of the Python binary discovered on the system PATH (i.e. **sys.which("python")**)

5. At customary locations for Python, e.g. */usr/local/bin/python*, */opt/local/bin/python*...

R Markdown :: CHEAT SHEET

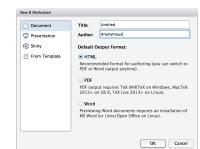
What is R Markdown?

.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

Embed code with knitr syntax

INLINE CODE

Insert with `r <code>` . Results appear as text without code.
Built with `r getVersion()` → 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 {{r}}

```  
{{r echo=TRUE}}  
getVersion()  
...````

getVersion()  
## [1] '3.2.3'

### 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)

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)

Options not listed above: R.options, anopts, 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



.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 wth Knit with parameters or the params argument of render():

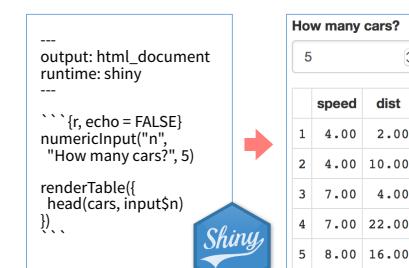
render("doc.Rmd", params = list(n = 1, d = as.Date("2015-01-01")))

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



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:
```

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

```
> 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] 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.
```

```
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 = mc2
```

```
block quote
```

Header1

Header 2

Header 3

Header 4

Header 5

Header 6

```
HTML ignored in pdfs
```

```
http://www.rstudio.com
```

```
link
```

```
Jump to Header 1
```

```
image:
```



```
Caption
```

- 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

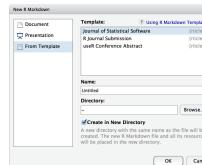
2. In the directory, Place a folder that contains:

`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

citation_package	The LaTeX package to process citations, natbib, biblatex or none	html	X	pdf		word		odt		rtf	X	md	X	gitbook		ioslides		slidy	X	beamer	
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
dev	Graphics device to use for figure output (e.g. "png")																		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
fig_height, fig_width	Default figure height and width (in inches) for document																		X	X	X
highlight	Syntax highlighting: "tango", "pygments", "kate", "zenburn", "textmate"																		X	X	X
includes	File of content to place in document (in_header, before_body, after_body)																		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
keep_tex	Save a copy of .tex file that contains knitr output																			X	
latex_engine	Engine to render latex, "pdflatex", "xelatex", or "luatex"																			X	
lib_dir	Directory of dependency files to use (Bootstrap, MathJax, etc.)																			X	X
mathjax	Set to local or a URL to use a local/URL version of MathJax to render equations																			X	X
md_extensions	Markdown extensions to add to default definition or R Markdown																		X	X	X
number_sections	Add section numbering to headers																		X	X	
pandoc_args	Additional arguments to pass to Pandoc																		X		
preserve_yaml	Preserve YAML front matter in final document?																				
reference_docx	docx file whose styles should be copied when producing docx output																		X		
self_contained	Embed dependencies into the doc																		X	X	X
slide_level	The lowest heading level that defines individual slides																				
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
theme	Bootswatch or Beamer theme to use for page																		X		
toc	Add a table of contents at start of document																		X	X	X
toc_depth	The lowest level of headings to add to table of contents																		X	X	X
toc_float	Float the table of contents to the left of the main content																				

Learn more in the `stargazer`, `xtable`, and `knitr` packages.

Table Suggestions

Several functions format R data into tables

Table with kable	
eruptions	waiting
3.600	79
1.800	54
3.333	74
2.282	62

Table with stargazer

eruptions		waiting
1	3.60	79.00
2	1.80	54.00
3	3.33	74.00
4	2.28	62.00

`data <- faithful[1:4,]`

````{r results = "asis"}`

`knitr::kable(data, caption = "Table with kable")`

````{r results = "asis"}`

`print(xtable::xtable(data, caption = "Table with xtable"),`

`type = "html", html.table.attributes = "border=0")`

````{r results = "asis"}`

`stargazer::stargazer(data, type = "html", title = "Table with stargazer")`

`````

3. Render. Bibliography will be added to end of document

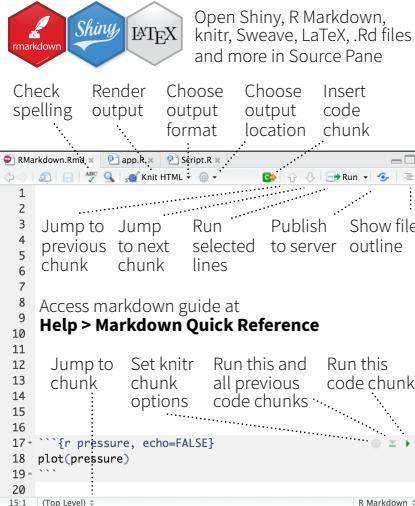
Smith cited (@smith04).
Smith cited without author (@smith04).
@smith04 cited in line.

Smith cited (Joe Smith 2004).
Smith cited without author (2004).
Joe Smith (2004) cited in line.

RStudio® is a trademark of RStudio, Inc. • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more at rmarkdown.rstudio.com • rmarkdown 1.6 • Updated: 2016-02

RStudio IDE :: CHEAT SHEET

Documents and Apps



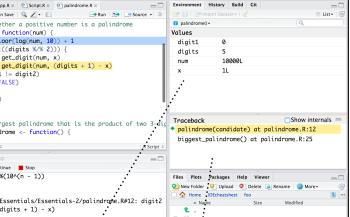
Debug Mode

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



Step through code one line at a time



Step into and out of functions to run



Resume execution



Quit debug mode



Write Code

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



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.

Change file type

Jump to function in file
Working Directory
Maximize, minimize panes
Press ↑ to see command history
Drag pane boundaries

Path to displayed directory

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

R Support

Import data with wizard History of past commands to run/copy Display .RPres slideshows

File > New File > R Presentation

Start new R Session in current project Close R Session in project Select R Version



Load workspace Save workspace Delete all saved objects Search inside environment

Choose environment to display from list of parent environments

Data 150 obs. of 5 variables Values 1 Functions foo

Displays saved objects by type with short description View in data viewer View function source code

Files Plots Packages Help Viewer

New Folder Upload Delete Rename More

Create folder Upload file Delete file Rename file Change directory

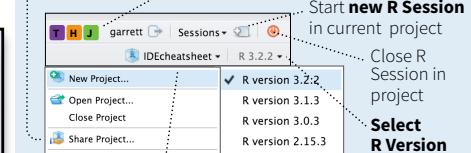
Path to displayed directory

19 B April 13, 2016, 11:17 AM

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

Pro Features

Share Project Active shared with Collaborators.. collaborators



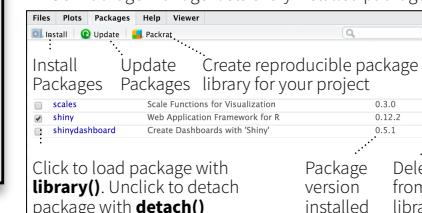
PROJECT SYSTEM File > New Project

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

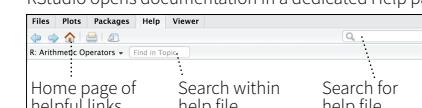
RStudio opens plots in a dedicated Plots pane



GUI Package manager lists every installed package



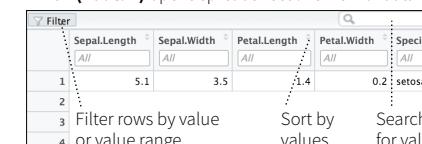
RStudio opens documentation in a dedicated Help pane



Viewer Pane displays HTML content, such as Shiny apps, R Markdown 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 |

4 WRITE CODE

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

| Windows /Linux | Mac |
|--------------------------|------------------------|
| Tab or Ctrl+Space | Tab or Cmd+Space |
| ↑/↓ | ↑/↓ |
| Enter, Tab, or → | Enter, Tab, or → |
| Esc | Esc |
| Ctrl+Z | Cmd+Z |
| Ctrl+Shift+Z | Cmd+Shift+Z |
| Ctrl+X | Cmd+X |
| Ctrl+C | Cmd+C |
| Ctrl+V | Cmd+V |
| Ctrl+A | Cmd+A |
| Ctrl+D | Cmd+D |
| Shift+[Arrow] | Shift+[Arrow] |
| Select Word | Select Word |
| Ctrl+Shift+←/→ | Option+Shift+←/→ |
| Cmd+Shift+← | Cmd+Shift+← |
| Alt+Shift+← | Alt+Shift+← |
| Alt+Shift+→ | Alt+Shift+→ |
| Select to Line Start | Select to Line Start |
| Select to Line End | Select to Line End |
| Select Page Up/Down | Select Page Up/Down |
| Select to Start/End | Select to Start/End |
| Delete Word Left | Delete Word Left |
| Delete Word Right | Delete Word Right |
| Delete to Line End | Delete to Line End |
| Delete to Line Start | Delete to Line Start |
| Indent | Indent |
| Outdent | Outdent |
| Yank line up to cursor | Yank line up to cursor |
| Yank line after cursor | Yank line after cursor |
| Insert yanked text | Insert yanked text |
| Insert <-> | Alt+- |
| Insert %>% | Option+- |
| Ctrl+Shift+M | Cmd+Shift+M |
| F1 | F1 |
| F2 | F2 |
| Ctrl+Shift+N | Cmd+Shift+N |
| New document (Chrome) | Ctrl+Alt+Shift+N |
| Open document | Ctrl+O |
| Save document | Ctrl+S |
| Close document | Ctrl+W |
| Close document (Chrome) | Ctrl+Alt+W |
| Close all documents | Ctrl+Shift+W |
| Extract function | Ctrl+Alt+X |
| Extract variable | Ctrl+Alt+V |
| Reindent lines | Ctrl+I |
| (Un)Comment lines | Ctrl+Shift+C |
| Reflow Comment | Ctrl+Shift+/ |
| Reformat Selection | Ctrl+Shift+A |
| Select within braces | Ctrl+Shift+E |
| Show Diagnostics | Ctrl+Shift+Alt+P |
| Transpose Letters | Ctrl+T |
| Move Lines Up/Down | Alt+↑/↓ |
| Copy Lines Up/Down | Shift+Alt+↑/↓ |
| Add New Cursor Above | Ctrl+Alt+Up |
| Add New Cursor Below | Ctrl+Alt+Down |
| Move Active Cursor Up | Ctrl+Alt+Shift+Up |
| Move Active Cursor Down | Ctrl+Alt+Shift+Down |
| Find and Replace | Ctrl+F |
| Use Selection for Find | Ctrl+F3 |
| Replace and Find | 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/



2 RUN CODE

Search command history

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

| Windows/Linux | Mac |
|---------------|--------------|
| Ctrl+↑ | Cmd+↑ |
| ↑/↓ | ↑/↓ |
| Home | Cmd+← |
| End | Cmd+→ |
| Ctrl+Shift+H | Ctrl+Shift+H |

Interrupt current command

Clear console

Quit Session (desktop only)

| Windows/Shift+F10 | Mac/Shift+F10 |
|-------------------|---------------|
| Ctrl+Enter | Cmd+Enter |

Restart R Session

Run current line/selection

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

| Windows/Shift+S | Mac/Shift+S |
|------------------|-----------------|
| Ctrl+Shift+Enter | Cmd+Shift+Enter |

3 NAVIGATE CODE

Goto File/Function

Fold Selected
 Unfold Selected
 Fold All
 Unfold All
 Go to line

| Windows /Linux |
|----------------|
| Ctrl+. |

| Mac |
|------------------------|
| Ctrl+. |
| Alt+L |
| Cmd+Option+L |
| Shift+Alt+L |
| Cmd+Shift+Option+L |
| Alt+O |
| Cmd+Option+O |
| Shift+Alt+O |
| Cmd+Shift+Option+O |
| Shift+Alt+G |
| Cmd+Shift+Option+G |
| Shift+Alt+J |
| Cmd+Shift+Option+J |
| Ctrl+Shift+. |
| Ctrl+Shift+. |
| Ctrl+F11 |
| Ctrl+F11 |
| Ctrl+F12 |
| Ctrl+F12 |
| Ctrl+Shift+F11 |
| Ctrl+Shift+F11 |
| Ctrl+Shift+F12 |
| Ctrl+Shift+F12 |
| Ctrl+F9 |
| Cmd+F9 |
| Ctrl+F10 |
| Cmd+F10 |
| Ctrl+P |
| Ctrl+P |
| Ctrl+Shift+Alt+E |
| Ctrl+Shift+Option+E |
| Ctrl+F3 |
| Cmd+E |
| Ctrl+Shift+F |
| Cmd+Shift+F |
| Win: F3, Linux: Ctrl+G |
| Cmd+G |
| W: Shift+F3, L: |
| Cmd+Shift+G |
| Ctrl+←/→ |
| Option+←/→ |
| Ctrl+↑/↓ |
| Cmd+↑/↓ |
| Ctrl+Shift+O |

Source the current file

Source with echo

4 WRITE CODE

Attempt completion

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

| Windows /Linux | Mac |
|--------------------------|------------------------|
| Tab or Ctrl+Space | Tab or Cmd+Space |
| ↑/↓ | ↑/↓ |
| Enter, Tab, or → | Enter, Tab, or → |
| Esc | Esc |
| Ctrl+Z | Cmd+Z |
| Ctrl+Shift+Z | Cmd+Shift+Z |
| Ctrl+X | Cmd+X |
| Ctrl+C | Cmd+C |
| Ctrl+V | Cmd+V |
| Ctrl+A | Cmd+A |
| Ctrl+D | Cmd+D |
| Shift+[Arrow] | Shift+[Arrow] |
| Select Word | Select Word |
| Ctrl+Shift+←/→ | Option+Shift+←/→ |
| Cmd+Shift+← | Cmd+Shift+← |
| Alt+Shift+← | Alt+Shift+← |
| Alt+Shift+→ | Alt+Shift+→ |
| Select to Line Start | Select to Line Start |
| Select to Line End | Select to Line End |
| Select Page Up/Down | Select Page Up/Down |
| Select to Start/End | Select to Start/End |
| Delete Word Left | Delete Word Left |
| Delete Word Right | Delete Word Right |
| Delete to Line End | Delete to Line End |
| Delete to Line Start | Delete to Line Start |
| Indent | Indent |
| Outdent | Outdent |
| Yank line up to cursor | Yank line up to cursor |
| Yank line after cursor | Yank line after cursor |
| Insert yanked text | Insert yanked text |
| Insert <-> | Alt+- |
| Insert %>% | Option+- |
| Ctrl+Shift+M | Cmd+Shift+M |
| F1 | F1 |
| F2 | F2 |
| Ctrl+Shift+N | Cmd+Shift+N |
| New document (Chrome) | Ctrl+Alt+Shift+N |
| Open document | Ctrl+O |
| Save document | Ctrl+S |
| Close document | Ctrl+W |
| Close document (Chrome) | Ctrl+Alt+W |
| Close all documents | Ctrl+Shift+W |
| Extract function | Ctrl+Alt+X |
| Extract variable | Ctrl+Alt+V |
| Reindent lines | Ctrl+I |
| (Un)Comment lines | Ctrl+Shift+C |
| Reflow Comment | Ctrl+Shift+/ |
| Reformat Selection | Ctrl+Shift+A |
| Select within braces | Ctrl+Shift+E |
| Show Diagnostics | Ctrl+Shift+Alt+P |
| Transpose Letters | Ctrl+T |
| Move Lines Up/Down | Alt+↑/↓ |
| Copy Lines Up/Down | Shift+Alt+↑/↓ |
| Add New Cursor Above | Ctrl+Alt+Up |
| Add New Cursor Below | Ctrl+Alt+Down |
| Move Active Cursor Up | Ctrl+Alt+Shift+Up |
| Move Active Cursor Down | Ctrl+Alt+Shift+Down |
| Find and Replace | Ctrl+F |
| Use Selection for Find | Ctrl+F3 |
| Replace and Find | Ctrl+Shift+J |

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

| Windows /Linux | Mac |
|----------------|---------------|
| Ctrl+Alt+D | Ctrl+Option+D |
| Ctrl+Alt+M | Ctrl+Option+M |
| Ctrl+↑/↓ | Ctrl+↑/↓ |
| Spacebar | Spacebar |
| Enter | Enter |

7 MAKE PACKAGES

Build and Reload
Load All (devtools)
Test Package (Desktop)
 Test Package (Web)
 Check Package
Document Package

| Windows /Linux | Mac |
|----------------|-------------|
| Ctrl+Shift+B | Cmd+Shift+B |
| Ctrl+Shift+L | Cmd+Shift+L |
| Ctrl+Shift+T | Cmd+Shift+T |
| Ctrl+Alt+F7 | Cmd+Opt+F7 |
| Ctrl+Shift+E | Cmd+Shift+E |
| Ctrl+Shift+D | Cmd+Shift+D |

8 DOCUMENTS AND APPS

Preview HTML (Markdown, etc.)
Knit Document (knitr)
 Compile Notebook
 Compile PDF (TeX and Sweave)
 Insert chunk (Sweave and Knitr)
 Insert code section
 Re-run previous region
 Run current document
Run from start to current line
Run the current code section
 Run previous Sweave/Rmd code
 Run the current chunk
 Run the next chunk
 Sync Editor & PDF Preview
 Previous plot
 Next plot
Show Keyboard Shortcuts

| Windows /Linux | Mac |
|---------------------|-----------------------|
| Ctrl+Shift+K | Cmd+Shift+K |
| Ctrl+Shift+K | Cmd+Shift+K |
| Ctrl+Shift+K | Cmd+Shift+K |
| Ctrl+Shift+K | Cmd+Shift+K |
| Ctrl+Shift+K | Cmd+Shift+K |
| Ctrl+Shift+R | Cmd+Shift+R |
| Ctrl+Shift+R | Cmd+Shift+R |
| Ctrl+Shift+P | Cmd+Shift+P |
| Ctrl+Alt+R | Cmd+Option+R |
| Ctrl+Alt+B | Cmd+Option+B |
| Ctrl+Alt+T | Cmd+Option+T |
| Ctrl+Alt+P | Cmd+Option+P |
| Ctrl+Alt+C | Cmd+Option+C |
| Ctrl+Alt+N | Cmd+Option+N |
| Ctrl+F8 | Cmd+F8 |
| Ctrl+Alt+F11 | Cmd+Option+F11 |
| Ctrl+Alt+F12 | Cmd+Option+F12 |
| Alt+Shift+K | Option+Shift+K |

Shiny :: CHEAT SHEET

Basics

A **Shiny** app is a web page (**UI**) connected to a computer running a live R session (**Server**)



APP TEMPLATE

Begin writing a new app with this template. Preview the app by running the code at the R command line.

```
library(shiny)
ui <- fluidPage()
server <- function(input, output){}
shinyApp(ui = ui, server = server)
```

- **ui** - nested R functions that assemble an HTML user interface for your app
- **server** - a function with instructions on how to build and rebuild the R objects displayed in the UI
- **shinyApp** - combines ui and server into an app. Wrap with `runApp()` if calling from a sourced script or inside a function.

SHARE YOUR APP

 The easiest way to share your app is to host it on shinyapps.io, a cloud based service from RStudio

1. Create a free or professional account at <http://shinyapps.io>
2. Click the **Publish** icon in the RStudio IDE or run:
`rsconnect::deployApp("<path to directory>")`

Build or purchase your own Shiny Server
www.rstudio.com/products/shiny-server/



Building an App

Complete the template by adding arguments to `fluidPage()` and a body to the `server` function.

Add inputs to the UI with `*Input()` functions

Add outputs with `*Output()` functions

Tell server how to render outputs with R in the server function. To do this:

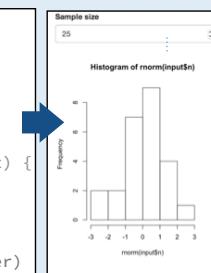
1. Refer to outputs with `output$<id>`
2. Refer to inputs with `input$<id>`
3. Wrap code in a `render*()` function before saving to output

Save your template as `app.R`. Alternatively, split your template into two files named `ui.R` and `server.R`.

```
library(shiny)
ui <- fluidPage(
  numericInput(inputId = "n",
    "Sample size", value = 25),
  plotOutput(outputId = "hist")
)
server <- function(input, output) {
  output$hist <- renderPlot({
    hist(rnorm(input$n))
  })
}
shinyApp(ui = ui, server = server)
```

```
# ui.R
fluidPage(
  numericInput(inputId = "n",
    "Sample size", value = 25),
  plotOutput(outputId = "hist")
)

# server.R
function(input, output) {
  output$hist <- renderPlot({
    hist(rnorm(input$n))
  })
}
```



`ui.R` contains everything you would save to ui.

`server.R` ends with the function you would save to server.

No need to call `shinyApp()`.

Save each app as a directory that holds an `app.R` file (or a `server.R` file and a `ui.R` file) plus optional extra files.



- The directory name is the name of the app
- (optional) defines objects available to both ui.R and server.R
- (optional) used in showcase mode
- (optional) data, scripts, etc.
- (optional) directory of files to share with web browsers (images, CSS, .js, etc.) Must be named "`www`"

Launch apps with
`runApp(<path to directory>)`

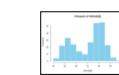
Outputs - `render*()` and `*Output()` functions work together to add R output to the UI



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

works with

`dataTableOutput(outputId, icon, ...)`



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



`renderPlot(expr, width, height, res, ..., env, quoted, func)`



`renderPrint(expr, env, quoted, func, width)`



`renderTable(expr, ..., env, quoted, func)`



`renderText(expr, env, quoted, func)`



`renderUI(expr, env, quoted, func)`

`imageOutput(outputId, width, height, click, dblclick, hover, hoverDelay, inline, hoverDelayType, brush, clickId, hoverId)`

`plotOutput(outputId, width, height, click, dblclick, hover, hoverDelay, inline, hoverDelayType, brush, clickId, hoverId)`

`verbatimTextOutput(outputId)`

`tableOutput(outputId)`

`textOutput(outputId, container, inline)`

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

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



Inputs

collect values from the user

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

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

`ActionLink(inputId, label, icon, ...)`

`checkboxGroupInput(inputId, label, choices, selected, inline)`

`checkboxInput(inputId, label, value)`

`dateInput(inputId, label, value, min, max, format, startview, weekstart, language)`

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

`fileInput(inputId, label, multiple, accept)`

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

`passwordInput(inputId, label, value)`

`radioButtons(inputId, label, choices, selected, inline)`

`selectInput(inputId, label, choices, selected, multiple, selectize, width, size) (also selectizeInput())`

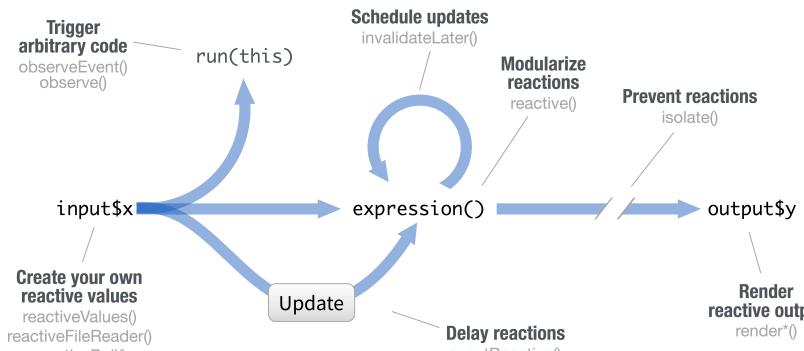
`sliderInput(inputId, label, min, max, value, step, round, format, locale, ticks, animate, width, sep, pre, post)`

`submitButton(text, icon) (Prevents reactions across entire app)`

`textInput(inputId, label, value)`

Reactivity

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



CREATE YOUR OWN REACTIVE VALUES

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

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

***Input()** functions
(see front page)
reactiveValues(...)

Each input function creates a reactive value stored as **input\$<inputId>**
reactiveValues() creates a list of reactive values whose values you can set.

RENDERS REACTIVE OUTPUT

```
library(shiny)
ui <- fluidPage(
 textInput("a","","A"),
  textOutput("b"))

server <-
function(input,output){
  output$b <-
    renderText({
      input$a
    })
}

shinyApp(ui, server)
```

render*() functions
(see front page)

Builds an object to display. Will rerun code in body to rebuild the object whenever a reactive value in the code changes.

Save the results to **output\$<outputId>**

PREVENT REACTIONS

```
library(shiny)
ui <- fluidPage(
  textInput("a","","A"),
  textOutput("b"))

server <-
function(input,output){
  output$b <-
    renderText({
      isolate({input$a})
    })
}

shinyApp(ui, server)
```

isolate(expr)

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

TRIGGER ARBITRARY CODE

```
library(shiny)
ui <- fluidPage(
  textInput("a","","A"),
  actionButton("go","Go"))

server <-
function(input,output){
  observeEvent(input$go, {
    print(input$a)
  })
}

shinyApp(ui, server)
```

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

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

MODULARIZE REACTIONS

```
ui <- fluidPage(
  textInput("a","","A"),
  textInput("z","","Z"),
  textOutput("b"))

server <-
function(input,output){
  re <- reactive({
    paste(input$a,input$z)
  })
  output$b <- renderText({
    re()
  })
}

shinyApp(ui, server)
```

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

Creates a **reactive expression** that
• **caches** its value to reduce computation
• can be called by other code
• notifies its dependencies when it has been invalidated
Call the expression with function syntax, e.g. **re()**

DELAY REACTIONS

```
library(shiny)
ui <- fluidPage(
  textInput("a","","A"),
  actionButton("go","Go"),
  textOutput("b"))

server <-
function(input,output){
  re <- eventReactive(
    input$go, {input$a}
  )
  output$b <- renderText({
    re()
  })
}

shinyApp(ui, server)
```

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

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

UI

An app's UI is an HTML document.

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

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

Returns HTML



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

| | | | | |
|------------------|-------------------|---------------|----------------|---------------|
| tags\$a | tags\$data | tags\$h6 | tags\$nav | tags\$span |
| tags\$abbr | tags\$datalist | tags\$head | tags\$noscript | tags\$strong |
| tags\$address | tags\$dd | tags\$header | tags\$object | tags\$style |
| tags\$area | tags\$del | tags\$hgroup | tags\$script | tags\$sub |
| tags\$article | tags\$details | tags\$hr | tags\$optgroup | tags\$summary |
| tags\$aside | tags\$dfn | tags\$HTML | tags\$option | tags\$sup |
| tags\$audio | tags\$div | tags\$iframe | tags\$output | tags\$table |
| tags\$b | tags\$dt | tags\$img | tags\$param | tags\$tbody |
| tags\$base | tags\$em | tags\$input | tags\$pre | tags\$thead |
| tags\$bdo | tags\$embed | tags\$ins | tags\$progress | tags\$tfoot |
| tags\$blockquote | tags\$eventsource | tags\$kbd | tags\$srq | tags\$th |
| tags\$body | tags\$fieldset | tags\$keygen | tags\$ruby | tags\$thead |
| tags\$button | tags\$figcaption | tags\$label | tags\$rp | tags\$title |
| tags\$canvas | tags\$footer | tags\$legend | tags\$rt | tags\$track |
| tags\$caption | tags\$form | tags\$li | tags\$ss | tags\$u |
| tags\$cite | tags\$mark | tags\$link | tags\$small | tags\$var |
| tags\$code | tags\$h2 | tags\$map | tags\$meta | tags\$video |
| tags\$col | tags\$h3 | tags\$menu | tags\$source | tags\$wbr |
| tags\$colgroup | tags\$h4 | tags\$select | | |
| tags\$command | tags\$h5 | tags\$section | | |
| | | tags\$small | | |
| | | tags\$source | | |

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

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

Header 1

bold
italic
code
link
Raw html

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

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

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

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

To include an image
1. Place the file in the **www** subdirectory
2. Link to it with **img(src=<file name>)**

Layouts

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

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

absolutePanel()
conditionalPanel()
fixedPanel()
headerPanel()
inputPanel()
mainPanel()

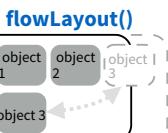
navlistPanel()
sidebarPanel()
tabPanel()
tabsetPanel()
titlePanel()
wellPanel()

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

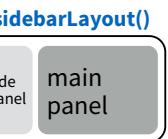
fluidRow()



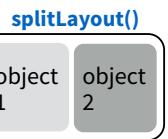
```
ui <- fluidPage(
  fluidRow(column(width = 4),
    column(width = 2, offset = 3)),
  fluidRow(column(width = 12))
)
```



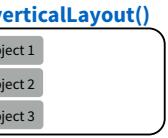
```
ui <- fluidPage(
  flowLayout(# object 1,
            # object 2,
            # object 3
  )
)
```



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



```
ui <- fluidPage(
  splitLayout(# object 1,
             # object 2
  )
)
```



```
ui <- fluidPage(
  verticalLayout(# object 1,
                 # object 2,
                 # object 3
  )
)
```



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

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

tab 1
tab 2
tab 3
contents

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

tab 1
tab 2
tab 3
contents

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

Page
tab 1
tab 2
tab 3
contents

Data Science in Spark with sparklyr :: CHEAT SHEET

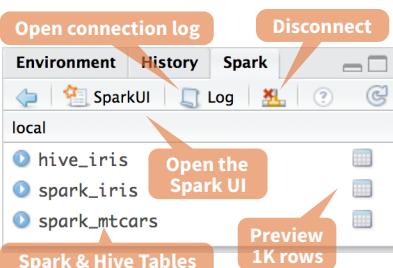


Intro

sparklyr is an R interface for Apache Spark™, it provides a complete **dplyr** backend and the option to query directly using **Spark SQL** statement. With sparklyr, you can orchestrate distributed machine learning using either **Spark's MLLib** or **H2O** Sparkling Water.

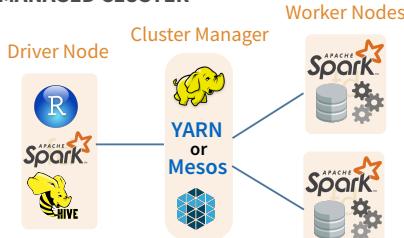
Starting with **version 1.044**, **RStudio Desktop, Server and Pro** include integrated support for the **sparklyr** package. You can create and manage connections to Spark clusters and local Spark instances from inside the IDE.

RStudio Integrates with sparklyr

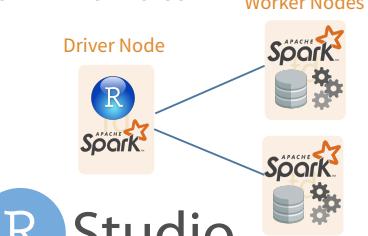


Cluster Deployment

MANAGED CLUSTER

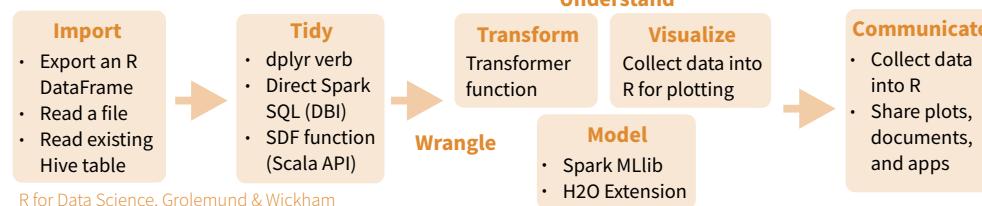


STAND ALONE CLUSTER



R Studio

Data Science Toolchain with Spark + sparklyr



Getting Started

LOCAL MODE (No cluster required)

1. Install a local version of Spark:
`spark_install("2.0.1")`
2. Open a connection
`sc <- spark_connect(master = "local")`

ON A MESOS MANAGED CLUSTER

1. Install RStudio Server or Pro on one of the existing nodes
2. Locate path to the cluster's Spark Home Directory, it normally is `"/usr/lib/spark"`
3. Open a connection
`spark_connect(master = "[mesos URL]", version = "1.6.2", spark_home = [Cluster's Spark path])`

USING LIVY (Experimental)

1. The Livy REST application should be running on the cluster
2. Connect to the cluster
`sc <- spark_connect(method = "livy", master = "http://host:port")`

Tuning Spark

EXAMPLE CONFIGURATION

```
config <- spark_config()
config$spark.executor.cores <- 2
config$spark.executor.memory <- "4G"
sc <- spark_connect(master = "yarn-client",
  config = config, version = "2.0.1")
```

IMPORTANT TUNING PARAMETERS with defaults

- `spark.yarn.am.cores`
- `spark.yarn.am.memory 512m`
- `spark.network.timeout 120s`
- `spark.executor.memory 1g`
- `spark.executor.cores 1`
- `spark.executor.instances`
- `spark.executor.extraJavaOptions`
- `spark.executor.heartbeatInterval 10s`
- `sparklyr.shell.executor-memory`
- `sparklyr.shell.driver-memory`

Using sparklyr

A brief example of a data analysis using Apache Spark, R and sparklyr in local mode

```
library(sparklyr); library(dplyr); library(ggplot2);
library(tidyr);
set.seed(100)
```

Install Spark locally

```
spark_install("2.0.1")
```

Connect to local version

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

```
import_iris <- copy_to(sc, iris, "spark_iris",
  overwrite = TRUE)
```

Copy data to Spark memory

```
partition_iris <- sdf_partition(
  import_iris, training=0.5, testing=0.5)
```

Partition data

```
sdf_register(partition_iris,
c("spark_iris_training", "spark_iris_test"))
```

Create a hive metadata for each partition

```
tidy_iris <-tbl(sc, "spark_iris_training") %>%
  select(Species, Petal_Length, Petal_Width)
```

Spark ML Decision Tree Model

```
model_iris <- tidy_iris %>%
  ml_decision_tree(response = "Species",
  features = c("Petal_Length", "Petal_Width"))
```

Create reference to Spark table

```
test_iris <-tbl(sc, "spark_iris_test")
```

```
pred_iris <- sdf_predict(
  model_iris, test_iris) %>%
  collect
```

Bring data back into R memory for plotting

```
pred_iris %>%
  inner_join(data.frame(prediction = 0:2,
  lab = model_iris$model.parameters$labels)) %>%
  ggplot(aes(Petal_Length, Petal_Width, col = lab)) +
  geom_point()
```

```
spark_disconnect(sc)
```

Disconnect

Reactivity

COPY A DATA FRAME INTO SPARK

```
sdf_copy_to(sc, iris, "spark_iris")
```

```
sdf_copy_to(sc, x, name, memory, repartition, overwrite)
```

IMPORT INTO SPARK FROM A FILE

Arguments that apply to all functions:

```
sc, name, path, options = list(), repartition = 0, memory = TRUE, overwrite = TRUE
```

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

JSON `spark_read_json()`

PARQUET `spark_read_parquet()`

SPARK SQL COMMANDS

```
DBI::dbWriteTable(sc, "spark_iris", iris)
```

```
DBI::dbWriteTable(conn, name, value)
```

FROM A TABLE IN HIVE

```
my_var <- tbl_cache(sc, name = "hive_iris")
```

`tbl_cache(sc, name, force = TRUE)`
Loads the table into memory

```
my_var <- dplyr::tbl(sc, name = "hive_iris")
```

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

Wrangle

SPARK SQL VIA DPLYR VERBS

Translates into Spark SQL statements

```
my_table <- my_var %>%>%  
  filter(Species == "setosa") %>%>%  
  sample_n(10)
```

DIRECT SPARK SQL COMMANDS

```
my_table <- DBI::dbGetQuery(sc, "SELECT *  
  FROM iris LIMIT 10")
```

```
DBI::dbGetQuery(conn, statement)
```

SCALA API VIA SDF FUNCTIONS

```
sdf_mutate(.data)
```

Works like dplyr mutate function

```
sdf_partition(x, ..., weights = NULL, seed =  
  sample(.Machine$integer.max, 1))
```

```
sdf_partition(x, training = 0.5, test = 0.5)
```

```
sdf_register(x, name = NULL)
```

Gives a Spark DataFrame a table name

```
sdf_sample(x, fraction = 1, replacement =  
  TRUE, seed = NULL)
```

```
sdf_sort(x, columns)
```

Sorts by >=1 columns in ascending order

```
sdf_with_unique_id(x, id = "id")
```

```
sdf_predict(object, newdata)
```

Spark DataFrame with predicted values

ML TRANSFORMERS

```
ft_binarizer(my_table, input.col = "Petal_Length", output.col = "petal_large", threshold = 1.2)
```

Arguments that apply to all functions:
`x, input.col = NULL, output.col = NULL`

```
ft_binarizer(threshold = 0.5)
```

Assigned values based on threshold

```
ft_bucketizer(splits)
```

Numeric column to discretized column

```
ft_discrete_cosine_transform(inverse = FALSE)
```

Time domain to frequency domain

```
ft_elementwise_product(scaling.col)
```

Element-wise product between 2 cols

```
ft_index_to_string()
```

Index labels back to label as strings

```
ft_one_hot_encoder()
```

Continuous to binary vectors

```
ft_quantile_discretizer(n.buckets = 5L)
```

Continuous to binned categorical values

```
ft_sql_transformer(sql)
```

```
ft_string_indexer(params = NULL)
```

Column of labels into a column of label indices.

```
ft_vectorAssembler()
```

Combine vectors into single row-vector

Visualize & Communicate

DOWNLOAD DATA TO R MEMORY

```
r_table <- collect(my_table)
```

```
plot(Petal_Width ~ Petal_Length, data = r_table)
```

```
dplyr::collect(x)
```

Download a Spark DataFrame to an R DataFrame

```
sdf_read_column(x, column)
```

Returns contents of a single column to R

SAVE FROM SPARK TO FILE SYSTEM

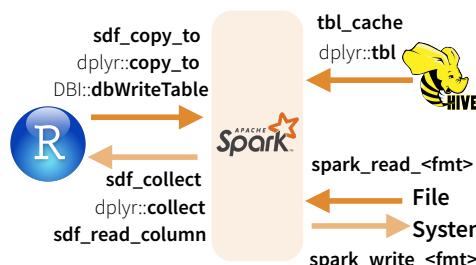
Arguments that apply to all functions: `x, path`

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

JSON `spark_read_json(mode = NULL)`

PARQUET `spark_read_parquet(mode = NULL)`

Reading & Writing from Apache Spark



Extensions

Create an R package that calls the full Spark API & provide interfaces to Spark packages.

CORE TYPES

```
spark_connection()
```

Connection between R and the Spark shell process

```
spark_jobobj()
```

Instance of a remote Spark object

```
spark_dataframe()
```

Instance of a remote Spark DataFrame object

CALL SPARK FROM R

```
invoke()
```

Call a method on a Java object

```
invoke_new()
```

Create a new object by invoking a constructor

```
invoke_static()
```

Call a static method on an object

MACHINE LEARNING EXTENSIONS

```
ml_create_dummy_variables()
```

```
ml_options()
```

```
ml_prepare_dataframe()
```

```
ml_model()
```

```
ml_prepare_response_features_intercept()
```

Model (MLlib)

```
ml_decision_tree(my_table, response = "Species", features = c("Petal_Length", "Petal_Width"))
```

```
ml_als_factorization(x, user.column = "user", rating.column = "rating", item.column = "item", rank = 10L, regularization.parameter = 0.1, iter.max = 10L, ml.options = ml_options())
```

```
ml_decision_tree(x, response, features, max.bins = 32L, max.depth = 5L, type = c("auto", "regression", "classification"), ml.options = ml_options())
```

```
ml_generalized_linear_regression(x, response, features, intercept = TRUE, family = gaussian(link = "identity"), iter.max = 100L, ml.options = ml_options())
```

```
ml_kmeans(x, centers, iter.max = 100, features = dplyr::tbl_vars(x), compute.cost = TRUE, tolerance = 1e-04, ml.options = ml_options())
```

```
ml_lda(x, features = dplyr::tbl_vars(x), k = length(features), alpha = (50/k) + 1, beta = 0.1 + 1, ml.options = ml_options())
```

```
ml_linear_regression(x, response, features, intercept = TRUE, alpha = 0, lambda = 0, iter.max = 100L, ml.options = ml_options())
```

```
ml_multilayer_perceptron(x, response, features, layers, iter.max = 100, seed = sample(Machine$integer.max, 1), ml.options = ml_options())
```

```
ml_naive_bayes(x, response, features, lambda = 0, ml.options = ml_options())
```

```
ml_one_vs_rest(x, classifier, response, features, ml.options = ml_options())
```

```
ml_pca(x, features = dplyr::tbl_vars(x), ml.options = ml_options())
```

```
ml_random_forest(x, response, features, max.bins = 32L, max.depth = 5L, num.trees = 20L, type = c("auto", "regression", "classification"), ml.options = ml_options())
```

```
ml_survival_regression(x, response, features, intercept = TRUE, censor = "censor", iter.max = 100L, ml.options = ml_options())
```

```
ml_binary_classification_eval(predicted_tbl_spark, label, score, metric = "areaUnderROC")
```

```
ml_classification_eval(predicted_tbl_spark, label, predicted_lbl, metric = "f1")
```

```
ml_tree_feature_importance(sc, model)
```



sparklyr

is an R
interface
for

Apache
Spark

sparklyr

ML

Extensions

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.

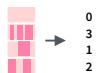
Detect Matches



`str_detect(string, pattern)` Detect the presence of a pattern match in a string.
`str_detect(fruit, "a")`



`str_which(string, pattern)` Find the indexes of strings that contain a pattern match.
`str_which(fruit, "a")`

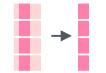


`str_count(string, pattern)` Count the number of matches in a string.
`str_count(fruit, "a")`

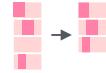


`str_locate(string, pattern)` Locate the positions of pattern matches in a string. Also `str_locate_all`.
`str_locate(fruit, "a")`

Subset Strings



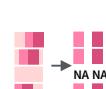
`str_sub(string, start = 1L, end = -1L)` Extract substrings from a character vector.
`str_sub(fruit, 1, 3); str_sub(fruit, -2)`



`str_subset(string, pattern)` Return only the strings that contain a pattern match.
`str_subset(fruit, "b")`



`str_extract(string, pattern)` Return the first pattern match found in each string, as a vector. Also `str_extract_all` to return every pattern match.
`str_extract(fruit, "[aeiou]")`



`str_match(string, pattern)` Return the first pattern match found in each string, as a matrix with a column for each () group in pattern. Also `str_match_all`.
`str_match(sentences, "(a|the) ([^]+)")`

Manage Lengths



`str_length(string)` The width of strings (i.e. number of code points, which generally equals the number of characters).
`str_length(fruit)`



`str_pad(string, width, side = c("left", "right", "both"), pad = " ")` Pad strings to constant width.
`str_pad(fruit, 17)`



`str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...")` Truncate the width of strings, replacing content with ellipsis.
`str_trunc(fruit, 3)`



`str_trim(string, side = c("both", "left", "right"))` Trim whitespace from the start and/or end of a string.
`str_trim(fruit)`

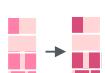
Mutate Strings



`str_sub() <- value`. Replace substrings by identifying the substrings with `str_sub()` and assigning into the results.
`str_sub(fruit, 1, 3) <- "str"`



`str_replace(string, pattern, replacement)` Replace the first matched pattern in each string.
`str_replace(fruit, "a", "-")`



`str_replace_all(string, pattern, replacement)` Replace all matched patterns in each string.
`str_replace_all(fruit, "a", "-")`

A STRING
↓
a string

a string
↓
A STRING

a string
↓
A String

Join and Split



`str_c(..., sep = "", collapse = NULL)` Join multiple strings into a single string.
`str_c(letters, LETTERS)`



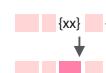
`str_c(..., sep = "", collapse = NULL)` Collapse a vector of strings into a single string.
`str_c(letters, collapse = "")`



`str_dup(string, times)` Repeat strings times times.
`str_dup(fruit, times = 2)`



`str_split_fixed(string, pattern, n)` Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also `str_split` to return a list of substrings.
`str_split_fixed(fruit, " ", n=2)`

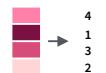


`str_glue(..., .sep = "", .envir = parent.frame())` Create a string from strings and {expressions} to evaluate.
`str_glue("Pi is {pi}")`



`str_glue_data(x, ..., .sep = "", .envir = parent.frame(), .na = "NA")` Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate.
`str_glue_data(mtcars, "rnames(mtcars)
has {hp} hp")`

Order Strings



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



`str_sort(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)`¹ Sort a character vector.
`str_sort(x)`

Helpers

`str_conv(string, encoding)` Override the encoding of a string.
`str_conv(fruit, "ISO-8859-1")`

apple
banana
pear

`str_view(string, pattern, match = NA)` View HTML rendering of first regex match in each string.
`str_view(fruit, "[aeiou]")`

apple
banana
pear

`str_view_all(string, pattern, match = NA)` View HTML rendering of all regex matches.
`str_view_all(fruit, "[aeiou]")`

`str_wrap(string, width = 80, indent = 0, exdent = 0)` Wrap strings into nicely formatted paragraphs.
`str_wrap(sentences, 20)`

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

Need to Know

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

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

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

| Special Character | Represents |
|-------------------|------------|
| \\\ | \ |
| \" | " |
| \n | new line |

Run ?"" to see a complete list

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

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

`writeLines("|\.")`

#|.

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

#| is a backslash

INTERPRETATION

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

`regexp(pattern, ignore_case = FALSE, multiline = FALSE, comments = FALSE, dotall = FALSE, ...)`

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

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

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

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



Regular Expressions -

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

MATCH CHARACTERS

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

see <- function(rx) str_view_all("abc ABC 123\ t.!?\ \ \n", rx)

example

| | |
|------------------|-----------------------|
| see("a") | abc ABC 123 .!?\ \ \n |
| see("\.") | abc ABC 123 .!?\ \ \n |
| see("\!") | abc ABC 123 .!?\ \ \n |
| see("\?") | abc ABC 123 .!?\ \ \n |
| see("\\\\") | abc ABC 123 .!?\ \ \n |
| see("\(\)") | abc ABC 123 .!?\ \ \n |
| see("\(\)") | abc ABC 123 .!?\ \ \n |
| see("\{\}") | abc ABC 123 .!?\ \ \n |
| see("\n") | abc ABC 123 .!?\ \ \n |
| see("\t") | abc ABC 123 .!?\ \ \n |
| see("\s") | abc ABC 123 .!?\ \ \n |
| see("\d") | abc ABC 123 .!?\ \ \n |
| see("\w") | abc ABC 123 .!?\ \ \n |
| see("\b") | abc ABC 123 .!?\ \ \n |
| see("[":digit:]) | abc ABC 123 .!?\ \ \n |
| see("[":alpha:]) | abc ABC 123 .!?\ \ \n |
| see("[":lower:]) | abc ABC 123 .!?\ \ \n |
| see("[":upper:]) | abc ABC 123 .!?\ \ \n |
| see("[":alnum:]) | abc ABC 123 .!?\ \ \n |
| see("[":punct:]) | abc ABC 123 .!?\ \ \n |
| see("[":graph:]) | abc ABC 123 .!?\ \ \n |
| see("[":space:]) | abc ABC 123 .!?\ \ \n |
| see("[":blank:]) | abc ABC 123 .!?\ \ \n |
| see(".") | abc ABC 123 .!?\ \ \n |

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

ALTERNATES

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

alt <- function(rx) str_view_all("abcde", rx)

ANCHORS

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

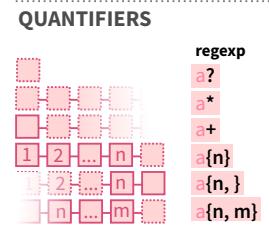
anchor <- function(rx) str_view_all("aaa", rx)

LOOK AROUNDS

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

look <- function(rx) str_view_all("bacad", rx)

QUANTIFIERS



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

quant <- function(rx) str_view_all(".aa.aaa", rx)

GRAPHS

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

| regexp | matches | example |
|---------|-----------------|---------|
| (ab d)e | sets precedence | abcde |

ref <- function(rx) str_view_all("abbaab", rx)

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

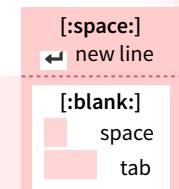
\\1

1 (etc.)

first () group, etc.

ref("(a)(b)\\2\\1")

abbaab



| [:space:] | new line |
|-----------|----------|
| [:blank:] | space |
| . | tab |

| [:graph:] | . |
|-----------|----|
| , | , |
| : | : |
| ; | ; |
| ? | ? |
| ! | ! |
| ? | ? |
| \ | \ |
| / | / |
| ` | ` |
| = | = |
| * | * |
| + | + |
| ^ | ^ |
| - | - |
| _ | _ |
| = | = |
| ~ | ~ |
| [| [|
|] |] |
| { | { |
| } | } |
| (| (|
|) |) |
| < | < |
| > | > |
| @ | @ |
| # | # |
| \$ | \$ |

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

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



Tidy evaluation with rlang :: CHEAT SHEET

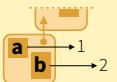


Vocabulary

Tidy Evaluation (Tidy Eval) is not a package, but a framework for doing non-standard evaluation (i.e. delayed evaluation) that makes it easier to program with tidyverse functions.

pi

Symbol - a name that represents a value or object stored in R. `is_symbol(expr(pi))`



Environment - a list-like object that binds symbols (names) to objects stored in memory. Each env contains a link to a second, **parent** env, which creates a chain, or search path, of environments. `is_environment(current_env())`

`rlang::caller_env(n = 1)` Returns calling env of the function it is in.

`rlang::child_env(.parent, ...)` Creates new env as child of .parent. Also **env**.

`rlang::current_env()` Returns execution env of the function it is in.

1

Constant - a bare value (i.e. an atomic vector of length 1). `is_bare_atomic(1)`

abs (1)

Call object - a vector of symbols/constants/calls that begins with a function name, possibly followed by arguments. `is_call(expr(abs(1)))`

pi — code
3.14 — result

Code - a sequence of symbols/constants/calls that will return a result if evaluated. Code can be:

1. Evaluated immediately (**Standard Eval**)
2. Quoted to use later (**Non-Standard Eval**)

`is_expression(expr(pi))`

e
a + b

Expression - an object that stores quoted code without evaluating it. `is_expression(expr(a + b))`

q
a + b, [a, b]

Quosure- an object that stores both quoted code (without evaluating it) and the code's environment. `is_quosure(quo(a + b))`

`[a, b]` `rlang::quo_get_env(quo)` Return the environment of a quosure.

`[a, b]` `rlang::quo_set_env(quo, expr)` Set the environment of a quosure.

`a + b` `rlang::quo_get_expr(quo)` Return the expression of a quosure.

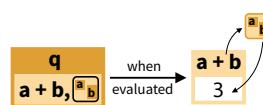
Expression Vector - a list of pieces of quoted code created by base R's `expression` and `parse` functions. Not to be confused with **expression**.



Quoting Code

Quote code in one of two ways (if in doubt use a quosure):

QUOSURES



Quosure- An expression that has been saved with an environment (aka a closure).

A quosure can be evaluated later in the stored environment to return a predictable result.

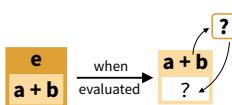
`rlang::quo(expr)` Quote contents as a quosure. Also **quos** to quote multiple expressions. `a <- 1; b <- 2; q <- quo(a + b); qs <- quos(a, b)`

`rlang::enquo(arg)` Call from within a function to quote what the user passed to an argument as a quosure. Also **enquos** for multiple args.
`quote_this <- function(x) enquo(x)`
`quote_these <- function(...) enquos(...)`

`rlang::new_quosure(expr, env = caller_env())` Build a quosure from a quoted expression and an environment.
`new_quosure(expr(a + b), current_env())`



EXPRESSION



Quoted Expression - An expression that has been saved by itself.

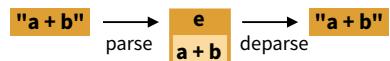
A quoted expression can be evaluated later to return a result that will depend on the environment it is evaluated in

`rlang::expr(expr)` Quote contents. Also **exprs** to quote multiple expressions. `a <- 1; b <- 2; e <- expr(a + b); es <- exprs(a, b, a + b)`

`rlang::enexpr(arg)` Call from within a function to quote what the user passed to an argument. Also **enexprs** to quote multiple arguments.
`quote_that <- function(x) enexpr(x)`
`quote_those <- function(...) enexprs(...)`

`rlang::ensym(x)` Call from within a function to quote what the user passed to an argument as a symbol, accepts strings. Also **ensyms**.
`quote_name <- function(name) ensym(name)`
`quote_names <- function(...) ensyms(...)`

Parsing and Deparsing



Parse - Convert a string to a saved expression.

...

`rlang::parse_expr(x)` Convert a string to an expression. Also **parse_exprs**, **sym**, **parse_quo**, **parse_quos**. `e <- parse_expr("a+b")`

Deparse - Convert a saved expression to a string.

...

`rlang::expr_text(expr, width = 60L, nlines = Inf)` Convert expr to a string. Also **quo_name**. `expr_text(e)`

Building Calls

`rlang::call2(fn, ..., .ns = NULL)` Create a call from a function and a list of args. Use **exec** to create and then evaluate the call. (See back page for **!!!**) `args <- list(x = 4, base = 2)`

`log(x = 4, base = 2)`

2

`call2("log", x = 4, base = 2)
 call2("log", !!!args)
 exec("log", x = 4, base = 2)
 exec("log", !!!args)`

Evaluation

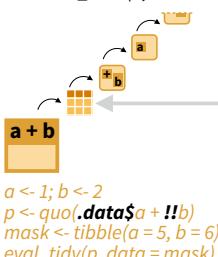
To evaluate an expression, R :

1. Looks up the symbols in the expression in the active environment (or a supplied one), followed by the environment's parents
2. Executes the calls in the expression

The result of an expression depends on which environment it is evaluated in.

QUOTED EXPRESSION

`rlang::eval_bare(expr, env = parent.frame())` Evaluate expr in env. `eval_bare(e, env = GlobalEnv)`



QUOSURES (and quoted exprs)

`rlang::eval_tidy(expr, data = NULL, env = caller_env())` Evaluate expr in env, using data as a **data mask**. Will evaluate quosures in their stored environment. `eval_tidy(q)`

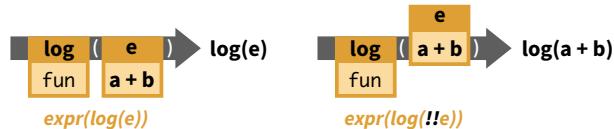
Data Mask - If data is non-NULL, `eval_tidy` inserts data into the search path before env, matching symbols to names in data.

Use the pronoun **.data\$** to force a symbol to be matched in data, and **!!** (see back) to force a symbol to be matched in the environments.

Quasiquotation (!! , !!!, :=)

QUOTATION

Storing an expression without evaluating it.
`e <- expr(a + b)`

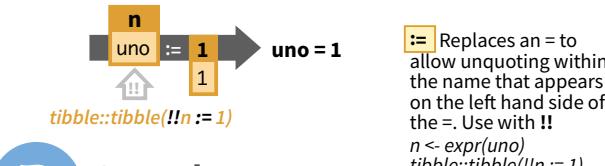
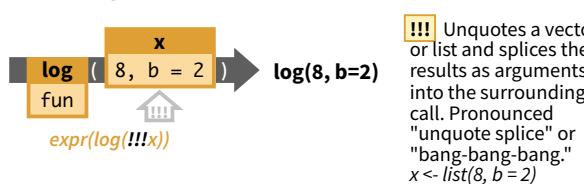
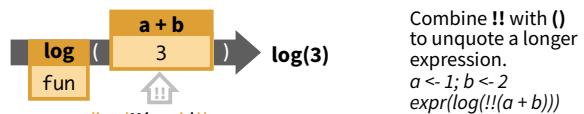
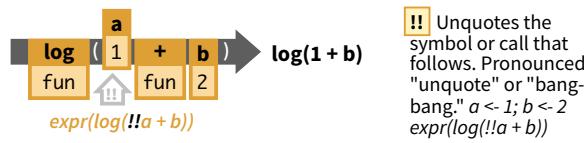


rlang provides !!, !!!, and := for doing quasiquotation.

!!, !!!, and := are not functions but syntax (symbols recognized by the functions they are passed to). Compare this to how

- . is used by `magrittr::%>%()`
- . is used by `stats::lm()`
- .x is used by `purrr::map()`, and so on.

!!, !!!, and := are only recognized by some rlang functions and functions that use those functions (such as tidyverse functions).



Programming Recipes

Quoting function- A function that quotes any of its arguments internally for delayed evaluation in a chosen environment. You must take **special steps to program safely** with a quoting function.

How to spot a quoting function?
A function quotes an argument if the argument returns an error when run on its own.

Many tidyverse functions are quoting functions: e.g. `filter`, `select`, `mutate`, `summarise`, etc.

```
dplyr::filter(cars, speed == 25)
          speed dist
          1     25    85
```

```
speed == 25
          Error!
```

PROGRAM WITH A QUOTING FUNCTION

```
data_mean <- function(data, var) {
  require(dplyr)
  var <- rlang::enquo(var) 1
  data %>%
    summarise(mean = mean (!!var)) 2
}
```

1. Capture user argument that will be quoted with `rlang::enquo`.
2. Unquote the user argument into the quoting function with !!!.

PASS MULTIPLE ARGUMENTS TO A QUOTING FUNCTION

```
group_mean <- function(data, var, ...) {
  require(dplyr)
  var <- rlang::enquo(var)
  group_vars <- rlang::enquos(...) 1
  data %>%
    group_by (!!group_vars) %>%
    summarise(mean = mean (!!var))
}
```

1. Capture user arguments that will be quoted with `rlang::enquos`.
2. Unquote splice the user arguments into the quoting function with !!!.

MODIFY USER ARGUMENTS

```
my_do <- function(f, v, df) {
  f <- rlang::enquo(f) 1
  v <- rlang::enquo(v)
  todo <- rlang::quo (!!f)(!!v) 2
  rlang::eval_tidy(todo, df) 3
}
```

1. Capture user arguments with `rlang::enquo`.
2. **Unquote** user arguments into a new expression or quoture to use
3. **Evaluate** the new expression/ quoture instead of the original argument

APPLY AN ARGUMENT TO A DATA FRAME

```
subset2 <- function(df, rows) {
  rows <- rlang::enquo(rows) 1
  vals <- rlang::eval_tidy(rows, data = df)
  df[vals, , drop = FALSE] 2
}
```

1. Capture user argument with `rlang::enquo`.
2. Evaluate the argument with `rlang::eval_tidy`. Pass the data frame to `data` to use as a data mask.
3. **Suggest** in your documentation that your users use the `.data` and `.env` pronouns.

WRITE A FUNCTION THAT RECOGNIZES QUASIQUOTATION (!! , !!!, :=)



1. Capture the quasiquotation-aware argument with `rlang::enquo`.
2. Evaluate the arg with `rlang::eval_tidy`.

```
add1 <- function(x) {
  q <- rlang::enquo(x)
  rlang::eval_tidy(q) + 1
}
```

PASS TO ARGUMENT NAMES OF A QUOTING FUNCTION

```
named_m <- function(data, var, name) {
  require(dplyr)
  var <- rlang::enquo(var)
  name <- rlang::ensym(name) 1
  data %>%
    summarise (!!name := mean (!!var)) 2
}
```

1. Capture user argument that will be quoted with `rlang::ensym`.
2. Unquote the name into the quoting function with !! and :=.

PASS CRAN CHECK

```
#' @importFrom rlang .data 1
mutate_y <- function(df) {
  dplyr::mutate(df, y = .data$a + 1) 2
}
```

Quoted arguments in tidyverse functions can trigger an **R CMD check** NOTE about undefined global variables. To avoid this:

1. Import `rlang::data` to your package, perhaps with the roxygen2 tag `@importFrom rlang .data`
2. Use the `.data$` pronoun in front of variable names in tidyverse functions

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 | | | Programming | | | While Loop | | |
|---------------------------|----------------------------------|-----------------------------|------------------------------|------------------------|----------|------------------------------------|----------------|--|
| Creating Vectors | | | For Loop | | | While Loop | | |
| c(2, 4, 6) | 2 4 6 | Join elements into a vector | for (variable in sequence){ | Do something | } | while (condition){ | Do something | } |
| 2:6 | 2 3 4 5 6 | An integer sequence | for (i in 1:4){ | j <- i + 10 | print(j) | while (i < 5){ | print(i) | i <- i + 1 |
| seq(2, 3, by=0.5) | 2.0 2.5 3.0 | A complex sequence | Example | | | Example | | |
| rep(1:2, times=3) | 1 2 1 2 1 2 | Repeat a vector | for (i in 1:4){ | j <- i + 10 | print(j) | while (i < 5){ | print(i) | i <- i + 1 |
| rep(1:2, each=3) | 1 1 1 2 2 2 | Repeat elements of a vector | | | | | | |
| Vector Functions | | | | | | | | |
| sort(x) | rev(x) | | if (condition){ | Do something | | functions_name <- function(var){ | Do something | |
| Return x sorted. | Return x reversed. | | } else { | Do something different | } | return(new_variable) | | |
| table(x) | unique(x) | | | | | | | |
| See counts of values. | See unique values. | | | | | | | |
| Selecting Vector Elements | | | | | | | | |
| By Position | | | If Statements | | | Functions | | |
| x[4] | The fourth element. | | if (condition){ | Do something | | function_name <- function(var){ | Do something | |
| x[-4] | All but the fourth. | | } else { | Do something different | } | return(new_variable) | | |
| 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. | | if (i > 3){ | print('Yes') | | square <- function(x){ | squared <- x*x | |
| x[x < 0] | All elements less than zero. | | } else { | print('No') | } | return(squared) | | |
| x[x %in% c(1, 2, 5)] | Elements in the set 1, 2, 5. | | | | | | | |
| Named Vectors | | | | | | | | |
| x['apple'] | Element with name 'apple'. | | Reading and Writing Data | | | Also see the readr package. | | |
| | | | Input | | | Output | | |
| | | | 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 |
| | | | a != b | Not equal | a < b | Less than | a <= b | Less than or equal to |
| | | | | | | | | is.na(a) |
| | | | | | | | | 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

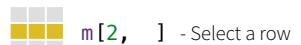
| | |
|-----------------|--|
| ls() | List all variables in the environment. |
| rm(x) | Remove x from the environment. |
| rm(list = ls()) | Remove all variables from the environment. |

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

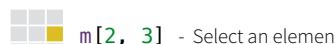
Matrices

```
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: $m^* x = n$

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,]`

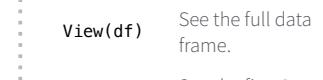


`df[2, 2]`

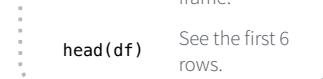
List subsetting



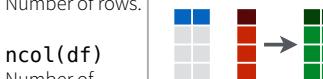
`df$x`



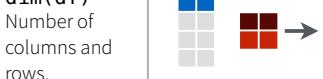
`df[[2]]`



`df$x`



`df$x`



`df$x`

Strings

Also see the `stringr` package.

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

Join multiple vectors together.

```
paste(x, collapse = ' ')
```

Join elements of a vector together.

```
grep(pattern, x)
```

Find regular expression matches in x.

```
gsub(pattern, replace, x)
```

Replace matches in x with a string.

```
toupper(x)
```

Convert to uppercase.

```
tolower(x)
```

Convert to lowercase.

```
nchar(x)
```

Number of characters in a string.

Factors

```
factor(x)
```

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

```
cut(x, breaks = 4)
```

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

Statistics

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

Linear model.

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

Generalised linear model.

```
summary
```

Get more detailed information out a model.

```
t.test(x, y)
```

Perform a t-test for difference between means.

```
pairwise.t.test
```

Perform a t-test for paired data.

```
prop.test
```

Test for a difference between proportions.

```
aov
```

Analysis of variance.

Distributions

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

Plotting

Also see the `ggplot2` package.



`plot(x)`

Values of x in order.



`plot(x, y)`

Values of x against y.



`hist(x)`

Histogram of x.

Dates

See the `lubridate` package.

caret Package

Cheat Sheet

Specifying the Model

Possible syntaxes for specifying the variables in the model:

```
train(y ~ x1 + x2, data = dat, ...)
train(x = predictor_df, y = outcome_vector, ...)
train(recipe_object, data = dat, ...)
```

- `rfe`, `sbf`, `gafs`, and `safs` only have the `x/y` interface.
- The `train` formula method will **always** create dummy variables.
- The `x/y` interface to `train` will not create dummy variables (but the underlying model function might).

Remember to:

- Have column names in your data.
- Use factors for a classification outcome (not 0/1 or integers).
- Have valid R names for class levels (not "0"/"1")
- Set the random number seed prior to calling `train` repeatedly to get the same resamples across calls.
- Use the `train` option `na.action = na.pass` if you will be imputing missing data. Also, use this option when predicting new data containing missing values.

To pass options to the underlying model function, you can pass them to `train` via the *ellipses*:

```
train(y ~ ., data = dat, method = "rf",
      # options to `randomForest`:
      importance = TRUE)
```

Parallel Processing

The `foreach` package is used to run models in parallel. The `train` code does not change but a "`do`" package must be called first.

```
# on MacOS or Linux      # on Windows
library(doMC)            library(doParallel)
registerDoMC(cores=4)    cl <- makeCluster(2)
registerDoParallel(cl)
```

The function `parallel::detectCores` can help too.

Preprocessing

Transformations, filters, and other operations can be applied to the *predictors* with the `preProc` option.

```
train(..., preProc = c("method1", "method2"), ...)
```

Methods include:

- `"center"`, `"scale"`, and `"range"` to normalize predictors.
- `"BoxCox"`, `"YeoJohnson"`, or `"expoTrans"` to transform predictors.
- `"knnImpute"`, `"bagImpute"`, or `"medianImpute"` to impute.
- `"corr"`, `"nzv"`, `"zv"`, and `"conditionalX"` to filter.
- `"pca"`, `"ica"`, or `"spatialSign"` to transform groups.

`train` determines the order of operations; the order that the methods are declared does not matter.

The `recipes` package has a more extensive list of preprocessing operations.

Adding Options

Many `train` options can be specified using the `trainControl` function:

```
train(y ~ ., data = dat, method = "cubist",
      trControl = trainControl(<options>))
```

Resampling Options

`trainControl` is used to choose a resampling method:

```
trainControl(method = <method>, <options>)
```

Methods and options are:

- `"cv"` for K-fold cross-validation (`number` sets the # folds).
- `"repeatedcv"` for repeated cross-validation (`repeats` for # repeats).
- `"boot"` for bootstrap (`number` sets the iterations).
- `"LGOCV"` for leave-group-out (`number` and `p` are options).
- `"L0O"` for leave-one-out cross-validation.
- `"oob"` for out-of-bag resampling (only for some models).
- `"timeslice"` for time-series data (options are `initialWindow`, `horizon`, `fixedWindow`, and `skip`).

Performance Metrics

To choose how to summarize a model, the `trainControl` function is used again.

```
trainControl(summaryFunction = <R function>,
             classProbs = <logical>)
```

Custom R functions can be used but `caret` includes several: `defaultSummary` (for accuracy, RMSE, etc), `twoClassSummary` (for ROC curves), and `prSummary` (for information retrieval). For the last two functions, the option `classProbs` must be set to `TRUE`.

Grid Search

To let `train` determine the values of the tuning parameter(s), the `tunelength` option controls how many values `per tuning` parameter to evaluate.

Alternatively, specific values of the tuning parameters can be declared using the `tuneGrid` argument:

```
grid <- expand.grid(alpha = c(0.1, 0.5, 0.9),
                      lambda = c(0.001, 0.01))
```

```
train(x = x, y = y, method = "glmnet",
      preProc = c("center", "scale"),
      tuneGrid = grid)
```

Random Search

For tuning, `train` can also generate random tuning parameter combinations over a wide range. `tuneLength` controls the total number of combinations to evaluate. To use random search:

```
trainControl(search = "random")
```

Subsampling

With a large class imbalance, `train` can subsample the data to balance the classes them prior to model fitting.

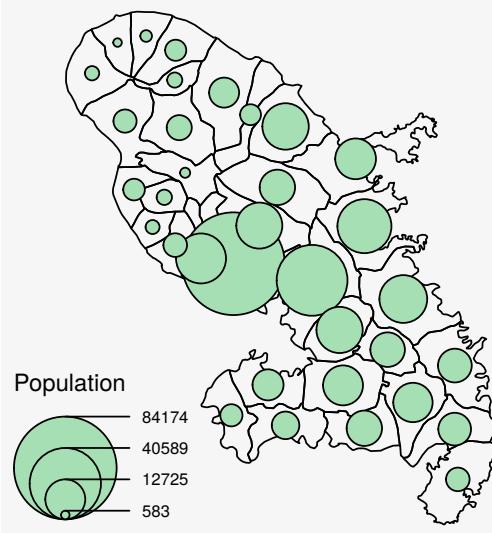
```
trainControl(sampling = "down")
```

Other values are `"up"`, `"smote"`, or `"rose"`. The latter two may require additional package installs.

Thematic maps with cartography :: CHEAT SHEET

Use cartography with spatial objects from sf or sp packages to create thematic maps.

```
library(cartography)
library(sf)
mtq <- st_read("martinique.shp")
plot(st_geometry(mtq))
propSymbolsLayer(x = mtq, var = "P13_POP",
  legend.title.txt = "Population",
  col = "#a7dfb4")
```



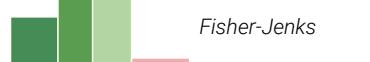
Classification

Available methods are: quantile, equal, q6, fisher-jenks, mean-sd, sd, geometric progression...

```
bks1 <- getBreaks(v = var, nclass = 6,
  method = "quantile")
bks2 <- getBreaks(v = var, nclass = 6,
  method = "fisher-jenks")
pal <- carto.pal("green.pal", 3, "wine.pal", 3)
hist(var, breaks = bks1, col = pal)
```



```
hist(var, breaks = bks2, col = pal)
```



Symbology

In most functions the x argument should be an sf object. sp objects are handled through spdf and df arguments.



Choropleth
choroLayer(x = mtq, var = "myvar",
method = "quantile", nclass = 8)



Typology
typoLayer(x = mtq, var = "myvar")



Proportional Symbols
propSymbolsLayer(x = mtq, var = "myvar",
inches = 0.1, symbols = "circle")



Colorized Proportional Symbols (relative data)
propSymbolsChoroLayer(x = mtq, var = "myvar",
var2 = "myvar2")



Colorized Proportional Symbols (qualitative data)
propSymbolsTypoLayer(x = mtq, var = "myvar",
var2 = "myvar2")



Double Proportional Symbols
propTrianglesLayer(x = mtq, var1 = "myvar",
var2 = "myvar2")



OpenStreetMap Basemap (see rosm package)
tiles <- getTiles(x = mtq, type = "osm")
tilesLayer(tiles)



Isopleth (see SpatialPosition package)
smoothLayer(x = mtq, var = "myvar",
typefcf = "exponential", span = 500,
beta = 2)



Discontinuities
discLayer(x = mtq.borders, df = mtq_df,
var = "myvar", threshold = 0.5)



Flows
propLinkLayer(x = mtq_link, df = mtq_df,
var = "fij")



Dot Density
dotDensityLayer(x = mtq, var = "myvar")



Labels
labelLayer(x = mtq, txt = "myvar",
halo = TRUE, overlap = FALSE)

Transformations

Polygons to Grid

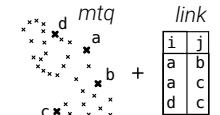
```
mtq_grid <- getGridLayer(x = mtq, cellsize = 3.6e+07,
  type = "hexagonal", var = "myvar")
```



Grids layers can be used by
choroLayer() or propSymbolsLayer().

Points to Links

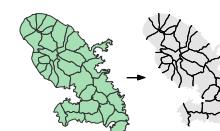
```
mtq_link <- getLinkLayer(x = mtq, df = link)
```



Links layers can be
used by *LinkLayer()

Polygons to Borders

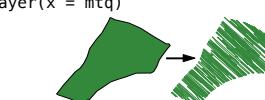
```
mtq_border <- getBorders(x = mtq)
```



Borders layers can be used by
discLayer() function

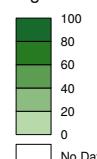
Polygons to Pencil Lines

```
mtq_pen <- getPencilLayer(x = mtq)
```



Legends

legendChoro()



```
legendChoro(pos = "topleft",
  title.txt = "legendChoro()",  
breaks = c(0,20,40,60,80,100),  
col = carto.pal("green.pal", 5),  
nodata = TRUE, nodata.txt = "No Data")
```

legendTypo()



```
legendTypo(title.txt = "legendTypo()",  
col = c("peru", "skyblue", "gray77"),  
categ = c("type 1", "type 2", "type 3"),  
nodata = FALSE)
```

legendCirclesSymbols()

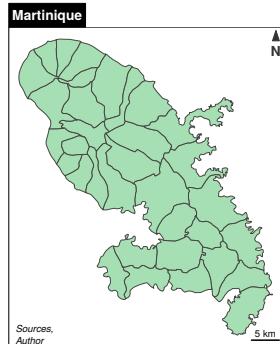


```
legendCirclesSymbols(var = c(10,100),  
title.txt = "legendCirclesSymbols()",  
col = "#a7dfb4ff", inches = 0.3)
```

See also legendSquaresSymbols(), legendBarsSymbols(),
legendGradLines(), legendPropLines() and legendPropTriangles().

Map Layout

North Arrow:
north(pos = "topright")



Scale Bar:
barscale(size = 5)

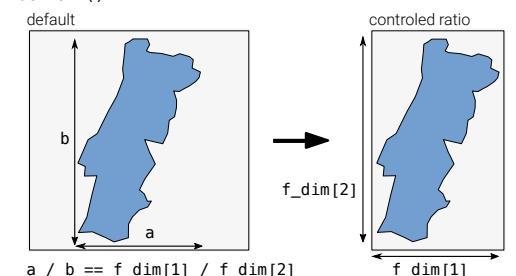
Full Layout:

```
layoutLayer(
  title = "Martinique",
  tabtitle = TRUE,
  frame = TRUE,
  author = "Author",
  sources = "Sources",
  north = TRUE,
  scale = 5)
```

Figure Dimensions

Get figure dimensions based on the dimension ratio of a spatial object, figure margins and output resolution.

```
f_dim <- getFigDim(x = sf_obj, width = 500,
  mar = c(0,0,0,0))
png("fig.png", width = 500, height = f_dim[2])
par(mar = c(0,0,0,0))
plot(sf_obj, col = "#729fcf")
dev.off()
```



Color Palettes

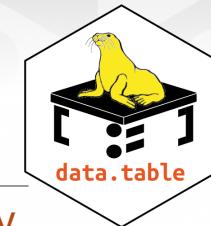
```
carto.pal(pal1 = "blue.pal", n1 = 5,
  pal2 = sand.pal, n2 = 3)
```



display.carto.all(n = 8)



Data Transformation with data.table :: CHEAT SHEET



Basics

data.table is an extremely fast and memory efficient package for transforming data in R. It works by converting R's native data frame objects into data.tables with new and enhanced functionality. The basics of working with data.tables are:

dt[i, j, by]

Take data.table **dt**,
subset rows using **i**
and manipulate columns with **j**,
grouped according to **by**.

data.tables are also data frames – functions that work with data frames therefore also work with data.tables.

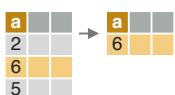
Create a data.table

data.table(a = c(1, 2), b = c("a", "b")) – create a data.table from scratch. Analogous to `data.frame()`.

setDT(df)* or as.data.table(df) – convert a data frame or a list to a data.table.

Subset rows using i

 **dt[1:2,]** – subset rows based on row numbers.

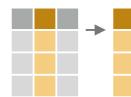
 **dt[a > 5,]** – subset rows based on values in one or more columns.

LOGICAL OPERATORS TO USE IN i

| | | | | | |
|---|----|----------|------|---|-----------|
| < | <= | is.na() | %in% | | %like% |
| > | >= | !is.na() | ! | & | %between% |

Manipulate columns with j

EXTRACT



dt[, c(2)] – extract columns by number. Prefix column numbers with “-” to drop.



dt[, .(b, c)] – extract columns by name.

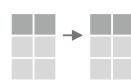
SUMMARIZE



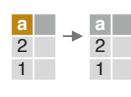
dt[, .(x = sum(a))] – create a data.table with new columns based on the summarized values of rows.

Summary functions like `mean()`, `median()`, `min()`, `max()`, etc. can be used to summarize rows.

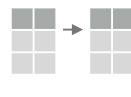
COMPUTE COLUMNS*



dt[, c := 1 + 2] – compute a column based on an expression.



dt[a == 1, c := 1 + 2] – compute a column based on an expression but only for a subset of rows.



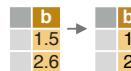
dt[, `:=` (c = 1, d = 2)] – compute multiple columns based on separate expressions.

DELETE COLUMN



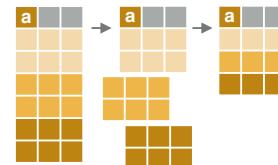
dt[, c := NULL] – delete a column.

CONVERT COLUMN TYPE



dt[, b := as.integer(b)] – convert the type of a column using `as.integer()`, `as.numeric()`, `as.character()`, `as.Date()`, etc..

Group according to by



dt[, j, by = .(a)] – group rows by values in specified columns.

dt[, j, keyby = .(a)] – group and simultaneously sort rows by values in specified columns.

COMMON GROUPED OPERATIONS

dt[, .(c = sum(b)), by = a] – summarize rows within groups.

dt[, c := sum(b), by = a] – create a new column and compute rows within groups.

dt[, .SD[1], by = a] – extract first row of groups.

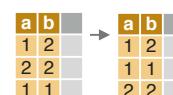
dt[, .SD[N], by = a] – extract last row of groups.

Chaining

dt[...][...] – perform a sequence of data.table operations by chaining multiple “[]”.

Functions for data.tables

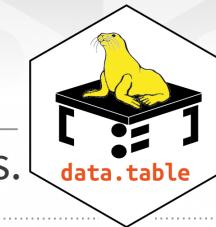
REORDER



setorder(dt, a, -b) – reorder a data.table according to specified columns. Prefix column names with “-” for descending order.

* SET FUNCTIONS AND :=

data.table's functions prefixed with “set” and the operator “:=” work without “-<” to alter data without making copies in memory. E.g., the more efficient “`setDT(df)`” is analogous to “`df <- as.data.table(df)`”.



Apply function to cols.

APPLY A FUNCTION TO MULTIPLE COLUMNS

| | | |
|-------|---|-------|
| a b | → | a b |
| 1 2 | | 1 2 |
| 2 2 | | 2 2 |

`dt[, lapply(.SD, mean), .SDcols = c("a", "b")]` – apply a function – e.g. `mean()`, `as.character()`, `which.max()` – to columns specified in `.SDcols` with `lapply()` and the `.SD` symbol. Also works with groups.

| | | |
|-------|---|-----------|
| a b | → | a a m |
| 1 4 | | 1 2 |
| 2 5 | | 2 2 |
| 3 6 | | 3 2 |

`cols <- c("a")`
`dt[, paste0(cols, "_m") := lapply(.SD, mean), .SDcols = cols]` – apply a function to specified columns and assign the result with suffixed variable names to the original data.

Sequential rows

ROW IDS

| | | |
|-------|---|-----------|
| a b | → | a b c |
| 1 a | | 1 a 1 |
| 2 a | | 2 a 2 |
| 3 b | | 3 b 1 |

`dt[, c := 1:N, by = b]` – within groups, compute a column with sequential row IDs.

LAG & LEAD

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

`dt[, c := shift(a, 1), by = b]` – within groups, duplicate a column with rows lagged by specified amount.

`dt[, c := shift(a, 1, type = "lead"), by = b]` – within groups, duplicate a column with rows leading by specified amount.

read & write files

IMPORT

`fread("file.csv")` – read data from a flat file such as `.csv` or `.tsv` into R.

`fread("file.csv", select = c("a", "b"))` – read specified columns from a flat file into R.

EXPORT

`fwrite(dt, "file.csv")` – write data to a flat file from R.

UNIQUE ROWS

| | | |
|-------|---|-------|
| a b | → | a b |
| 1 2 | | 1 2 |
| 2 2 | | 2 2 |

`unique(dt, by = c("a", "b"))` – extract unique rows based on columns specified in “by”. Leave out “by” to use all columns.

RENAME COLUMNS

| | | |
|-------|---|-------|
| a b | → | x y |
| 1 2 | | 1 2 |
| 2 2 | | 2 2 |

`setnames(dt, c("a", "b"), c("x", "y"))` – rename columns.

SET KEYS

`setkey(dt, a, b)` – set keys to enable fast repeated lookup in specified columns using `dt[.(value),]` or for merging without specifying merging columns using `dt_a[dt_b]`.

Combine data.tables

JOIN

| | | | | |
|-------|---|-------|---|-----------|
| a b | + | x y | = | a b x |
| 1 c | | 3 b | | 3 b 3 |
| 2 a | | 2 c | | 1 c 2 |
| 3 b | | 1 a | | 2 a 1 |

| | | | | |
|-----------|---|-----------|---|----------------|
| a b c | + | x y z | = | a b c x |
| 1 c 7 | | 3 b 4 | | 3 b 4 3 |
| 2 a 5 | | 2 c 5 | | 1 c 5 2 |
| 3 b 6 | | 1 a 8 | | NA a 8 1 |

`dt_a[dt_b, on = .(b = y, c > z)]` – join data.tables on rows with equal and unequal values.

ROLLING JOIN

| | | | | |
|--------------------|---|--------------------|---|------------------------|
| a id date | + | b id date | = | a id date b |
| 1 A 01-01-2010 | | 1 A 01-01-2013 | | 1 A 01-01-2013 1 |
| 2 A 01-01-2012 | | 1 B 01-01-2013 | | 2 B 01-01-2013 1 |
| 3 A 01-01-2014 | | | | |
| 1 B 01-01-2010 | | | | |
| 2 B 01-01-2012 | | | | |

`dt_a[dt_b, on = .(id = id, date = date), roll = TRUE]` – join data.tables on matching rows in id columns but only keep the most recent preceding match with the left data.table according to date columns. “roll = -Inf” reverses direction.

BIND

| | | | | |
|-------|---|-------|---|-------|
| a b | + | a b | = | a b |
| 1 2 | | 1 2 | | 1 2 |
| 2 2 | | 2 2 | | 2 2 |

`rbind(dt_a, dt_b)` – combine rows of two data.tables.

| | | | | |
|-------|---|-------|---|---------------|
| a b | + | x y | = | a b x y |
| 1 2 | | 1 2 | | 1 2 1 2 |
| 2 2 | | 2 2 | | 2 2 2 2 |

`cbind(dt_a, dt_b)` – combine columns of two data.tables.

RESHAPE TO WIDE FORMAT

| | | |
|----------------|---|------------------------------------|
| id y a b | → | id a x a z b x b z |
| A x 1 3 | | A 1 2 3 4 |
| A z 2 4 | | B 1 2 3 4 |
| B x 1 3 | | |
| B z 2 4 | | |

`dcast(dt, id ~ y, value.var = c("a", "b"))`

Reshape a data.table from long to wide format.

`dt` A data.table.
`id ~ y` Formula with a LHS: ID columns containing IDs for multiple entries. And a RHS: columns with values to spread in column headers.
`value.var` Columns containing values to fill into cells.

RESHAPE TO LONG FORMAT

| | | |
|------------------------------------|---|----------------|
| id a x a z b x b z | → | id y a b |
| A 1 2 3 4 | | A 1 3 |
| B 1 2 3 4 | | B 1 3 |
| A 2 2 4 | | A 2 4 |
| B 2 2 4 | | B 2 4 |

`melt(dt, id.vars = c("id"), measure.vars = patterns("^a", "^b"), variable.name = "y", value.name = c("a", "b"))`

Reshape a data.table from wide to long format.

`dt` A data.table.
`id.vars` ID columns with IDs for multiple entries.
`measure.vars` Columns containing values to fill into cells (often in pattern form).
`variable.name` Names of new columns for variables and values derived from old headers.
`value.name` Names of new columns for variables and values derived from old headers.

DeclareDesign:: CHEAT SHEET

Model

What is your model of the world, including how outcomes respond to interventions in the world?

Population

Define the size of the population, hierarchical structure (if any), and background variables.

Simple dataset with no background variables

```
pop <- declare_population(N = 100)
pop()
```

Simple dataset with background variables

```
declare_population(N = 100,
                   X = rnorm(N))
```

Two-level dataset

```
declare_population(
  schools =
    add_level(N = 10,
              funding = rnorm(N)),
  students =
    add_level(N = 100,
              scores = rnorm(N))
)
```

Outcomes

Outcomes that depend on a treatment (Z)

Using a formula

```
declare_potential_outcomes(
  Y ~ .5 * Z + rnorm(N))
```

As separate variables

```
declare_potential_outcomes(
  Y_Z_0 = rnorm(N),
  Y_Z_1 = Y_Z_0 + .5)
```

Outcomes that do not depend on treatment

```
declare_potential_outcomes(
  Y = rnorm(N))
```

Inquiry

What is the research question you want to answer?

Causal inquiries

```
declare_estimand(
  ATE = mean(Y_Z_1 - Y_Z_0))
```

Descriptive inquiries

```
declare_estimand(
  Y_median = median(Y))
```

Conditional estimands

```
declare_estimand(
  LATE = mean(Y_Z_1 - Y_Z_0),
  subset = complier == TRUE)
```

Data Strategy

How will you generate data to answer your inquiry?

Sampling

```
declare_sampling(n = 100)
```

```
declare_sampling(
  strata_n = 20,
  strata = urban_area)
```

Treatment assignment

```
declare_assignment(m = 100)
```

```
declare_assignment(
  clusters = villages,
  m = 10)
```

Answer Strategy

How will you generate an answer to your inquiry?

OLS with robust standard errors

```
declare_estimator(
  Y ~ Z, model = lm_robust)
```

2SLS instrumental variables regression with robust SEs

```
declare_estimator(
  Y ~ D | Z, model = iv_robust)
```

Difference-in-means

```
declare_estimator(
  Y ~ Z,
  model = difference_in_means)
```

DeclareDesign is a software implementation of the MIDA framework, according to which research designs have a **Model** of the world, an **Inquiry** about that model, a **Data strategy** that generates information about the world, and an **Answer** strategy that uses data to make a guess about the **Inquiry**. Declared designs can be “diagnosed” to calculate the properties of the design such as power and bias using Monte Carlo simulation.

All `declare_*` functions return *functions*. Most functions take a `data.frame` and return a `data.frame`.

Design Declaration

Put together all the steps into a declared design using the `+` operator

```
design <-
  declare_population(N = 200, X = rnorm(N)) +
  declare_potential_outcomes(Y ~ .5 * Z + X) +
  declare_estimand(ATE = mean(Y_Z_1 - Y_Z_0)) +
  declare_sampling(n = 100) +
  declare_assignment(m = 50) +
  declare_estimator(Y ~ Z, model = lm_robust)
```

```
draw_data(design)
draw_estimates(design)
get_estimates(design, data = real_data)
draw_estimands(design)
run_design(design)
summary(design)
compare_designs(design_1, design_2)
```

Design Diagnosis

Diagnose the properties of your design

```
diagnosis <- diagnose_design(
  design, sims = 100, bootstrap_sims = 100)
```

```
summary(diagnosis)
get_diagnosands(diagnosis)
get_simulations(diagnosis)
```

Custom diagnosands

```
diagnose_design(
  design,
  diagnosands = declare_diagnosands(
    sig_pos = mean(p.value < .05 & estimate > 0)))
```

estimatr:: CHEAT SHEET

OLS with lm_robust()

lm_robust() is lm() with robust SEs. HC2 is the default.

```
lm_robust(mpg ~ hp, data = mtcars)
lm_robust(mpg ~ hp, se_type = "HC1",
          data = mtcars)
lm_robust(mpg ~ hp, se_type = "classical",
          data = mtcars)
```

Indicate clusters to get clustered SEs. CR2 is the default.

```
lm_robust(mpg ~ hp, clusters = carb,
          data = mtcars)
lm_robust(mpg ~ hp, clusters = carb,
          se_type = "stata", data = mtcars)
```

Fixed effects two ways:

```
# FEs as "dummies"
lm_robust(mpg ~ hp + as.factor(am),
          data = mtcars)

# "Absorbing" FEs (substantially faster)
lm_robust(mpg ~ hp,
          fixed_effects = ~ am,
          data = mtcars)
```

post-estimation commands:

```
fit <- lm_robust(mpg ~ hp, data = mtcars)
summary(fit)
print(fit)
tidy(fit)
vcov(fit)
confint(fit)
nobs(fit)
predict(fit, newdata = mtcars)
```

estimatr is part of the DeclareDesign suite of packages for designing, implementing, and analyzing social science research designs.

2SLS with iv_robust()

iv_robust() is AER::ivreg() with robust SEs.

```
iv_robust(mpg ~ hp | am, data = mtcars)
iv_robust(mpg ~ hp | am,
          clusters = carb, data = mtcars)
```

Two-group estimators

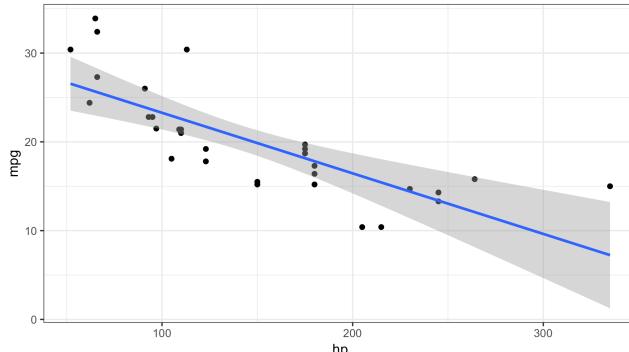
difference_in_means() and horvitz_thompson()
compare two groups

```
difference_in_means(mpg ~ am, data = mtcars)
horvitz_thompson(mpg ~ am, data = mtcars)
```

ggplot2 integration

Use robust variance estimates for drawing confidence intervals:

```
library(ggplot2)
ggplot(mtcars, aes(mpg, hp)) +
  geom_point() +
  stat_smooth(method = "lm_robust") +
  theme_bw()
```



Multiple models

Same outcome, different subsets:

```
library(tidyverse)
mtcars %>%
  split(~cyl) %>%
  map(~lm_robust(mpg ~ hp, data = .)) %>%
  map(tidy) %>%
  bind_rows(.id = "cyl")
```

Different outcomes, same subset:

```
c("mpg", "disp") %>%
  map(~formula(paste0(., " ~ hp"))) %>%
  map(~lm_robust(., data = mtcars)) %>%
  map(tidy) %>%
  bind_rows
```

Extras

```
# Lin (2013) covariate adjustment
lm_lin(mpg ~ am, covariates = ~ hp,
        data = mtcars)
```

```
# regression tables with texreg
fit <- lm_robust(mpg ~ hp, data = mtcars)
texreg::texreg(fit, include.ci = FALSE)
```

estimatr-to-Stata dictionary

estimatr

```
lm_robust(y ~ z,
           data = dat)
```

Stata

```
reg y z, vce(hc2)
```

```
lm_robust(y ~ z,
           clusters = cl,
           se_type = "stata",
           data = dat)
```

```
reg y z, vce(cluster cl)
```

```
lm_robust(mpg ~ hp,
           fixed_effects = ~ am,
           se_type = "stata",
           data = mtcars)
```

```
areg mpg hp, absorb(am)
vce(robust)
```

```
iv_robust(mpg ~ hp | am,
           se_type = "HC1",
           data = mtcars)
```

```
ivregress 2sls mpg (hp =
am), vce(robust) small
```

The eurostat package

R tools to access open data from Eurostat database

Search and download

Data in the Eurostat database is stored in tables. Each table has an identifier, a short table_code, and a description (e.g. tsdtr420 - People killed in road accidents).

Key eurostat functions allow to find the table_code, download the eurostat table and polish labels in the table.

Find the table code

The `search_eurostat(pattern,...)` function scans the directory of Eurostat tables and returns codes and descriptions of tables that match pattern.

```
library("eurostat")
query <- search_eurostat("road", type = "table")
query[1:3,1:2]
##          title      code
## 1   Goods transport by road ttr00005
## 2  People killed in road accidents tsdtr420
## 3 Enterprises with broadband access tin00090
```

Download the table

The `get_eurostat(id, time_format = "date", filters = "none", type = "code", cache = TRUE, ...)` function downloads the requested table from the Eurostat bulk download facility or from The Eurostat Web Services JSON API (if `filters` are defined). Downloaded data is cached (if `cache=TRUE`). Additional arguments define how to read the time column (`time_format`) and if table dimensions shall be kept as codes or converted to labels (`type`).

```
dat <- get_eurostat(id="tsdtr420", time_format="num")
head(dat)
##    unit sex geo time values
## 1    NR   T AT 1999 1079
## 2    NR   T BE 1999 1397
## 3    NR   T CZ 1999 1455
## 4    NR   T DK 1999 514
## 5    NR   T EL 1999 2116
## 6    NR   T ES 1999 5738
```

Add labels

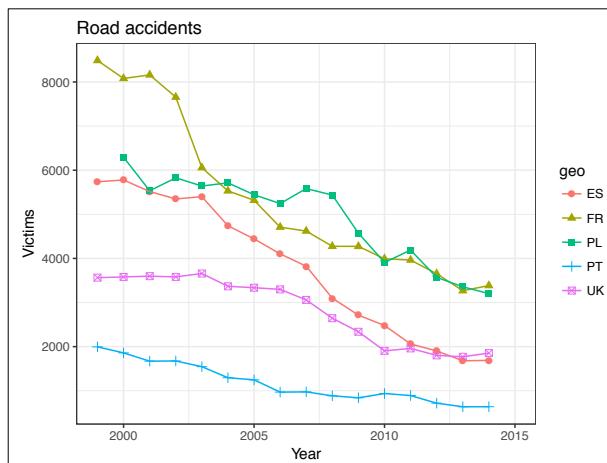
The `label_eurostat(x, lang = "en", ...)` gets definitions for Eurostat codes and replace them with labels in given language ("en", "fr" or "de").

```
dat <- label_eurostat(dat)
head(dat)
##    unit sex      geo time values
## 1 Number Total Austria 1999 1079
## 2 Number Total Belgium 1999 1397
## 3 Number Total Czech Republic 1999 1455
## 4 Number Total Denmark 1999 514
## 5 Number Total Greece 1999 2116
## 6 Number Total Spain 1999 5738
```

eurostat and plots

The `get_eurostat()` function returns tibbles in the long format. Packages `dplyr` and `tidyverse` are well suited to transform these objects. The `ggplot2` package is well suited to plot these objects.

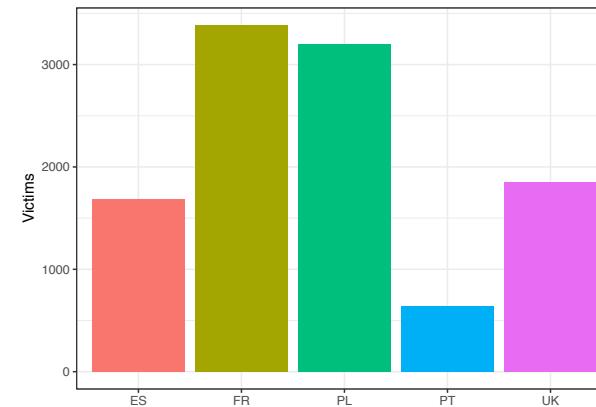
```
t1 <- get_eurostat("tsdtr420", filters =
  list(geo = c("UK", "FR", "PL", "ES", "PT")))
library("ggplot2")
ggplot(t1, aes(x = time, y = values, color = geo,
  group = geo, shape = geo)) +
  geom_point(size = 2) +
  geom_line() + theme_bw() +
  labs(title="Road accidents", x = "Year", y = "Victims")
```



Road accidents

```
library("dplyr")
t2 <- t1 %>% filter(time == "2014-01-01")
ggplot(t2, aes(geo, values, fill=geo)) +
  geom_bar(stat = "identity") + theme_bw() +
  theme(legend.position = "none")+
  labs(title="Road accidents in 2014", x="", y="Victims")
```

Road accidents in 2014



eurostat and maps

Fetch and process data

There are three function to work with geospatial data from GISCO. The `get_eurostat_geospatial()` returns preprocessed spatial data as sp-objects or as data frames. The `merge_eurostat_geospatial()` both downloads and merges the geospatial data with a preloaded tabular data. The `cut_to_classes()` is a wrapper for `cut()` - function and is used for categorizing data for maps with tidy labels.

```
library("eurostat")
library("dplyr")
```

```
fertility <- get_eurostat("demo_r_frate3") %>%
  filter(time == "2014-01-01") %>%
  mutate(cat = cut_to_classes(values, n=7, decimals=1))
```

```
mapdata <- merge_eurostat_geodata(fertility,
  resolution = "20")
```

```
head(select(mapdata,geo,values,cat,long,lat,order,id))
##   geo values   cat long lat order id
## 1 AT124 1.3 ~< 1.5 15.54245 48.90770 214 10
## 2 AT124 1.39 1.3 ~< 1.5 15.75363 48.85218 215 10
## 3 AT124 1.39 1.3 ~< 1.5 15.88763 48.78511 216 10
## 4 AT124 1.39 1.3 ~< 1.5 15.81535 48.69270 217 10
## 5 AT124 1.39 1.3 ~< 1.5 15.94094 48.67173 218 10
## 6 AT124 1.39 1.3 ~< 1.5 15.90833 48.59815 219 10
```

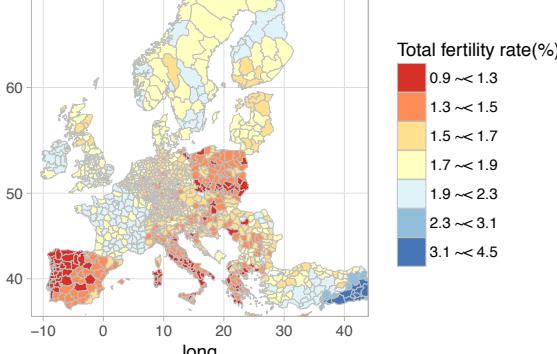
Draw a cartogram

The object returned by `merge_eurostat_geospatial()` are ready to be plotted with `ggplot2` package. The `coord_map()` function is useful to set the projection while `labs()` adds annotations o the plot.

```
library("ggplot2")
ggplot(mapdata, aes(x = long, y = lat, group = group))+
  geom_polygon(aes(fill=cat), color="grey", size = .1)+ 
  scale_fill_brewer(palette = "RdYlBu") +
  labs(title="Fertility rate, by NUTS-3 regions, 2014",
       subtitle="Avg. number of live births per woman",
       fill="Total fertility rate(%)") + theme_light()+
  coord_map(xlim=c(-12,44), ylim=c(35,67))
```

Fertility rate, by NUTS-3 regions, 2014

Avg. number of live births per woman



Animate ggplots with ganimate :: CHEAT SHEET

Core Concepts

ganimate builds on ggplot2's grammar of graphics to provide functions for animation. You add them to plots created with `ggplot()` the same way you add a geom.

Main Function Groups

- `transition_*`(): What variable controls change and how?
- `view_*`(): Should the axes change with the data?
- `enter/exit_*`(): How does new data get added the plot? How does old data leave?
- `shadow_*`(): Should previous data be "remembered" and shown with current data?
- `ease_aes()`: How do you want to handle the pace of change between transition values?

Note: you only need a `transition_*`() or `view_*`() to make an animation. The other function groups enable you to add features or alter ganimate's default settings.

Starting Plots

```
library(tidyverse)
library(ganimate)

a <- ggplot(diamonds,
            aes(carat, price)) +
  geom_point()

b <- ggplot(txhousing,
            aes(month, sales)) +
  geom_col()

c <- ggplot(economics,
            aes(date, psavert)) +
  geom_line()
```

transition_*

transition_states()

```
a + transition_states(color, transition_length = 3, state_length = 1)
```

We're cycling between values of `color`, ...

... and spending 3 times as long going to the next cut as we do pausing there.

transition_time()

```
b + transition_time(year, range = c(2002L, 2006L))
```

We're cycling through each `year` of the data...

...from 2002 to 2006 (range is optional; default is the whole time frame). Unlike `transition_states()`, `transition_time()` treats the data as continuous and so the transition length is based on the actual values. Using 2002L instead of 2002 because the underlying data is an integer.

transition_reveal()

```
c + transition_reveal(date)
```

We're adding each `date` of the data on top of 'old' data

`transition_length` and `filter_length` work the same as `transition/state_length()` in `transition_states()`...

transition_filters()

```
a + transition_filter(transition_length = 3,
                      filter_length = 1,
                      cut == "Ideal",
                      Deep = depth >= 60)
```

... but now we're cycling between these two filtering conditions. **Names** are optional, but can be useful (see "Label variables" on next page).

Other transitions

- `transition_manual()`: Similar to `transition_states()`, but without intermediate states.
- `transition_layers()`: Add layers (geoms) one at time.
- `transition_components()`: Transition elements independently from each other.
- `transition_events()`: Each element's duration can be controlled individually.

Baseline Animation

```
anim_a <- a + transition_states(color, transition_length = 1)
```

view_*

view_follow()

```
anim_a +  
  view_follow(fixed_x = TRUE,  
              fixed_y = c(2500, NA))
```

x-axis shows **full range**, y shows [2500, as much is needed for that frame]. Default is for both axis to vary as needed.

view_step()

```
anim_a +  
  view_step(pause_length = 2,  
            step_length = 1,  
            nstep = 7)
```

We're spending **twice** as long moving between views as staying at them...

... and we're cycling between **seven** views. Seven is the number of steps in the transition, so the view is changing when the points are static, and visa versa. Views are determined by what data is in the current frame.

view_zoom()

`view_zoom()` works similarly to `view_step()`, except it changes the view by zooming and panning.

Note: both `view_step()` and `view_zoom()` have `view_*_manual()` versions for setting views directly instead of inferring it from frame data.

Animate ggplots with ganimate :: CHEAT SHEET

enter/exit_*

Every `enter_*`() function has a corresponding `exit_*`() function, and visa versa.

`enter/exit_fade()`

When new points need to be added, they will start transparent and become opaque.

`enter_grow()/exit_shrink()`

`anim_a + enter_fade()`

When extra points need to be removed, they will shrink in size before disappearing.

`enter/exit_fly()`

`anim_a + enter_fly(x_loc = 0, y_loc = 0)`

When new points need to be added, they will fly in from (0, 0).

`enter/exit_drift()`

`anim_a + exit_drift(x_mod = 3, y_mod = -2)`

When extra points need to be removed, they drift 3 units to the right and down 2 units before disappearing.

`enter/exit_recolour() (or enter/exit_recolor())`

`anim_a + enter_recolour(color = "red")`

When new points need to be added, they start as red before transitioning to their correct color.

Note: `enter/exit_*`() functions can be combined so that you can have old data fade away and shrink to nothing by adding `exit_fade()` and `exit_shrink()` to the plot.

shadow_*

`shadow_wake()`

`anim_a + shadow_wake(wake_length = 0.05)`

Points have a wake of points with the data from the last 5% of frames.

`shadow_trail()`

`anim_a + shadow_trail(distance = 0.05)`

Animation will keep the points from 5% of the frames, spaced as evenly as possible.

`shadow_mark()`

`anim_a + shadow_mark(color = "red")`

Animation will keep past states plotted in red (but not the intermediate frames).

ease_aes()

`ease_aes()` allows you to set an easing function to control the rate of change between transition states. See `?ease_aes` for the full list.

Compare:

```
anim_a  
anim_a + ease_aes("cubic-in") # Change easing of all aesthetics  
anim_a + ease_aes(x = "elastic-in") # Only change `x` (others remain "linear")
```

Saving animations

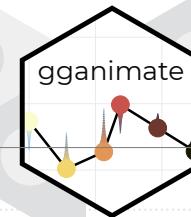
```
animation_to_save <- anim_a + exit_shrink()  
anim_save("first_saved_animation.gif", animation = animation_to_save)
```

Since the `animation` argument uses your last rendered animation by default, this also works:

```
anim_a + exit_shrink()  
anim_save("second_saved_animation.gif")
```

`anim_save()` uses `gifsicle` to render the animation as a .gif file by default. You can use the `renderer` argument for other output types including video files (`av_renderer()` or `ffmpeg_renderer()`) or spritesheets (`sprite_renderer()`):

```
# requires you to have the av package installed  
anim_save("third_saved_animation.mp4",  
          renderer = av_renderer())
```



Label variables

`ganimate`'s `transition_*`() functions create label variables you can pass to (sub)titles and other labels with the `glue` package. For example, `transition_states()` has `next_state`, which is the name of the state the animation is transitioning towards. Label variables are different between transitions, and details are included in the documentation of each.

```
anim_a + labs(subtitle = "Moving to {next_state}")
```

We're using the `next_state` label variable to tell the viewer where we're going.

| Label variable | Description | Transitions |
|--|---|--------------------------|
| <code>transitioning</code> | TRUE if the current frame is an transition frame, FALSE otherwise | states, layers, filter |
| <code>previous_state/layer</code> | Last shown state/layer | states, layers |
| <code>next_state/layer</code> | State/layer that will been shown next | states, layers |
| <code>closest_state/layer</code> | State/layer that current frame is closest to (if between states/layers, either next or closest). | states, layers |
| <code>previous/closest/next_filter/expression</code> | Similar to their state/layer analogs.
<code>*_filter</code> variables return the name of the filter, <code>*_expression</code> variables return the condition. | filter |
| <code>frame_time</code> | Time of current frame | time, components, events |
| <code>frame_along</code> | Current frame's value for the dimension we're transitioning over | reveal |
| <code>nlayers</code> | Number of layers (total, not just currently shown) | layer |



GWAS Catalog access with gwasrapidd

Introduction

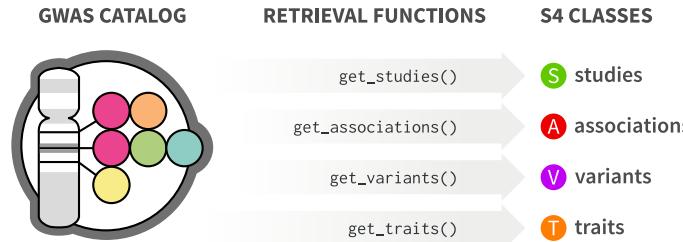
The **GWAS Catalog** is a service provided by the EMBL-EBI and NHGRI that offers a manually curated and freely available database of published genome-wide association studies (GWAS).

The GWAS Catalog data provided by the **RESTful API** is organized around four core entities:

- **studies**
- **associations**
- **variants**
- **traits**

Get GWAS Catalog Entities

gwasrapidd facilitates the access to the Catalog via the RESTful API, allowing you to programmatically retrieve data directly into R. Each of the four entities is mapped to an S4 object of a class of the same name.



| Search by | Example | S A V T |
|------------------|--|---------|
| study_id | "GCST000858" | ● ● ● ● |
| association_id | "24300113" | ● ● ● ● |
| variant_id | "rs12752552" | ● ● ● ● |
| efo_id | "EFO_0005543" | ● ● ● ● |
| pubmed_id | "21626137" | ● ● ● ● |
| user_requested | TRUE | ● ● ● ● |
| full_pvalue_set | FALSE | ● ● ● ● |
| efo_uri | "http://www.ebi.ac.uk/efo/EFO_0004761" | ● ● ● ● |
| genomic_range | list(chromosome = "22", start = 1L, end = 15473564L) | ● ● ● ● |
| gene_name | "BRCA1" | ● ● ● ● |
| efo_trait | "lung adenocarcinoma" | ● ● ● ● |
| reported_trait | "Breast cancer" | ● ● ● ● |
| cytogenetic_band | "1p36.33" | ● ● ● ● |

S4 Representation of GWAS Catalog Entities

S4 class studies

The **studies** object consists of eight slots, each a table (tibble). Each study is an observation (row) in the studies table — main table. All tables have the column `study_id` as primary key.

For details about the studies S4 class: `class?studies`.

| studies | genotyping_techs | countries_of_recruitment |
|--|--------------------------------------|---------------------------------|
| • <code>study_id</code> | • <code>study_id</code> | • <code>study_id</code> |
| • <code>reported_trait</code> | • genotyping technology | • <code>ancestry_id</code> |
| • <code>initial_sample_size</code> | • platforms | • <code>country_name</code> |
| • <code>replication_sample_size</code> | • <code>study_id</code> | • <code>major_area</code> |
| • <code>gxe</code> | • manufacturer | • <code>region</code> |
| • <code>gxg</code> | • ancestries | • countries_of_origin |
| • <code>snp_count</code> | • <code>study_id</code> | • <code>study_id</code> |
| • <code>qualifier</code> | • <code>ancestry_id</code> | • <code>ancestry_id</code> |
| • <code>imputed</code> | • <code>type</code> | • <code>country_name</code> |
| • <code>pooled</code> | • <code>number_of_individuals</code> | • <code>major_area</code> |
| • <code>study_design_comment</code> | • <code>ancestral_groups</code> | • <code>region</code> |
| • <code>full_pvalue_set</code> | • <code>study_id</code> | • publications |
| • <code>user_requested</code> | • <code>ancestry_id</code> | • <code>study_id</code> |
| | • <code>ancestral_group</code> | • <code>pubmed_id</code> |
| | | • <code>publication_date</code> |
| | | • <code>publication</code> |
| | | • <code>title</code> |
| | | • <code>author_fullname</code> |
| | | • <code>author_orcid</code> |

S4 class associations

The **associations** object consists of six slots, each a table (tibble). Each association is an observation (row) in the associations table — main table. All tables have the column `association_id` as primary key.

For details about the associations S4 class: `class?associations`.

| associations | loci | genes |
|---------------------------------------|------------------------------------|-------------------------------|
| • <code>association_id</code> | • <code>association_id</code> | • <code>association_id</code> |
| • <code>pvalue</code> | • <code>locus_id</code> | • <code>locus_id</code> |
| • <code>pvalue_description</code> | • <code>haplotype.snp_count</code> | • <code>gene_name</code> |
| • <code>pvalue_mantissa</code> | • <code>description</code> | • ensembl_ids |
| • <code>pvalue_exponent</code> | • risk_alleles | • <code>association_id</code> |
| • <code>multiple.snp.haplotype</code> | • <code>association_id</code> | • <code>locus_id</code> |
| • <code>snp.interaction</code> | • <code>locus_id</code> | • <code>gene_name</code> |
| • <code>snp.type</code> | • <code>variant_id</code> | • <code>ensembl_id</code> |
| • <code>standard.error</code> | • <code>risk allele</code> | • entrez_ids |
| • <code>range</code> | • <code>risk_frequency</code> | • <code>association_id</code> |
| • <code>or.per.copy.number</code> | • <code>genome.wide</code> | • <code>locus_id</code> |
| • <code>beta.number</code> | • <code>limited_list</code> | • <code>gene_name</code> |
| • <code>beta.unit</code> | | • <code>entrez_id</code> |
| • <code>beta.direction</code> | | |
| • <code>beta.description</code> | | |
| • <code>last.mapping.date</code> | | |
| • <code>last.update.date</code> | | |

S4 class variants

The **variants** object consists of four slots, each a table (tibble). Each variant is an observation (row) in the variants table — main table. All tables have the column `variant_id` as primary key.

For details about the variants S4 class: `class?variants`.

| variants | genomic_contexts | ensembl_ids |
|----------------------------------|------------------------------------|---------------------------|
| • <code>variant_id</code> | • <code>variant_id</code> | • <code>variant_id</code> |
| • <code>merged</code> | • <code>merged</code> | • <code>gene_name</code> |
| • <code>functional.class</code> | • <code>chromosome.name</code> | • <code>ensembl_id</code> |
| • <code>chromosome.name</code> | • <code>chromosome.position</code> | • <code>entrez_ids</code> |
| • <code>chromosome.region</code> | • <code>chromosome.region</code> | • <code>variant_id</code> |
| • <code>last.update.date</code> | • <code>last.update.date</code> | • <code>gene_name</code> |
| | | • <code>entrez_id</code> |

V

S4 class traits

The **traits** object consists of one slot only, a table (tibble) of GWAS Catalog EFO traits. Each EFO trait is an observation (row) in the traits table — main table.

For details about the traits S4 class: `class?traits`.

| traits |
|-----------------------|
| • <code>efo_id</code> |
| • <code>trait</code> |
| • <code>uri</code> |

T

Manipulate Cases

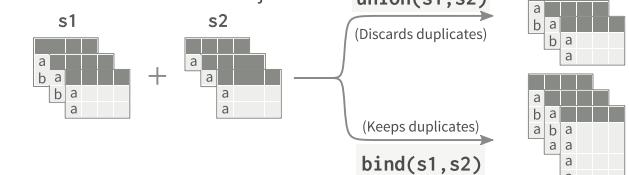
Get a **studies** object **s** of two GWAS studies:

```
s <- get_studies(study_id = c('a', 'b'))
```

Subset object **s** by either identifier or position using '`[`:



Combine two studies' objects:



h2o:: CHEAT SHEET

H₂O.ai

Dataset Operations

DATA IMPORT / EXPORT

h2o.uploadFile: Upload a file into H2O from a client-side path, and parse it.

h2o.downloadCSV: Download a H2O dataset to a client-side CSV file.

h2o.importFile: Import a file into H2O from a server-side path, and parse it.

h2o.exportFile: Export an H2O Data Frame to a server-side file.

h2o.parseRaw: Parse a raw data file.

NATIVE R TO H2O COERCION

as.h2o: Convert a R object to an H2O object

H2O TO NATIVE R COERCION

as.data.frame: Check if an object is a data frame, and coerce it if possible.

DATA GENERATION

h2o.createFrame: Creates a data frame in H2O with real-valued, categorical, integer, and binary columns specified by the user, with optional randomization.

h2o.runif: Produce a vector of random uniform numbers.

h2o.interaction: Create interaction terms between categorical features of an H2O Frame.

h2o.target_encode_apply: Target encoding map to an H2O Data Frame, which can improve performance of supervised learning models for high cardinality categorical columns.

DATA SAMPLING / SPLITTING

h2o.splitFrame: Split an existing H2O dataset according to user-specified ratios.

MISSING DATA HANDLING

h2o.impute: Impute a column of data using the mean, median, or mode.

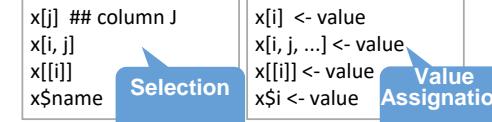
h2o.insertMissingValues: Replaces a user-specified fraction of entries in an H2O dataset with missing values.

h2o.na.omit: Remove Rows With NAs.

General Operations

SUBSCRIPTING

Subscripting example to pull (/push) pieces from (/to) a H2O Parsed Data object.



SUBSETTING

h2o.head, h2o.tail: Object's Start or End.

DATA ATTRIBUTES

h2o.names: Return column names for an H2O Frame. Also: **h2o.colnames**

names<-: Set the row or column names of a H2O Frame. Also: **colnames<-**

h2o.dim: Retrieve object dimensions.

h2o.length: Length of vector, list or factor.

h2o.nrow: Number of H2O Frame rows.

h2o.ncol: Number of H2O Frame columns.

h2o.anyFactor: Check if an H2O Frame object has any categorical data columns.

is.factor, is.character, is.numeric: Check Column's Data Type.

DATA TYPE COERCION

Convert to:

h2o.asfactor, as.factor: Factor.

h2o.as_date, as.Date: Date.

h2o.ascharacter, as.character: Character.

h2o.asnumeric, as.numeric: Numeric.

BASIC DATA MANIPULATION

c: Combine Values into a Vector or List.

h2o.cbind; h2o.rbind: Combine a sequence of H2O datasets by column (cbind) or rows (rbind).

h2o.merge: Merges 2 H2OFrames.

h2o.arrange: Sorts an H2OFrame by columns.

ELEMENT INDEX SELECTION

h2o.which: True Condition's Row Numbers

CONDITIONAL VALUE SELECTION

h2o.ifelse: Apply conditional statements to numeric vectors in an H2O Frame.

Math Operations

(math) vectorized function

MATH

h2o.abs: Compute the absolute value of x.

h2o.sqrt: Principal Square Root of x, \sqrt{x} .

h2o.ceiling: Take a single numeric argument x and return a numeric vector containing the smallest integers not less than the corresponding elements of x.

h2o.floor: Take a single numeric argument x and return a numeric vector containing the largest integers not greater than the corresponding elements of x.

h2o.trunc: Take a single numeric argument x and return a numeric vector containing the integers formed by truncating the values in x toward 0.

h2o.log: Compute natural logarithms. See also: **h2o.log10, h2o.log2, h2o.log1p**

h2o.exp: Compute the exponential function

h2o.cos, h2o.cosh, h2o.acos, h2o.sin, h2o.tan, h2o.tanh, Math: ?groupGeneric

sign: Return a vector with the signs of the corresponding elements of x (the sign of a real number is 1, 0, or -1 if the number is positive, zero, or negative, respectively).

&& (Vectorized AND), || (Vectorized OR), !x, %in%, Ops: +, -, *, /, ^, %%, %/, ==, !=, <, <=, >=, >, &, |, !

CUMULATIVE

h2o.cummax: Vector of the cumulative maxima of the elements of the argument.

h2o.cummin: Vector of the cumulative minima of the elements of the argument.

h2o.cumprod: Vector of the cumulative products of the elements of the argument.

h2o.cumsum: Vector of the cumulative sums of the elements of the argument.

PRECISION

h2o.round: Round values to the specified number of decimal places. The default is 0.

h2o.signif: Round values to the specified number of significant digits.

Group By Summaries

(group by) summary function

nrow: Count the number of rows.

max: All input argument's Maximum.

min: All input argument's Minimum.

sum: All argument values Sum.

mean: (Trimmed) arithmetic mean.

sd: Calculate the standard deviation of a column of continuous real valued data.

var: Compute the variance of x.

Generic Summaries

NON-GROUP_BY SUMMARIES

h2o.median: Calculate the median of x.

h2o.range: Input argument's Min/Max Vector

h2o.cor: Correlation Matrix of H2O Frames.

h2o.quantile: Obtain and display quantiles for an H2O Frame Column.

h2o.hist: Compute a histogram over a numeric H2O Frame Column.

h2o.prod: Product of all arguments values.

h2o.any: Given a set of logical vectors, determine if at least one of the values is true.

h2o.all: Given a set of logical vectors, determine if all of the values are true.

NON-GROUP_BY SUMMARIES: GENERIC

h2o.summary: Produce result summaries of the results of various model fitting functions.

Aggregations

ROW / COLUMN AGGREGATION

apply: Apply a function over an H2O parsed data object (an array) margins.

GROUP BY AGGREGATION

h2o.group_by: Apply an aggregate function to each group of an H2O dataset.

TABULATION

h2o.table: Use the cross-classifying factors to build a table of counts at each combination of factor levels.

h2o:: CHEAT SHEET



Data Modeling

MODEL TRAINING: SUPERVISED LEARNING

h2o.deeplearning: Perform Deep Learning Neural Networks on an H2OFrame.

h2o.gbm: Build Gradient Boosted Regression Trees or Classification Trees.

h2o.glm: Fit a Generalized Linear Model, specified by a response variable, a set of predictors, and the error distribution.

h2o.naiveBayes: Compute Naive Bayes classification probabilities on an H2O Frame.

h2o.randomForest: Perform Random Forest Classification on an H2O Frame.

h2o.xgboost: Build an Extreme Gradient Boosted Model using the XGBoost backend.

h2o.stackedEnsemble: Build a stacked ensemble (aka. Super Learner) using the specified H2O base learning algorithms.

h2o.automl: Automates the Supervised Machine Learning Model Training Process: Automatically Trains and Cross-validates a set of Models, and trains a Stacked Ensemble.

MODEL TRAINING: UNSUPERVISED LEARNING

h2o.prcomp: Perform Principal Components Analysis on the given H2O Frame.

h2o.kmeans: Perform k-means Clustering on the given H2O Frame.

h2o.anomaly: Detect anomalies in a H2O Frame using a H2O Deep Learning Model with Auto-Encoding.

h2o.deepfeatures: Extract the non-linear features from a H2O Frame using a H2O Deep Learning Model.

h2o.glrn: Builds a Generalized Low Rank Decomposition of an H2O Frame.

h2o.svd: Singular value decomposition of an H2O Frame using the power method.

h2o.word2vec: Trains a word2vec model on a String column of an H2O data frame.

SURVIVAL MODELS: TIME-TO-EVENT

h2o.coxph: Trains a Cox Proportional Hazards Model (CoxPH) on an H2O Frame.

Data Munging

GENERAL COLUMN MANIPULATION

is.na: Display missing elements.

FACTOR LEVEL MANIPULATIONS

h2o.levels: Display a list of the unique values found in a categorical data column.

h2o.relevel: Reorders levels of an H2O factor, similarly to standard R's relevel.

h2o.setLevels: Set Levels of H2O Factor.

NUMERIC COLUMN MANIPULATIONS

h2o.cut: Convert H2O Numeric Data to Factor by breaking it into Intervals.

CHARACTER COLUMN MANIPULATIONS

h2o.strsplit: "String Split": Splits the given factor column on the input split.

h2o.tolower: Convert the characters of a character vector to lower case.

h2o.toupper: Convert the characters of a character vector to upper case.

h2o.trim: "Trim spaces": Remove leading and trailing white space.

h2o.gsub: Match a pattern & replace **all** instances (occurrences) of the matched pattern with the replacement string globally.

h2o.sub: Match a pattern & replace the **first** instance (occurrence) of the matched pattern with the replacement string.

DATE MANIPULATIONS

h2o.month: Convert Milliseconds to Months in H2O Datasets (Scale: 0 to 11).

h2o.year: Convert Milliseconds to Years in H2O Datasets, indexed starting from 1900.

h2o.day: Convert Milliseconds to Day of Month in H2O Datasets (Scale: 1 to 31).

h2o.hour: Convert Milliseconds to Hour of Day in H2O Datasets (Scale: 0 to 23).

h2o.dayOfWeek: Convert Milliseconds to Day of Week in a H2OFrame (Scale: 0 to 6)

MATRIX OPERATIONS

%*%: Multiply two conformable matrices.

t: Returns the transpose of an H2OFrame.

Cluster Operations

H2O KEY VALUE STORE ACCESS

h2o.assign: Assign H2O hex.keys to R objects.

h2o.getFrame: Get H2O dataset Reference.

h2o.getModel: Get H2O model reference.

h2o.ls: Display a list of object keys in the running instance of H2O.

h2o.rm: Remove specified H2O Objects from the H2O server, but not from the R environment.

h2o.removeAll: Remove All H2O Objects from the H2O server, but not from the R environment.

H2O MODEL IMPORT / EXPORT

h2o.loadModel: Load H2OModel from disk.

h2o.saveModel: Save H2OModel object to disk.

h2o.download_pojo: Download the Scoring POJO (Plain Old Java Object) of an H2O Model.

h2o.download_mojo: Download the model in MOJO format.

H2O CLUSTER CONNECTION

h2o.init: Connect to a running H2O instance using all CPUs on the host.

h2o.shutdown: Shut down the specified H2O instance. All data on the server will be lost!

H2O CLUSTER INFORMATION

h2o.clusterInfo: Display the name, version, uptime, total nodes, total memory, total cores and health of a cluster running H2O.

h2o.clusterStatus: Retrieve information on the status of the cluster running H2O.

H2O LOGGING

h2o.clearLog: Clear all H2O R command and error response logs from the local disk.

h2o.downloadAllLogs: Download all H2O log files to the local disk.

h2o.logAndEcho: Write a message to the H2O Java log file and echo it back.

h2o.openLog: Open existing logs of H2O R POST commands and error responses on disk.

h2o.getLogPath: Get the file path for the H2O R command and error response logs.

h2o.startLogging: Begin logging H2O R POST commands and error responses.

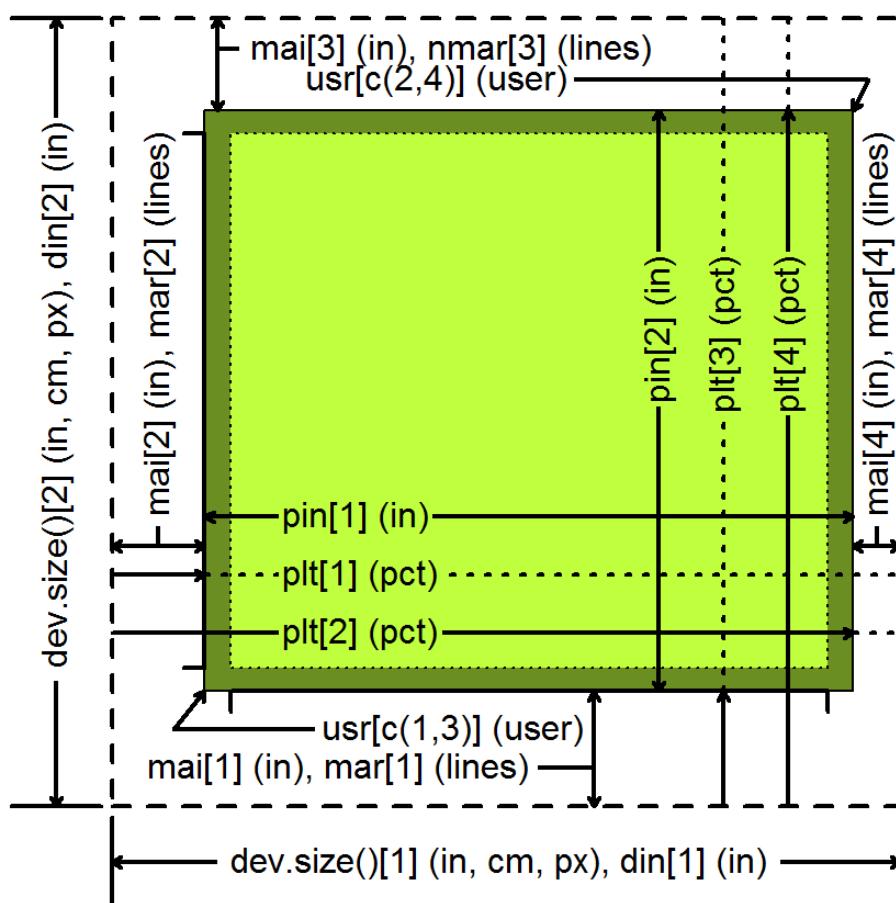
h2o.stopLogging: Stop logging H2O R POST commands and error responses.

How Big is Your Graph?

An R Cheat Sheet

Introduction

All functions that open a device for graphics will have **height** and **width** arguments to control the size of the graph and a **pointsize** argument to control the relative font size. In **knitr**, you control the size of the graph with the chunk options, **fig.width** and **fig.height**. This sheet will help you with calculating the size of the graph and various parts of the graph within R.



Your graphics device

`dev.size()` (width, height)
`par("din")` (r.o.) (width, height) in inches

Both the `dev.size` function and the `din` argument of `par` will tell you the size of the graphics device. The `dev.size` function will report the size in

1. inches (`units="in"`), the default
2. centimeters (`units="cm"`)
3. pixels (`units="px"`)

Like several other `par` arguments, `din` is read only (r.o.) meaning that you can ask its current value (`par("din")`) but you cannot change it (`par(din=c(5,7))` will fail).

Your plot margins

`par("mai")` (bottom, left, top, right) in inches
`par("mar")` (bottom, left, top, right) in lines

Margins provide you space for your axes, axis labels, and titles.

A "line" is the amount of vertical space needed for a line of text.

If your graph has no axes or titles, you can remove the margins (and maximize the plotting region) with

```
par(mar=rep(0,4))
```

Your plotting region

`par("pin")` (width, height) in inches
`par("plt")` (left, right, bottom, top) in pct

The `pin` argument `par` gives you the size of the plotting region (the size of the device minus the size of the margins) in inches.

The `plt` argument gives you the percentage of the device from the left/bottom edge up to the left edge of the plotting region, the right edge, the bottom edge, and the top edge. The first and third values are equivalent to the percentage of space devoted to the left and bottom margins. Subtract the second and fourth values from 1 to get the percentage of space devoted to the right and top margins.

Your x-y coordinates

`par("usr")` (xmin, ymin, xmax, ymax)

Your x-y coordinates are the values you use when plotting your data. This normally is not the same as the values you specified with the `xlim` and `ylim` arguments in `plot`. By default, R adds an extra 4% to the plotting range (see the dark green region on the figure) so that points right up on the edges of your plot do not get partially clipped. You can override this by setting `xaxs="i"` and/or the `yaxs="i"` in `par`.

Run `par("usr")` to find the minimum X value, the maximum X value, the minimum Y value, and the maximum Y value. If you assign new values to `usr`, you will update the x-y coordinates to the new values.

Getting a square graph

`par("pty")`

You can produce a square graph manually by setting the width and height to the same value and setting the margins so that the sum of the top and bottom margins equal the sum of the left and right margins. But a much easier way is to specify `pty="s"`, which adjusts the margins so that the size of the plotting region is always square, even if you resize the graphics window.

Converting units

For many applications, you need to be able to translate user coordinates to pixels or inches. There are some cryptic shortcuts, but the simplest way is to get the range in user coordinates and measure the proportion of the graphics device devoted to the plotting region.

```
user.range <- par("usr")[c(2,4)] -  
           par("usr")[c(1,3)]
```

```
region.pct <- par("plt")[c(2,4)] -  
              par("plt")[c(1,3)]
```

```
region.px <-  
           dev.size(units="px") * region.pct
```

```
px.per.xy <- region.px / user.range
```

To convert a horizontal or distance from the x-coordinate value to pixels, multiply by `px.per.xy[1]`. To convert a vertical distance, multiply by `region.px.per.xy[2]`. To convert a diagonal distance, you need to invoke Pythagoras.

```
a.px <- x.dist*px.per.xy[1]  
b.px <- y.dist*px.per.xy[2]  
c.px <- sqrt(a.px^2+b.px^2)
```

To rotate a string to match the slope of a line segment, you need to convert the distances to pixels, calculate the arctangent, and convert from radians to degrees.

```
segments(x0, y0, x1, y1)  
delta.x <- (x1 - x0) * px.per.xy[1]  
delta.y <- (y1 - y0) * px.per.xy[2]  
angle.radians <- atan2(delta.y, delta.x)  
angle.degrees <- angle.radians * 180 / pi  
text(x1, y1, "TEXT", srt=angle.degrees)
```

Panels

`par("fig")` (width, height) in pct
`par("fin")` (width, height) in inches

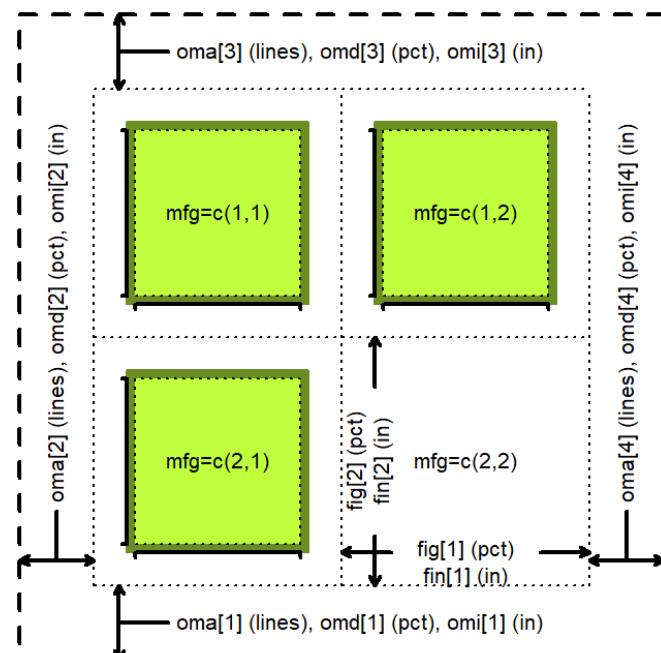
If you display multiple plots within a single graphics window (e.g., with the `mfrow` or `mfcol` arguments of `par` or with the `layout` function), then the `fig` and `fin` arguments will tell you the size of the current subplot window in percent or inches, respectively.

`par("oma")` (bottom, left, top, right) in lines
`par("omd")` (bottom, left, top, right) in pct
`par("omi")` (bottom, left, top, right) in inches

Each subplot will have margins specified by `mai` or `mar`, but no outer margin around the entire set of plots, unless you specify them using `oma`, `omd`, or `omi`. You can place text in the outer margins using the `mtext` function with the argument `outer=TRUE`.

`par("mfg")` (r, c) or (r, c, maxr, maxc)

The `mfg` argument of `par` will allow you to jump to a subplot in a particular row and column. If you query with `par("mfg")`, you will get the current row and column followed by the maximum row and column.



Character and string sizes

`strheight()`

The `strheight` functions will tell you the height of a specified string in inches (`units="inches"`), x-y user coordinates (`units="user"`) or as a percentage of the graphics device (`units="figure"`).

For a single line of text, `strheight` will give you the height of the letter "M". If you have a string with one or more linebreaks ("\n"), the `strheight` function will measure the height of the letter "M" plus the height of one or more additional lines. The height of a line is dependent on the line spacing, set by the `lheight` argument of `par`. The default line height (`lheight=1`), corresponding to single spaced lines, produces a line height roughly 1.5 times the height of "M".

`strwidth()`

The `strwidth` function will produce different widths to individual characters, representing the proportional spacing used by most fonts (a "W" using much more space than an "i"). For the width of a string, the `strwidth` function will sum up the lengths of the individual characters in the string.

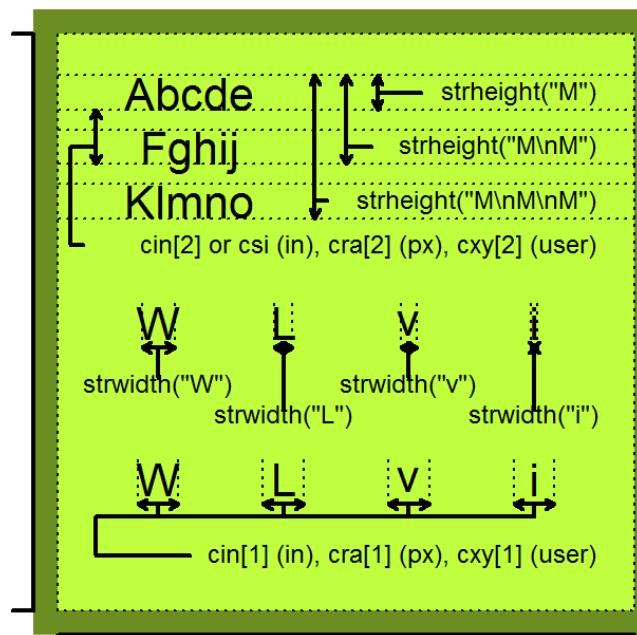
`par("cin")` (r.o.) (width, height) in inches
`par("csi")` (r.o.) height in inches
`par("cra")` (r.o.) (width, height) in pixels
`par("cxy")` (r.o.) (width, height) in xy coordinates

The single value returned by the `csi` argument of `par` gives you the height of a line of text in inches. The second of the two values returned by `cin`, `cra`, and `cxy` gives you the height of a line, in inches, pixels, or xy (user) coordinates.

The first of the two values returned by the `cin`, `cra`, and `cxy` arguments to `par` gives you the approximate width of a single character, in inches, pixels, or xy (user) coordinates. The width, very slightly smaller than the actual width of the letter "W", is a rough estimate at best and ignores the variable width of individual letters.

These values are useful, however, in providing fast ratios of the relative sizes of the differing units of measure

`px.per.in <- par("cra") / par("cin")`
`px.per.xy <- par("cra") / par("cxy")`
`xy.per.in <- par("cxy") / par("cin")`



If your fonts are too big or too small

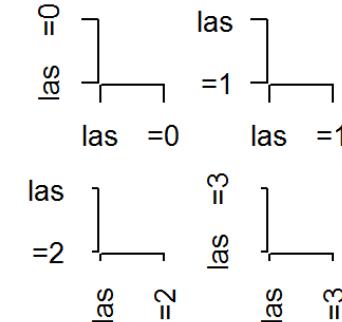
Fixing this takes a bit of trial and error.

- Specify a larger/smaller value for the `pointsize` argument when you open your graphics device.
- Try opening your graphics device with different values for `height` and `width`. Fonts that look too big might be better proportioned in a larger graphics window.
- Use the `cex` argument to increase or decrease the relative size of your fonts.

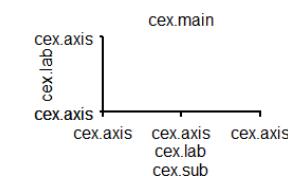
If your axes don't fit

There are several possible solutions.

- You can assign wider margins using the `mar` or `mai` argument in `par`.
- You can change the orientation of the axis labels with `las`. Choose among
 - `las=0` both axis labels parallel
 - `las=1` both axis labels horizontal
 - `las=2` both axis labels perpendicular
 - `las=3` both axis labels vertical.



- change the relative size of the font
 - `cex.axis` for the tick mark labels.
 - `cex.lab` for `xlab` and `ylab`.
 - `cex.main` for the main title
 - `cex.sub` for the subtitle.



Leaflet Cheat Sheet



for 

an open-source JavaScript library for mobile-friendly interactive maps

Quick Start

Installation

Use `install.packages("leaflet")` to install the package or directly from Github `devtools::install_github("rstudio/leaflet")`.

First Map

```
m <- leaflet() %>% # leaflet works with the pipe operator
addTiles() %>% # setup the default OpenStreetMap map tiles
addMarkers(lng = 174.768, lat = -36.852, popup = "The birthplace of R")
# add a single point layer
```



Map Widget

Initialization

```
m <- leaflet(options = leafletOptions(...))
center           Initial geographic center of the map
zoom            Initial map zoom level
minZoom         Minimum zoom level of the map
maxZoom         Maximum zoom level of the map
```

Map Methods

```
m %>% setView(lng, lat, zoom, options = list())
  Set the view of the map (center and zoom level)
m %>% fitBounds(lng1, lat1, lng2, lat2)
  Fit the view into the rectangle [lng1, lat1] - [lng2, lat2]
m %>% clearBounds()
  Clear the bound, automatically determine from the map elements
```

Data Object

Both `leaflet()` and the `map` layers have an optional data parameter that is designed to receive spatial data with the following formats:

Base R

The arguments of all layers take normal R objects:

```
df <- data.frame(lat = ..., lng = ...)
leaflet(df) %>% addTiles() %>% addCircles()
```

`library(sp)` Useful functions:

SpatialPoints, SpatialLines, SpatialPolygons, ...

`library(maps)` Build a map of states with colors:

```
mapStates <- map("state", fill = TRUE, plot = FALSE)
leaflet(mapStates) %>% addTiles() %>%
  addPolygons(fillColor = topo.colors(10, alpha = NULL), stroke = FALSE)
```

sp package

maps package

Markers

Use markers to call out points, express locations with latitude/longitude coordinates, appear as icons or as circles.

Data come from vectors or assigned data frame, or `sp` package objects.

Icon Markers

Regular Icons: default and simple

```
addMarkers(lng, lat, popup, label) add basic icon markers
makeIcon(Icons) (iconUrl, iconWidth, iconHeight, iconAnchorX, iconAnchorY,
  shadowUrl, shadowWidth, shadowHeight, ...) customize marker icons
iconList() create a list of icons
```

Awesome Icons: customizable with colors and icons

```
addAwesomeMarkers, makeAwesomeIcon, awesomeIcons, awesomeIconList
```

Marker Clusters: option of `addMarkers()`

```
clusterOptions = markerClusterOptions()
```

`freezeAtZoom` Freeze the cluster at assigned zoom level

Circle Markers

```
addCircleMarkers(color, radius, stroke, opacity, ...)
```

Customize their color, radius, stroke, opacity

Popups and Labels

`addPopups(lng, lat, ...content..., options)` Add standalone popups

options = `popupOptions(closeButton=FALSE)`

`addMarkers(..., popup, ...)` Show popups with markers or shapes

`addMarkers(..., label, labelOptions...)` Show labels with markers or shapes

labelOptions = `labelOptions(noHide, textOnly, textSize, direction, style)`

`addLabelOnlyMarkers()` Add labels without markers

Lines and Shapes

Polygons and Polylines

`addPolygons(color, weight=1, smoothFactor=0.5, opacity=1.0, fillOpacity=0.5,
 fillColor = ~colorQuantile(~YOrRd, ALAND)(ALAND), highlightOptions, ...)`

`highlightOptions(color, weight=2, bringToFront=TRUE)` highlight shapes

Use `rmapshaper::ms_simplify` to simplify complex shapes

Circles `addCircles(lng, lat, weight=1, radius, ...)`

Rectangles `addRectangles(lng1, lat1, lng2, lat2, fillColor="transparent", ...)`

Basemaps

Default Tiles



Default Tiles

Third-Party Tiles addProviderTiles()



Use `addTiles()` to add a custom map tile URL template, use `addWMSTiles()` to add WMS (Web Map Service) tiles

GeoJSON and TopoJSON

There are two options to use the GeoJSON/TopoJSON data.

* To read into `sp` objects with the `geojsonio` or `rgdal` package:
`geojsonio::geojson_read(..., what="sp") rgdal::readOGR(..., "OGRGeoJSON")`

* Or to use the `addGeoJSON()` and `addTopoJSON()` functions:
`addTopoJSON/addGeoJSON(... weight, color, fill, opacity, fillOpacity...)`
 Styles can also be tuned separately with a `style: {}` object.

Other packages including `RJSONIO` and `jsonlite` can help fast parse or generate the data needed.

Shiny Integration

To integrate a Leaflet map into an app:

* In the UI, call `leafletOutput("name")`

* On the server side, assign a `renderLeaflet(...)` call to the output

* Inside the `renderLeaflet` expression, return a Leaflet map object

Modification

To modify an existing map or add incremental changes to the map, you can use `leafletProxy()`. This should be performed in an observer on the server side.

Other useful functions to edit your map:

`fitBounds(o, 0, 11, 11)` similar to `setView`

fit the view to within these bounds

`addCircles(1:10, 1:10, layerId = LETTERS[1:10])`

create circles with layerIds of "A", "B", "C"...

`removeShape(c("B", "F"))` remove some of the circles

`clearShapes()` clear all circles (and other shapes)

Inputs/Events

Object Events

Object event names generally use this pattern:

`input$MAPID_OBJECTCATEGORY_EVENTNAME`

Trigger an event changes the value of the Shiny input at this variable.

Valid values for `OBJECTCATEGORY` are `marker`, `shape`, `geojson` and `topojson`.

Valid values for `EVENTNAME` are `click`, `mouseover` and `mouseout`.

All of these events are set to either `NULL` if the event has never happened, or a `list()` that includes:

* `lat` The latitude of the object, if available; otherwise, the mouse cursor

* `lng` The longitude of the object, if available; otherwise, the mouse cursor

* `id` The layerId, if any

GeoJSON events also include additional properties:

* `featureId` The feature ID, if any

* `properties` The feature properties

Map Events

`input$MAPID_click` when the map background or basemap is clicked
`value -- a list with lat and lng`

`input$MAPID_bounds` provide the lat/lng bounds of the visible map area
`value -- a list with north, east, south and west`

`input$MAPID_zoom` an integer indicates the zoom level

Machine Learning Modelling in R :: CHEAT SHEET

Supervised & Unsupervised Learning

SUPERVISED LEARNING

| ALGORITHM | DESCRIPTION | R PACKAGE::FUNCTION | SAMPLE CODE |
|-----------|---|----------------------|---|
| NBC | A classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. | e1071::naiveBayes | naiveBayes(class ~ ., data = x) |
| kNN | A non-parametric method used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether k-NN is used for classification or regression. | class::knn | knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE) |
| REG | Model the linear relationship between a scalar dependent variable Y and one or more explanatory variables (or independent variables) denoted X | stats::lm | lm(dist ~ speed, data=cars) |
| LREG | Used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables. | stats::glm | glm(Y ~ ., family = binomial(link = "logit"), data = X) |
| TM | The idea is to consecutively divide (branch) the training dataset based on the input features until an assignment criterion with respect to the target variable into a "data bucket" (leaf) is reached. | rpart::rpart | rpart(Kyphosis ~ Age + Number + Start, data = kyphosis) |
| ANN | Neural networks are built from units called perceptrons. Perceptrons take some inputs, an activation function and an output. An ANN model is built up by combining perceptrons in structured layers. | neuralnet::neuralnet | neuralnet(f,data=train,,hidden=c(5,3),linear.output=T) |
| SVM | A data classification method that separates data using hyperplanes | e1071::svm | svm(formula, data = NULL, ..., subset, na.action = na.omit, scale = TRUE) |

UNSUPERVISED LEARNING

| ALGORITHM | DESCRIPTION | R PACKAGE::FUNCTION | SAMPLE CODE |
|-----------|---|--|--|
| PCA | A procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components. | stats::prcomp
stats::princomp
FactoMineR::PCA
ade4::dudi.pca
amap::acp | stats::prcomp(formula, data = NULL, subset, na.action,...)
stats::princomp(formula, data = NULL, subset, na.action,...)
FactoMineR::PCA(decathlon, quanti.sup = 11:12, qual.sup=13)
ade4::dudi.pca(deug\$stab, center = deug\$cent, scale = FALSE, scan = FALSE)
amap::acp(lubinski) |
| KMC | Aims at partitioning n observations into k clusters in which each observation belongs to the cluster with the nearest mean | stats::kmeans | kmeans(x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"), trace=FALSE) |
| HCL | An approach which builds a hierarchy from the bottom-up, and doesn't require the number of clusters to be specified beforehand. | stats::hclust | hclust(d, method = "complete", members = NULL) |

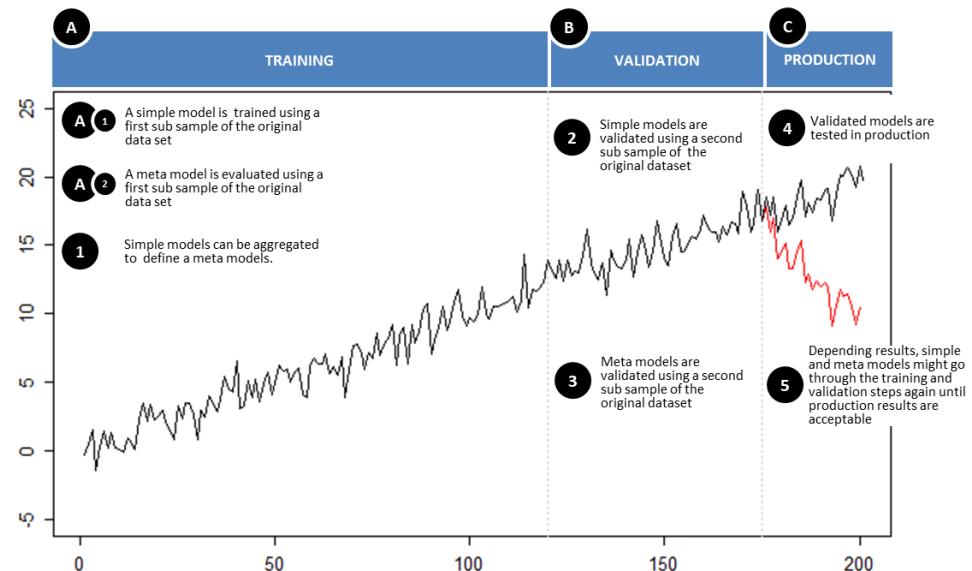
Meta-Algorithm, Time Series & Model Validation

| ALGORITHM | DESCRIPTION | R PACKAGE::FUNCTION | SAMPLE CODE |
|-----------|---|---|---|
| REGU | Regularisation L1 (Lasso) L2 (Ridge) | glmnet::glmnet | L1 : glmnet(myMatrixA, myMatrixB, family = "gaussian", alpha = 1)
L2 : glmnet(myMatrixA, myMatrixB, family = "gaussian", alpha = 0) |
| BOO | Boosting | gbmboost::gbmboost | gbmboost(Y ~ ., data = curr1[trnidxs,]) |
| BAG | Bagging | foreach::foreach
Tree models: ipred::bagging | foreach: d <- data.frame(x=1:10, y=rnorm(10))
s <- foreach(d=iter(d, by=1, .combine=rbind)
%>% parMap(.id = 1, identical, s, d)
ipred::bagging(formula, data, subset, na.action=na.rpart, .dots) |
| PRU | Pruning | rpart::rpart | prune(x, cp = 0.1) |
| RFO | Random Forrest | randomForest::randomForest | randomForest(X ~ ., data = Y, subset = mySub) |
| TS | Time Series
Lead-lag analysis, Auto-correlation, Spectral analysis, Time series clustering, Seasonality, Trend.... | xts::forecast
spectral::TTR
.... | Auto-correlation: acf(x, lag.max = NULL, type = c("correlation", "covariance", "partial"))
Spectral Analysis: spec.pgram(myTs, spans = NULL)
Seasonal Decomposition of Time Series - stl(x, s.window = 7, t.window = 50, t.jump = 1)..... |
| PM | Performance metrics | stats::outlierTest
stats::qqPlot
Classification:ROCR::Tree:caret::confusionMatrix | Regression: fit <- lm(Y~X,data=myData)
outlierTest(fit)
qqPlot(fit, main="QQ Plot") |
| BVT | Bias-Variance Tradeoff | tailored::tailored | Tailored to the analysis |
| CV | Cross validation | caret::createDataPartition
caret::createFolds | createDataPartition(classes, p = 0.8, list = FALSE) |
| LC | Learning Curves | caret::learning_curve_dat | learning_curve_dat(data, outcome = NULL, proportion = (1:10)/10, test_prop = 0, verbose = TRUE, ...) |

Standard Modelling Workflow



Time Series View



Intro stats with mosaic

ggformula version

Loading packages

```
library(mosaic)
```

Essential R syntax

Names in R are case sensitive

Function and arguments

```
rflip(10)
```

Optional arguments

```
rflip(10, prob = 0.8)
```

Assignment

```
x <- rflip(10, prob = 0.8)
```

Getting help on any function

```
help(mean)
```

Arithmetic operations

`+` `-` `*` `/` basic operations

`^` exponentiation

`()` grouping

`sqrt(x)` square root

`abs(x)` absolute value

`log10(x)` logarithm, base 10

`log(x)` natural logarithm, base e

`exp(x)` exponential function e^x

`factorial(k)` $k! = k(k - 1) \dots 1$

Logical operators

`==` is equal to (note double equal sign)

`!=` is not equal to

`<` is less than

`<=` is less than or equal to

`>` is greater than

`>=` is greater than or equal to

`&` A & B ("A and B") is TRUE if both A and B are TRUE

`|` A | B ("A or B") is TRUE if one or both of A and B are TRUE

`%in%` inclusion; for example

"C" %in% c("A", "B") is FALSE

Formula interface

Use for graphics, statistics, inference, and modeling operations.

```
goal(y ~ x, data = mydata)
```

Read as "Calculate `goal` for `y` using `mydata` "broken down by" `x`, or "modeled by" `x`.

```
mean(age ~ sex, data = HELPrc)
```

For graphics:

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

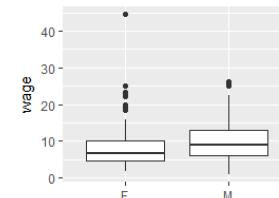
`y` : y-axis variable (*optional*)

`x` : x-axis variable (*required*)

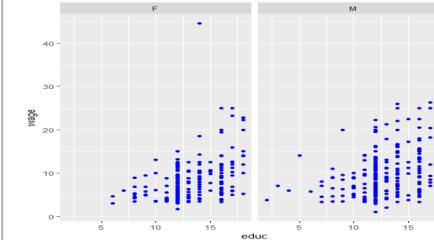
`z` : panel-by variable (*optional*)

`w` : color-by formula (*optional*)

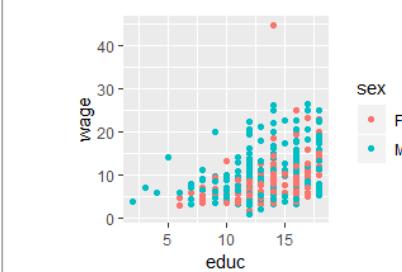
```
gf_boxplot(wage ~ sex,
           data = CPS85)
```



```
gf_point(wage ~ educ | sex,
          data = CPS85, color = "blue")
```



```
gf_point(wage ~ educ,
          data = CPS85, color = ~ sex)
```



Examining data

Print short summary of all variables

```
inspect(HELPrc)
```

Number of rows and columns

```
dim(HELPrc)
```

```
nrow(HELPrc)
```

```
ncol(HELPrc)
```

Print first rows or last rows

```
head(KidsFeet)
```

```
tail(KidsFeet, 10)
```

Names of variables

```
names(HELPrc)
```

One categorical variable

Counts by category

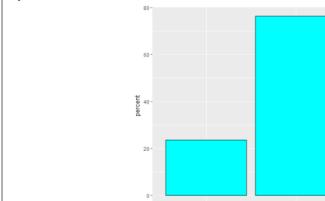
```
tally(~ sex, data = HELPrc)
```

Percentages by category

```
tally(~ sex, data = HELPrc,
      format = "percent")
```

Bar graph of percentages

```
gf_percent(~ sex,
           data = HELPrc, fill = "cyan",
           color = "black")
```



Tests and confidence intervals

Exact test

```
result1 <-
  binom.test(~ homeless ==
             "homeless", data = HELPrc)
```

Approximate test (large samples)

```
result2 <-
  prop.test(~ homeless ==
            "homeless", data = HELPrc,
            alternative = "less",
            p = 0.4)
```

Extract confidence intervals and p-values

```
confint(result1)
```

```
pval(result2)
```

One quantitative variable

Make output more readable

```
options(digits = 3)
```

Compute summary statistics

```
mean(~ cesd, data = HELPrc)
```

Other summary statistics work similarly

```
median() iqr() max() min()
```

```
fivenum() sd() var() sum()
```

Table of summary statistics

```
favstats(~ cesd, data = HELPrc)
```

Summary statistics by group

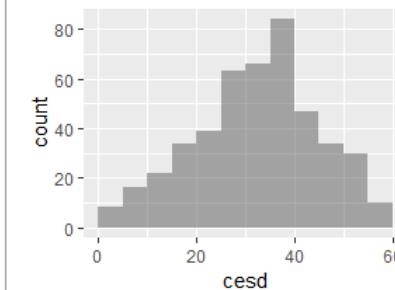
```
favstats(cesd ~ sex,
         data = HELPrc)
```

Quantiles

```
quantile(~ cesd, data = HELPrc,
         prob = c(0.25, 0.5, 0.8))
```

Histogram

```
gf_histogram(~ cesd,
            data = HELPrc, binwidth = 5,
            center = 2.5)
```



Normal probability plot

```
gf_qq(~ cesd, data = HELPrc)
```

Density plot

```
gf_dens(~ cesd, data = HELPrc,
        color = "blue", size = 1.25)
```

One-sample t-test

```
result <- t_test(~ cesd,
                  data = HELPrc, mu = 34)
```

Extract confidence intervals and p-values

```
confint(result)
```

```
pval(result)
```

Paired t-test

```
t_test(extra ~ group,
       data = sleep, paired = TRUE)
```

Data wrangling

Drop, rename, or reorder variables
`df <- select(HELPrc, c(id, age, gender = sex))`

Create new variables from existing ones
`KidsFeet <- mutate(KidsFeet, width_in = 0.394 * width)`

Extract specific rows from data
`girls_feet <- filter(KidsFeet, sex == "G")`

Sort data rows by value in column
`df <- arrange(KidsFeet, length)`

Compute summary statistics by group
`group_by(KidsFeet, sex) %>% summarize(mean_width = mean(width))`

For more, see [Tidyverse cheatsheet](#)

Importing data

Import data from file or URL

```
MustangPrice <-  
  read.file("C:/MustangPrice.csv")  
Note: R uses forward slashes in file paths  
kidsfeet <-  
  read.file("http://www.mosaic-  
web.org/go/datasets/kidsfeet.csv")
```

Randomization and simulation

Fix random number sequence
`set.seed(42)`

Toss coins
`rflip(10) # default prob is 0.5`

Do something repeatedly
`do(5) * rflip(10, prob = 0.75)`

Draw a simple random sample
`sample(LETTERS, 10)`

deal(Cards, 5) # poker hand

Resample with replacement
`Small <- sample(KidsFeet, 10)`

resample(� Small)

Random permutation (shuffling)
`shuffle(Cards)`

Random values from distributions
`rbinom(5, size = 10, prob = 0.7)`

`rnorm(5, mean = 10, sd = 2)`

Two categorical variables

Contingency table with margins
`tally(~ substance + sex,
 data = HELPrc, margins = TRUE)`

Percentages by column
`tally(~ sex | substance,
 data = HELPrc,
 format = "percent")`

Mosaic plot

```
my_tbl <- tally(substance ~ sex,  
                 data = HELPrc)  
mosaicplot(my_tbl, color = TRUE)
```



Test for proportions (approximate)

```
prop.test(homeless ~ sex,  
         success = "homeless",  
         data = HELPrc)
```

Distributions

Normal distribution function

`pnorm(13, mean = 10, sd = 2)`

Normal distribution function with graph

`xpnorm(1.645, mean = 0, sd = 1)`

Normal distribution quantiles

`qnorm(0.95) # mean = 0, sd = 1`

Normal distribution quantiles with graph

`xqnorm(0.85, mean = 10, sd = 2)`

Binomial density function ("size" means n)

`dbinom(5, size = 8, prob = 0.65)`

Binomial distribution function

`pbinom(5, size = 8, prob = 0.65)`

Central portion of distribution

```
cdist("norm", 0.95)  
cdist("t", c(0.90, 0.99), df = 5)
```

Plotting distributions

```
plotDist("binom", size = 8,  
        prob = 0.65, xlim = c(-1, 9))  
plotDist("norm", mean = 10,  
        sd = 2)
```

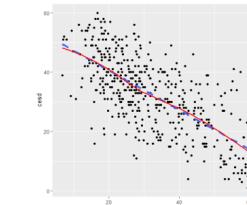
Two quantitative variables

Correlation coefficient

`cor(cesd ~ mcs, data = HELPrc)`

Scatterplot with regression line and smooth

```
gf_point(cesd ~ mcs,  
         data = HELPrc) %>%  
  gf_lm(size = 1.5, linetype =  
    "dashed") %>%  
  gf_smooth(color = "red")
```



Simple linear regression

```
cesdmodel <- lm(cesd ~ mcs,  
                 data = HELPrc)
```

`msummary(cesdmodel)`

Prediction

```
lm_fun <- makeFun(cesdmodel)  
lm_fun(mcs = 35)
```

Extract useful quantities

```
anova(cesdmodel)  
coef(cesdmodel)  
confint(cesdmodel)  
rsquared(cesdmodel)
```

Diagnostics; plot residuals

```
gf_dhistogram(~resid(cesdmodel))  
gf_qq(~resid(cesdmodel))
```

Diagnostics; plot residuals vs. fitted

```
gf_point(resid(cesdmodel) ~  
         fitted(cesdmodel)) %>%  
  gf_lm(size = 2)
```

Categorical response, quantitative predictor

Logistic regression

```
logit_mod <- glm(homeless ~ age,  
                  data = HELPrc,  
                  family = binomial)  
msummary(logit_mod)
```

Odds ratios and confidence intervals

```
exp(coef(logit_mod))  
exp(confint(logit_mod))
```

Quantitative response, categorical predictor

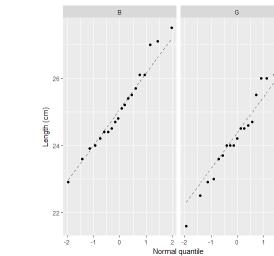
Two-level predictor: two-sample t test

Numeric summaries

```
favstats(~length | sex,  
         data = KidsFeet)
```

Graphic summaries

```
gf_qq(~ length | sex,  
      data = KidsFeet) %>%  
  gf_qqline() %>%  
  gf_labs(x = "Normal quantile",  
          y = "Length (cm)")
```



Two-sample t -test and confidence interval

```
result <- t_test(cesd ~ sex,  
                  data = HELPrc)  
result # view results  
confint(result)  
pval(result)
```

More than two levels (Analysis of variance)

Numeric and graphic summaries

```
favstats(cesd ~ substance,  
         data = HELPrc)
```

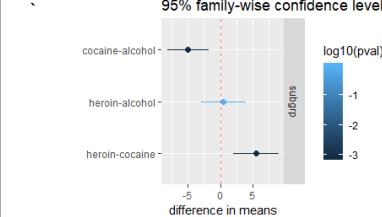
`gf_boxplot(cesd ~ substance,` `data = HELPrc)`

Fit and summarize model

`mod <- lm(age ~ substance,` `data = HELPrc)``anova(mod)`

Anova

Which differences are significant?

`mplot(TukeyHSD(mod))`

nardl Package

An R package to estimate the nonlinear cointegrating autoregressive distributed lag model

Specifying the Model

Possible syntaxes for specifying the variables in the model:

- **nardl with fixed p and q lags**

```
nardl(fod~inf,p,q,data=fod,ic="aic",maxlags = FALSE,graph = FALSE,case=3)
```

- **Auto selected lags (maxlags=TRUE)**

```
nardl(food~inf,data=fod,ic="aic",maxlags = TRUE,graph = FALSE,case=3)
```

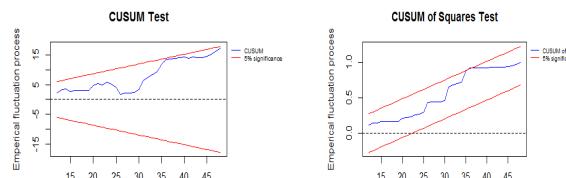
The formula:

- $y \sim x | z_1 + z_2 + \dots$
- **y** the dependent variable
- **x** the decomposed variable (this package version can't assume more than one decomposed variable)
- **$z_1 + z_2 + \dots$** independent variables
- **Data** is the dataframe
- **p** number of lags of the dependent variable
- **q** number of lags of the independent variables
- **ic** : c("aic", "bic", "ll", "R2") criteria model selection
- **maxlags** if **TRUE** auto lags selection
- **case** case number 3 for (unrestricted intercept, no trend) and 5 (unrestricted intercept, unrestricted trend), 1 2 and 4 not supported

Cusum and CusumQ plot

Cusum and CusumQ plot (graph=TRUE)

```
nardl(food~inf,data=fod,ic="aic",maxlags = TRUE,graph = TRUE,case=3)
```



Cointegration bounds test

```
pssbounds(obs, fstat, tstat = NULL, case, k)
```

pssbounds specification include:

- **Case** case number 3 for (unrestricted intercept, no trend) and 5 (unrestricted intercept, unrestricted trend), 1 2 and 4 not supported
- **fstat** represent the value of the F-statistic
- **obs** represent the number of observation
- **k** number of regressors appearing in lag levels

Example:

```
reg<-nardl(food~inf,fod,ic="aic",maxlags = TRUE,graph = TRUE,case=3)
```

```
pssbounds(case=reg$case,fstat=reg$fstat,obs=reg$obs,k=reg$k)
```

Dynamic multipliers plot

Dynamic multiplier plot

```
plotmultiplier(model, np, k, h)
```

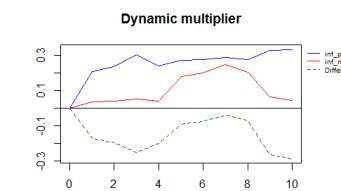
Methods and options are:

- **model** the fitted model
- **np** the selected number of lags
- **k** number of decomposed independent variables
- **h** is the horizon over which multipliers will be computed

Example

```
reg<-nardl(food~inf,p=4,q=4,fod,ic="aic",maxlags = FALSE,graph = TRUE,case=3)
```

```
plotmultiplier(reg,reg$np,1,10)
```



LM test for serial correlation

LM test for serial correlation

```
bp2(object, nlags, fill = NULL, type = c("F", "Chi2"))
```

Methods and options are:

- **object** fitted lm model
- **nlags** positive integer number of lags
- **fill** starting values for the lagged residuals in the auxiliary regression. By default 0.
- **type** Fisher or Chisquare statistics

Example :

```
reg<-nardl(food~inf,fod,ic="aic",maxlags = TRUE,graph = TRUE,case=3)
```

```
bp2(reg$fit,reg$np,fill=0,type="F")
```

Lagrange multiplier test

Lagrange multiplier test for conditional heteroscedasticity of Engle (1982), as described by Tsay (2005, pp. 101-102)

```
ArchTest(x, lags = 12, demean = FALSE)
```

Methods and options are:

- **x** numeric vector
- **lags** positive integer number of lags.
- **demean** logical: If TRUE, remove the mean before computing the test statistic.

Example :

```
reg<-nardl(food~inf,fod,ic="aic",maxlags = TRUE,graph = TRUE,case=3)
```

```
x<-reg$selresidu
```

```
nlag<-reg$np
```

```
ArchTest(x, lags=nlag)
```

pssbounds

pssbound function display the necessary critical values to conduct the Pesaran, Shin and Smith 2001 bounds test for cointegration. See <http://andyphilips.github.io/pssbounds/>.

```
pssbounds(obs, fstat, tstat = NULL, case, k)
```

Methods and options are:

- **obs** number of observations
- **fstat** value of the F-statistic
- **tstat** value of the t-statistic
- **case** case number
- **k** number of regressors appearing in lag levels

Example

```
reg<-nardl(food~inf,fod,ic="aic",maxlags = TRUE,graph = TRUE,case=3)
```

```
pssbounds(case=reg$case,fstat=reg$fstat,obs=reg$obs,k=reg$k)
```

```
# F-stat concludes I(1) and cointegrating, t-stat concludes I(0).
```

Parallel Computing :: CHEAT SHEET

Splitting :

Splitting a code by :

1. **Task** (different tasks on same data)
2. **Data** (one task on different data)

Hardware needs :

CPU (+2 cores)

RAM (shared memory vs distributed memory)

2 ideas in parallel computing :

1. Map-Reduced Models :

(distributed data; physically on different devices)

- Hadoop
- Spark

R Packages:

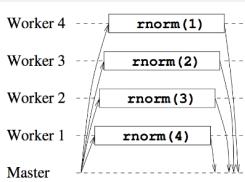
- sparklyr, iotools
- pbdr (programming with big data in R)

2. Master - Worker Models :

(M tasks on C cores; usually $1 < C \ll M$)

R Packages:

- snow, snowFT, snowfall
- foreach
- future, future.apply



Not always parallel computing:

stop/start cluster takes time

overhead (communication time b/w master and workers ; not good for repeatedly sending big data!)

Sequential vs Parallel:

```
library(microbenchmark)
microbenchmark( FUN1(...), FUN2(...),
times = 10)
```

parallel.R : core package

```
library(parallel)
ncores <- detectCores(logical=F) # physical cores
cl <- makeCluster(ncores)
clusterApply(cl, x = c(...), fun = FUN) # FUN(x,...)
stopCluster(cl)
```

Initialization of workers :

```
clusterCall(cl,FUN) # calls FUN on workers
clusterEvalQ(cl, exp) # eval an exp. on workers
## clusterEvalQ(cl, library(foo))

clusterExport(cl, varlist) # varlist on workers
## clusterExport(cl, c("mean")) where mean = 10
```

Data Chunk on workers :

1. generated on workers
`# clusterApply(cl,x, FUN) e.g. FUN(){ rnorm()}`
2. generated on master and pass to workers
`# ind <- splitIndices(200, 5)
clusterApply(cl, ind, FUN)
(-) : not efficient in Big Data : heavy`
3. chunk on workers # copy of original Data on all workers
`# clusterExport(cl, M) e.g. M is a matrix
clusterApply(cl, x, FUN) FUN contains subset M`

foreach.R : Sequential

```
library(foreach)      # by default return a list
foreach(n = rep(5,3), m = 10^(0:2)) %do% FUN(n,m)
foreach(n, .packages = "X") %do% FUN(n)
# FUN needs package X to be run
foreach(n, .export = c("Y") ) %do% FUN(n,b=Y)
# FUN needs outside object/function "Y"
foreach(n,.combine = rbind) %do% FUN(n) #row bind
foreach(n,.combine = '+') %do% FUN(n) #rbind + colSum
foreach(n,.combine = c) %do% FUN(n) # vector
foreach(n,.combine = c) %:% when(n > 2) %do% FUN(n)
```

future.R : asynchronously

```
library(future) (variables run as soon as created)
plan(multicore)

# plans : sequential, cluster, multicore, multiprocess
x %<-% mean(rnorm(100))
y %<-% mean(rnorm(100))
```

future.apply.R : parallel_apply

```
library(future.apply) (parallel _apply functions)
plan(multicore)      # can be other plans
future_apply(n,FUN),future_lapply(...),future_sapply(...)
```

foreach.R : Parallel

needs backend packages support parallel computing

- doParallel(parallel.R), doFuture (future.R), doSEQ

doParallel.R : backend of foreach

```
library(doParallel)
cl <- makeCluster(ncores)      # ncores = 2,3,..
registerDoParallel(cl)        # register the backend
foreach(...) %dopar% FUN(...)
```

doFuture.R : backend of foreach

```
library(doFuture)
registerDoFuture()

plan(cluster , workers = 3) # can be other plans
foreach(...) %dopar% FUN(...)
```

Load Balancing: for uneven task times

```
clusterApplyLB(cl,x,FUN) # not for small task time
clusterApply(cl, x = splitIndices(10,2), FUN)
library(itertools)
foreach(s=isplitVector(1:10, chunks = 2))%dopar% FUN
# e.g. FUN = sapply(s,"*",100)
future_sapply(..., future.scheduling = 1)
```

quanteda Cheat Sheet

General syntax

- **corpus_*** manage text collections/metadata
 - **tokens_*** create/modify tokenized texts
 - **dfm_*** create/modify doc-feature matrices
 - **fcm_*** work with co-occurrence matrices
 - **textstat_*** calculate text-based statistics
 - **textmodel_*** fit (un-)supervised models
 - **textplot_*** create text-based visualizations
- Consistent grammar:**
- **object()** constructor for the object type
 - **object_verb()** inputs & returns object type

Extensions

- quanteda** works well with these companion packages:
- **readtext**: an easy way to read text data
 - **spacyr**: NLP using the spaCy library
 - **quanteda.corpora**: additional text corpora
 - **stopwords**: multilingual stopword lists in R

Create a corpus from texts (corpus_*)

Read texts (txt, pdf, csv, doc, docx, json, xml)

```
my_texts <- readtext::readtext("~/link/to/path/*")
```

Construct a corpus from a character vector

```
x <- corpus(data_char_ukimmig2010, text_field = "text")
```

Explore a corpus

```
summary(data_corpus_ inaugural, n = 2)
# Corpus consisting of 58 documents, showing 2 documents:
#   Text Types Tokens Sentences Year President FirstName
# 1 1789-Washington 625 1538 23 1789 Washington George
# 2 1793-Washington 96 147 4 1793 Washington George
#
# Source: Gerhard Peters and John T. Woolley. The American Presidency Project.
# Created: Tue Jun 13 14:51:47 2017
# Notes: http://www.presidency.ucsb.edu/inaugurals.php
```

Extract or add document-level variables

```
party <- docvars(data_corpus_ inaugural, "Party")
docvars(x, "serial_number") <- 1:nDoc(x)
```

Bind or subset corpora

```
corpus(x[1:5]) + corpus(x[7:9])
corpus_subset(x, Year > 1990)
```

Change units of a corpus

```
corpus_reshape(x, to = c("sentences", "paragraphs"))
```

Segment texts on a pattern match

```
corpus_segment(x, pattern, valuetype, extract_pattern = TRUE)
```

Take a random sample of corpus texts

```
corpus_sample(x, size = 10, replace = FALSE)
```

Extract features (dfm_*; fcm_*)

Create a document-feature matrix (dfm) from a corpus

```
x <- dfm(data_corpus_ inaugural,
           tolower = TRUE, stem = FALSE, remove_punct = TRUE,
           remove = stopwords("english"))
```

head(x, n = 2, nf = 4)

```
## Document-feature matrix of: 2 documents, 4 features (41.7% sparse).
##               features
## docs      fellow-citizens senate house representatives
## 1 1789-Washington          1     1     2             2
## 2 1793-Washington          0     0     0             0
```

Create a dictionary

```
dictionary(list(negative = c("bad", "awful", "sad"),
                positive = c("good", "wonderful", "happy")))
```

Apply a dictionary

```
dfm_lookup(x, dictionary = data_dictionary_LSD2015)
```

Select features

```
dfm_select(x, dictionary = data_dictionary_LSD2015)
```

Randomly sample documents or features

```
dfm_sample(x, what = c("documents", "features"))
```

Weight or smooth the feature frequencies

```
dfm_weight(x, type = "prop") | dfm_smooth(x, smoothing = 0.5)
```

Sort or group a dfm

```
dfm_sort(x, margin = c("features", "documents", "both"))
dfm_group(x, groups = "President")
```

Combine identical dimension elements of a dfm

```
dfm_compress(x, margin = c("both", "documents", "features"))
```

Create a feature co-occurrence matrix (fcm)

```
x <- fcm(data_corpus_ inaugural, context = "window", size = 5)
fcm_compress/remove/select/toupper/tolower are also available
```

Useful additional functions

Locate keywords-in-context

```
kwic(data_corpus_ inaugural, "america*")
```

Utility functions

| | |
|--------------------------|--------------------------|
| texts(corpus) | Show texts of a corpus |
| ndoc(corpus/dfm/tokens) | Count documents/features |
| nfeat(corpus/dfm/tokens) | Count features |
| summary(corpus/dfm) | Print summary |
| head(corpus/dfm) | Return first part |
| tail(corpus/dfm) | Return last part |

Tokenize a set of texts (tokens_*)

Tokenize texts from a character vector or corpus

```
x <- tokens("Powerful tool for text analysis.",  
            remove_punct = TRUE, stem = TRUE)
```

Convert sequences into compound tokens

```
mysegs <- phrase(c("powerful", "tool", "text analysis"))  
tokens_compound(x, mysegs)
```

Select tokens

```
tokens_select(x, c("powerful", "text"), selection = "keep")
```

Create ngrams and skipgrams from tokens

```
tokens_ngrams(x, n = 1:3)
```

```
tokens_skipgrams(toks, n = 2, skip = 0:1)
```

Convert case of tokens

```
tokens_tolower(x) | tokens_topupper(x)
```

Stem the terms in an object

```
tokens_wordstem(x)
```

Calculate text statistics (textstat_*)

Tabulate feature frequencies from a dfm

```
textstat_frequency(x) | topfeatures(x)
```

Identify and score collocations from a tokenized text

```
toks <- tokens(c("quanteda is a pkg for quant text analysis",  
                 "quant text analysis is a growing field"))  
textstat_collocations(toks, size = 3, min_count = 2)
```

Calculate readability of a corpus

```
textstat_readability(data_corpus_ inaugural, measure = "Flesch")
```

Calculate lexical diversity of a dfm

```
textstat_lexdiv(x, measure = "TTR")
```

Measure distance or similarity from a dfm

```
textstat_simil(x, "2017-Trump", method = "cosine")
```

```
textstat_dist(x, "2017-Trump", margin = "features")
```

Calculate keyness statistics

```
textstat_keyness(x, target = "2017-Trump")
```

Fit text models based on a dfm (textmodel_*)

Correspondence Analysis (CA)

```
textmodel_ca(x, threads = 2, sparse = TRUE, residual_floor = 0.1)
```

Naïve Bayes classifier for texts

```
textmodel_nb(x, y = training_labels, distribution = "multinomial")
```

Wordscores text model

```
refscores <- c(seq(-1.5, 1.5, .75), NA)
```

```
textmodel_wordscores(data_dfm_lbgexample, refscores)
```

Wordfish Poisson scaling model

```
textmodel_wordfish(dfm(data_corpus_irishbudget2010), dir = c(6,5))
```

Textmodel methods: predict(), coef(), summary(), print()

Plot features or models (textplot_*)

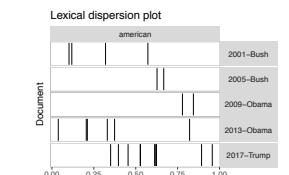
Plot features as a wordcloud

```
data_corpus_ inaugural %>%  
corpus_subset(President == "Obama") %>%  
dfm(remove = stopwords("english")) %>%  
textplot_wordcloud()
```



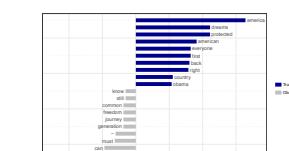
Plot the dispersion of key word(s)

```
data_corpus_ inaugural %>%  
corpus_subset(Year > 1945) %>%  
kwic("american") %>%  
textplot_xray()
```



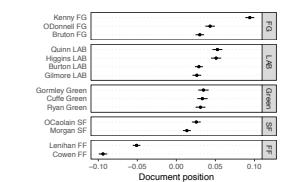
Plot word keyness

```
data_corpus_ inaugural %>%  
corpus_subset(President %in%  
              c("Obama", "Trump")) %>%  
dfm(groups = "President",  
     remove = stopwords("english")) %>%  
textstat_keyness(target = "Trump") %>%  
textplot_keyness()
```



Plot Wordfish, Wordscores or CA models

```
textplot_scale1d(scaling_model,  
                  groups = party,  
                  margin = "documents")
```



Convert dfm to a non-quanteda format

```
convert(x, to = c("lda", "tm", "stm", "austin", "topicmodels",  
               "lsa", "matrix", "data.frame"))
```

randomizr:: CHEAT SHEET

Two Arm Trials

Simple random assignment is like flipping coins for each unit separately.

```
simple_ra(N = 100, prob = 0.5)
```

Complete random assignment allocates a fixed number of units to each condition.

```
complete_ra(N = 100, m = 50)
complete_ra(N = 100, prob = 0.5)
```

Block random assignment conducts complete random assignment separately for groups of units.

```
blocks <- rep(c("A", "B", "C"),
              c(50, 100, 200))

# defaults to half of each block
block_ra(blocks = blocks)

# can change with block_m
block_ra(blocks = blocks,
         block_m = c(20, 30, 40))
```

Cluster random assignment allocates whole groups of units to conditions together.

```
clusters <- rep(letters, times = 1:26)
cluster_ra(clusters = clusters)
```

Block and cluster random assignment conducts cluster random assignment separately for groups of clusters.

```
clusters <- rep(letters, times = 1:26)
blocks <- rep(paste0("block_", 1:5),
             c(15, 40, 65, 90, 141))
block_and_cluster_ra(blocks = blocks,
                     clusters = clusters)
```

randomizr is part of the DeclareDesign suite of packages for designing, implementing, and analyzing social science research designs.

Multi Arm Trials

Set the number of arms with num_arms or with conditions.

```
complete_ra(N = 100, num_arms = 3)
complete_ra(N = 100, conditions = c("control",
                                    "placebo", "treatment"))
```

The *_each arguments in randomizr functions specify design parameters for each arm separately.

```
complete_ra(N = 100, m_each = c(20, 30, 50))
complete_ra(N = 100,
           prob_each = c(0.2, 0.3, 0.5))
```

If the design is the **same** for all blocks, use prob_each:

```
blocks <- rep(c("A", "B", "C"),
              c(50, 100, 200))
block_ra(blocks = blocks,
         prob_each = c(.1, .1, .8))
```

If the design is **different** in different blocks, use block_m_each or block_prob_each:

```
block_m_each <- rbind(c(10, 20, 20),
                       c(30, 50, 20),
                       c(50, 75, 75))
block_ra(blocks = blocks,
         block_m_each = block_m_each)

block_prob_each <- rbind(c(.1, .1, .8),
                         c(.2, .2, .6),
                         c(.3, .3, .4))
block_ra(blocks = blocks,
         block_prob_each = block_prob_each)
```

If conditions is numeric, the output will be **numeric**.

If conditions is not numeric, the output will be a **factor** with levels in the order provided to conditions.

```
complete_ra(N = 100, conditions = -2:2)
complete_ra(N = 100, conditions = c("A", "B"))
```

Declaration

Learn about assignment procedures by “declaring” them with declare_ra()

```
declaration <-
  declare_ra(N = 100, m_each = c(30, 30, 40))
```

```
declaration # print design information
```

Conduct a random assignment:

```
conduct_ra(declaration)
```

Obtain observed condition probabilities (useful for inverse probability weighting if probabilities of assignment are not constant)

```
Z <- conduct_ra(declaration)
obtain_condition_probabilities(declaration, Z)
```

Sampling

All assignment functions have sampling analogues: Sampling is identical to a two arm trial where the treatment group is sampled.

Assignment

```
simple_ra()
```

```
complete_ra()
```

```
block_ra()
```

```
cluster_ra()
```

```
block_and_cluster_ra()
```

```
declare_ra()
```

```
conduct_ra()
```

Sampling

```
simple_rs()
```

```
complete_rs()
```

```
strata_rs()
```

```
cluster_rs()
```

```
strata_and_cluster_rs()
```

```
declare_rs()
```

```
draw_rs()
```

Stata

A Stata version of randomizr is available, with the same arguments but different syntax:

```
ssc install randomizr
set obs 100
complete_ra, m(50)
```

Basic Regular Expressions in R

Cheat Sheet

Character Classes

| | |
|-------------------|---|
| [:digit:]] or \d | Digits; [0-9] |
| \D | Non-digits; [^0-9] |
| [:lower:]] | Lower-case letters; [a-z] |
| [:upper:]] | Upper-case letters; [A-Z] |
| [:alpha:]] | Alphabetic characters; [A-z] |
| [:alnum:]] | Alphanumeric characters [A-z0-9] |
| \w | Word characters; [A-z0-9_] |
| \W | Non-word characters |
| [:xdigit:]] or \x | Hexadecimal digits; [0-9A-Fa-f] |
| [:blank:]] | Space and tab |
| [:space:]] or \s | Space, tab, vertical tab, newline, form feed, carriage return |
| \S | Not space; [^[:space:]] |
| [:punct:]] | Punctuation characters; !#\$%&'()*,-./;:<>?@[]^_`{ ~ |
| [:graph:]] | Graphical characters; [:alnum:][:punct:]] |
| [:print:]] | Printable characters; [:alnum:][:punct:]\s] |
| [:cntrl:]] or \c | Control characters; \n, \r etc. |

Special Metacharacters

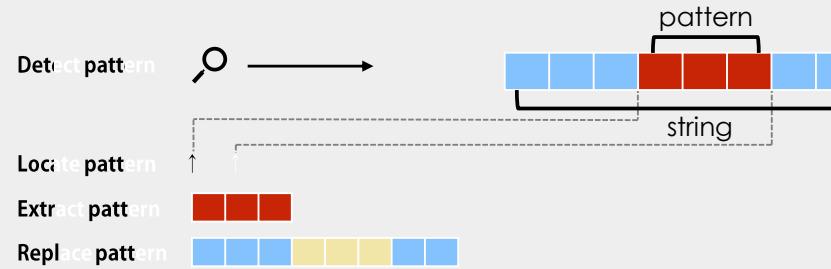
| | |
|----|-----------------|
| \n | New line |
| \r | Carriage return |
| \t | Tab |
| \v | Vertical tab |
| \f | Form feed |

Lookarounds and Conditionals*

| | |
|----------------|---|
| (?=) | Lookahead (requires PERL = TRUE), e.g. (?=yx): position followed by 'xy' |
| (?!) | Negative lookahead (PERL = TRUE); position NOT followed by pattern |
| (?=<=) | Lookbehind (PERL = TRUE), e.g. (?<=yx): position following 'xy' |
| (?<=) | Negative lookbehind (PERL = TRUE); position NOT following pattern |
| ?(if)then | If-then-condition (PERL = TRUE); use lookahead, optional char. etc in if-clause |
| ?(if)then else | If-then-else-condition (PERL = TRUE) |

*see, e.g. <http://www.regular-expressions.info/lookaround.html>
<http://www.regular-expressions.info/conditional.html>

Functions for Pattern Matching



```
> string <- c("Hipopopotamus", "Rhymenoceros", "time for bottomless lyrics")
> pattern <- "t.m"
```

Detect Patterns

grep(pattern, string)

```
[1] 1 3
```

grep(pattern, string, value = TRUE)

```
[1] "Hipopopotamus"
[2] "time for bottomless lyrics"
```

grepl(pattern, string)

```
[1] TRUE FALSE TRUE
```

stringr::str_detect(string, pattern)

```
[1] TRUE FALSE TRUE
```

Locate Patterns

regexpr(pattern, string)

find starting position and length of first match

gregexpr(pattern, string)

find starting position and length of all matches

stringr::str_locate(string, pattern)

find starting and end position of first match

stringr::str_locate_all(string, pattern)

find starting and end position of all matches

Split a String using a Pattern

strsplit(string, pattern) or stringr::str_split(string, pattern)

Character Classes and Groups

. Any character except \n

| Or, e.g. (a|b)

[...] List permitted characters, e.g. [abc]

[a-z] Specify character ranges

[^...] List excluded characters

(...) Grouping, enables back referencing using \\N where N is an integer

Anchors

^ Start of the string

\$ End of the string

\b Empty string at either edge of a word

\B NOT the edge of a word

\B Beginning of a word

\B End of a word

Quantifiers

* Matches at least 0 times

+ Matches at least 1 time

? Matches at most 1 time; optional string

{n} Matches exactly n times

{n,} Matches at least n times

{n,m} Matches between n and m times

General Modes

By default R uses *extended regular expressions*. You can switch to *PCRE regular expressions* using PERL = TRUE for base or by wrapping patterns with perl() for stringr.

All functions can be used with literal searches using fixed = TRUE for base or by wrapping patterns with fixed() for stringr.

All base functions can be made case insensitive by specifying ignore.cases = TRUE.

Escaping Characters

Metacharacters (. * + etc.) can be used as literal characters by escaping them. Characters can be escaped using \\ or by enclosing them in \\Q...\\E.

Case Conversions

Regular expressions can be made case insensitive using (?i). In backreferences, the strings can be converted to lower or upper case using \\L or \\U (e.g. \\L\\1). This requires PERL = TRUE.

Greedy Matching

By default the asterisk * is greedy, i.e. it always matches the longest possible string. It can be used in lazy mode by adding ?, i.e. *?.

Greedy mode can be turned off using (?U). This switches the syntax, so that (?U)a* is lazy and (?U)a*? is greedy.

Note

Regular expressions can conveniently be created using e.g. the packages **rex** or **rebus**.

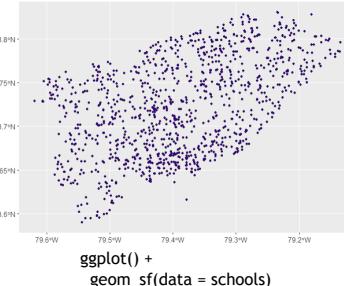
Spatial manipulation with sf: : CHEAT SHEET

The sf package provides a set of tools for working with geospatial vectors, i.e. points, lines, polygons, etc.



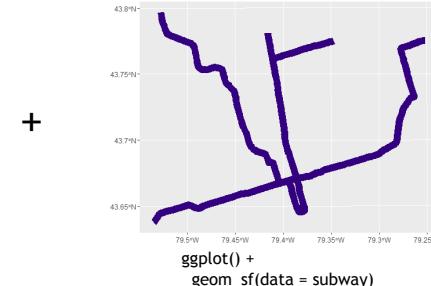
Geometric confirmation

- ☒ ☐ st_contains(x, y, ...) Identifies if x is within y (i.e. point within polygon)
- ☒ ☒ st_covered_by(x, y, ...) Identifies if x is completely within y (i.e. polygon completely within polygon)
- ☒ ☒ st_covers(x, y, ...) Identifies if any point from x is outside of y (i.e. polygon outside polygon)
- ☒ ☒ st_crosses(x, y, ...) Identifies if any geometry of x have commonalities with y
- ☒ ☒ st_disjoint(x, y, ...) Identifies when geometries from x do not share space with y
- ☒ ☒ st_equals(x, y, ...) Identifies if x and y share the same geometry
- ☒ ☒ st_intersects(x, y, ...) Identifies if x and y geometry share any space
- ☒ ☒ st_overlaps(x, y, ...) Identifies if geometries of x and y share space, are of the same dimension, but are not completely contained by each other
- ☒ ☒ st_touches(x, y, ...) Identifies if geometries of x and y share a common point but their interiors do not intersect
- ☒ ☒ st_within(x, y, ...) Identifies if x is in a specified distance to y

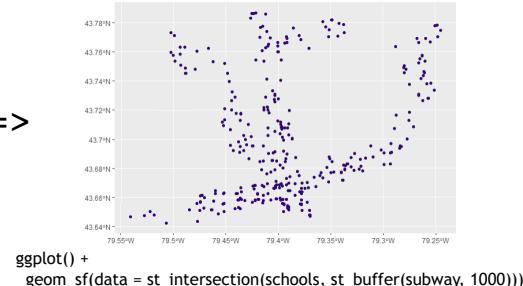


Geometric operations

- ☒ ☐ st_boundary(x) Creates a polygon that encompasses the full extent of the geometry
- ☒ ☒ st_buffer(x, dist, nQuadSegs) Creates a polygon covering all points of the geometry within a given distance
- ☒ ☒ st_centroid(x, ..., of_largest_polygon) Creates a point at the geometric centre of the geometry
- ☒ ⇒ ☐ st_convex_hull(x) Creates geometry that represents the minimum convex geometry of x
- ☒ ⇒ ☐ st_line_merge(x) Creates linestring geometry from sewing multi linestring geometry together
- ☒ ⇒ ☐ st_node(x) Creates nodes on overlapping geometry where nodes do not exist
- ☒ ☐ st_point_on_surface(x) Creates a point that is guaranteed to fall on the surface of the geometry
- ☒ ⇒ ☐ st_polygonize(x) Creates polygon geometry from linestring geometry
- ☒ ⇒ ☐ st_segmentize(x, dfMaxLength, ...) Creates linestring geometry from x based on a specified length
- ☒ ⇒ ☐ st_simplify(x, preserveTopology, dTolerance) Creates a simplified version of the geometry based on a specified tolerance



=>



Geometry creation

- ☒ ⇒ ☐ st_triangulate(x, dTolerance, bOnlyEdges) Creates polygon geometry as triangles from point geometry
- ☒ ⇒ ☐ st_voronoi(x, envelope, dTolerance, bOnlyEdges) Creates polygon geometry covering the envelope of x, with x at the centre of the geometry
- ☒ ☐ st_point(x, c(numeric vector), dim = "XYZ") Creating point geometry from numeric values
- ☒ ☐ st_multipoint(x = matrix(numeric values in rows), dim = "XYZ") Creating multi point geometry from numeric values
- ☒ ☐ st_linestring(x = matrix(numeric values in rows), dim = "XYZ") Creating linestring geometry from numeric values
- ☒ ☐ st_multilinestring(x = list(numeric matrices in rows), dim = "XYZ") Creating multi linestring geometry from numeric values
- ☒ ☐ st_polygon(x = list(numeric matrices in rows), dim = "XYZ") Creating polygon geometry from numeric values
- ☒ ☐ st_multipolygon(x = list(numeric matrices in rows), dim = "XYZ") Creating multi polygon geometry from numeric values

Spatial manipulation with sf: : CHEAT SHEET

The sf package provides a set of tools for working with geospatial vectors, i.e. points, lines, polygons, etc.



Geometry operations

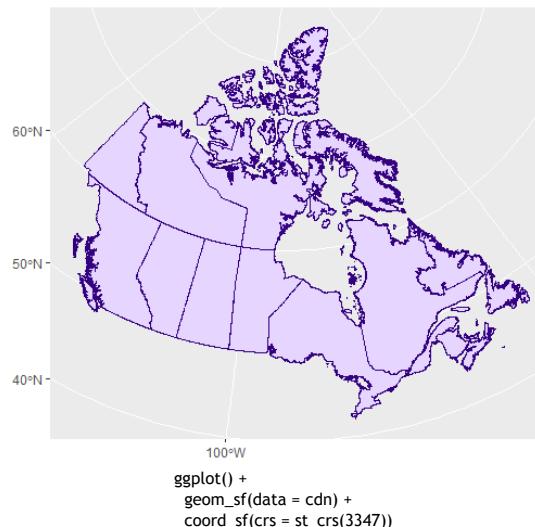
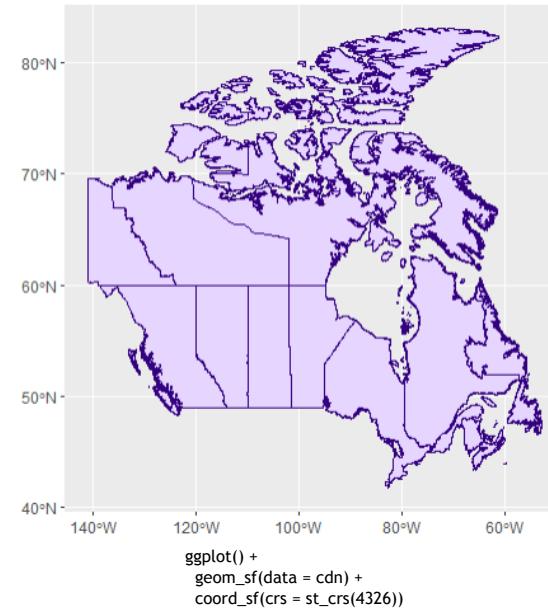
- ↳ ↳ `st_contains(x, y, ...)` Identifies if x is within y (i.e. point within polygon)
- ↳ ↳ `st_crop(x, y, ..., xmin, ymin, xmax, ymax)` Creates geometry of x that intersects a specified rectangle
- ↳ ↳ `st_difference(x, y)` Creates geometry from x that does not intersect with y
- ↳ ↳ `st_intersection(x, y)` Creates geometry of the shared portion of x and y
- ↳ ↳ `st_sym_difference(x, y)` Creates geometry representing portions of x and y that do not intersect
- ↳ ↳ `st_snap(x, y, tolerance)` Snap nodes from geometry x to geometry y
- ↳ ↳ `st_union(x, y, ..., by_feature)` Creates multiple geometries into a single geometry, consisting of all geometry elements

Geometric measurement

- `st_area(x)` Calculate the surface area of a polygon geometry based on the current coordinate reference system
- `st_distance(x, y, ..., dist_fun, by_element, which)` Calculates the 2D distance between x and y based on the current coordinate system
- `st_length(x)` Calculates the 2D length of a geometry based on the current coordinate system

Misc operations

- `st_as_sf(x, ...)` Create a sf object from a non-geospatial tabular data frame
- `st_cast(x, to, ...)` Change x geometry to a different geometry type
- `st_coordinates(x, ...)` Creates a matrix of coordinate values from x
- `st_crs(x, ...)` Identifies the coordinate reference system of x
- `st_join(x, y, join, FUN, suffix, ...)` Performs a spatial left or inner join between x and y
- `st_make_grid(x, cellsize, offset, n, crs, what)` Creates rectangular grid geometry over the bounding box of x
- `st_nearest_feature(x, y)` Creates an index of the closest feature between x and y
- `st_nearest_points(x, y, ...)` Returns the closest point between x and y
- `st_read(dsn, layer, ...)` Read file or database vector dataset as a sf object
- `st_transform(x, crs, ...)` Convert coordinates of x to a different coordinate reference system



Data & Variable Transformation with sjmisc Cheat Sheet



sjmisc complements dplyr, and helps with data transformation tasks and recoding *variables*.

sjmisc works together seamlessly with dplyr and pipes. All functions are designed to support labelled data.



Design Philosophy

The design of sjmisc functions follows the tidyverse-approach: first argument is always the data (either a *data frame* or *vector*), followed by variable names to be processed by the functions.

The returned object for each function *equals the type of the data-argument*.

Vector input

- If the data-argument is a *vector*, functions return a *vector*.



`rec(mtcars$carb, rec = "1,2=1; 3,4=2; else=3")`

Data frame input

- If the data-argument is a *data frame*, functions return a *data frame*.



`rec(mtcars, carb, rec = "1,2=1; 3,4=2; else=3")`

The ...-ellipses Argument

Apply functions to a single variable, selected variables or to a complete data frame.

Variable selection is powered by dplyr `select()`: Separate variables with comma, or use dplyr select-helpers to select variables, e.g. `?rec`:

`rec(mtcars, one_of(c("gear", "carb")))`
`rec = "min:3=1; 4:max=2"`

`rec(mtcars, gear, carb, rec = "min:3=1; 4:max=2")`

Descriptives and Summaries

Most of the sjmisc functions (including recode-functions) also work on grouped data frames:

```
library(dplyr)
etc %>%
  group_by(e16sex, c172code) %>%
  frq(e42dep)
```

Frequency Tables

`frq(x, ..., sort.frq = c("none", "asc", "desc"), weight.by = NULL, auto.grp = NUL)`

Print frequency tables of (labelled) vectors. Uses variable labels as table header.

```
data(etc); frq(etc, e42dep, c161sex)
```

Use this data set in examples!

```
flat_table(data, ..., margin = c("counts", "cell", "row", "col"), digits = 2, show.values = FALSE)
```

Print contingency tables of (labelled) vectors. Uses value labels.

```
flat_table(etc, e42dep, c172code, e16sex)
```

`count_na(x, ...)`

Print frequency table of tagged NA values.

```
library(haven); x <- labelled(c(1:3, tagged_na("a", "a", "z")), labels = c("Refused" = tagged_na("a"), "N/A" = tagged_na("z")))
count_na(x)
```

Descriptive Summary

`descr(x, ..., max.length = NULL)`

Descriptive summary of data frames, including variable labels in output.

```
descri(etc, contains("cop"), max.length = 20)
```

Finding Variables in a Data Frame

Use `find_var()` to search for variables by names, value or variable labels. Returns vector/data frame.

```
# variables with "cop" in n
find_var(etc, pattern = "cop", out = "df")
# variables with "level" in names and value labels
find_var(etc, "level", search = "name_value")
```

Recode and Transform Variables

Recode functions add a *suffix* to new variables, so original variables are preserved.

By default, original input data frame and new created variables are returned. Use `append = FALSE` to return the recoded variables only.

```
rec(x, ..., rec, as.num = TRUE, var.label = NULL, val.labels = NULL, append = TRUE, suffix = "_r")
```

Recode values, return result as numeric, character or categorical (factor).

```
rec(mtcars, carb, rec = "1,2=1; 3,4=2; else=3")
```

```
dicho(x, ..., dich.by = "median", as.num = FALSE, var.label = NULL, val.labels = NULL, append = TRUE, suffix = "_d")
```

Dichotomise variable by median, mean or specific value.

```
dicho(mtcars, disp)
```

```
split_var(x, ..., n, as.num = FALSE, val.labels = NULL, var.label = NULL, inclusive = FALSE, append = TRUE, suffix = "_g")
```

Split variable into equal sized groups. Unlike `dplyr::ntile()`, does not split original categories into different values (see examples in `?split_var`).

```
split_var(mtcars, mpg, disp, n = 3)
```

```
group_var(x, ..., size = 5, as.num = TRUE, right.interval = FALSE, n = 30, append = TRUE, suffix = "_gr")
```

Split variable into groups with equal value range, or into a max. # of groups (value range per group is adjusted to match # of groups).

```
group_var(mtcars, mpg, disp, size = 5)
group_var(mtcars, mpg, size = "auto", n = 4)
```

```
std(x, ..., robust = "sd", include.fac = FALSE, append = TRUE, suffix = "_z")
```

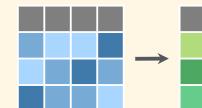
Z-standardise variables. Also `center()`.

```
std(etc, e17age, c160age)
```

```
recode_to(x, ..., lowest = 0, highest = -1, append = TRUE, suffix = "_r0")
Shift ("renumber") c
recode_to(mtcars$gear)
```

Summarise Variables and Cases

The summary functions mostly mimic base R equivalents, but are designed to work together with pipes and dplyr.



```
row_sums(x, ..., na.rm = TRUE, var = "rowsums", append = FALSE)
```

Row sums of data frames.

```
row_sums(etc, c82cop1:c90cop9)
```

```
row_means(x, ..., n, var = "rowmeans", append = FALSE)
```

Row means, for at least `n` valid (non-NA) values.

```
row_means(etc, c82cop1:c90cop9, n = 7)
```

```
row_count(x, ..., count, var = "rowcount", append = FALSE)
```

Row-wise count # of values in data frames. Also `col_count()`.

```
row_count(etc, c82cop1:c90cop9, count = 2)
```

Other Useful Functions

- `add_columns()` and `replace_columns()` to combine data frames, but either replace or preserve existing columns.
- `set_na()` and `replace_na()` to convert regular into missing values, or vice versa. `replace_na()` also replaces specific `tagged NA` values only.
- `remove_var()` and `var_rename()` to remove variables from data frames, or rename variables.
- `group_str()` to group similar string values. Useful for variables with similar, but not identically spelled string values th:
- `merge_df()` to full join data frames and preserve value and variable labels.
- `to_long()` to gather multiple columns in data frames from wide into long format.

Use with %>% and dplyr

```
# use sjmisc-functions in pipes
mtcars %>% select(gear, carb) %>%
  rec(rec = "min:3=1; 4:max=2")
```

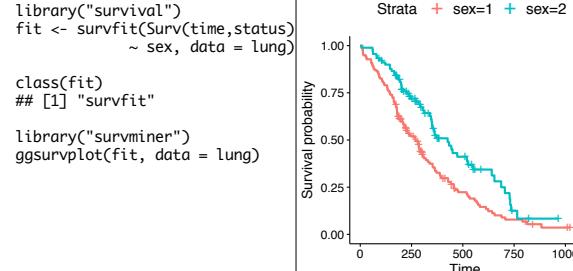
```
# use sjmisc-function inside mutate
mtcars %>% select(gear, carb) %>% mutate(
  carb2 = rec(carb, re
  gear2 = rec(gear, re
```

Creating Survival Plots

Informative and Elegant with *survminer*

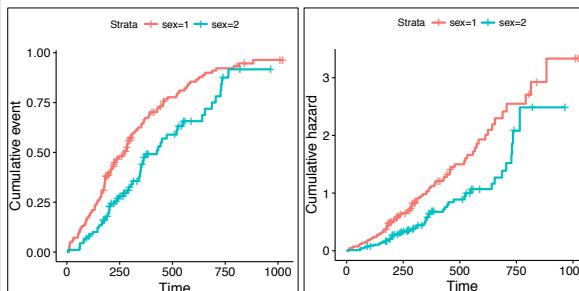
Survival Curves

The `ggsurvplot()` function creates `ggplot2` plots from `survfit` objects.



Use the `fun` argument to set the transformation of the survival curve. E.g. "`event`" for cumulative events, "`cumhaz`" for the cumulative hazard function or "`pct`" for survival probability in percentage.

```
ggsurvplot(fit, data = lung, fun = "event")
ggsurvplot(fit, data = lung, fun = "cumhaz")
```



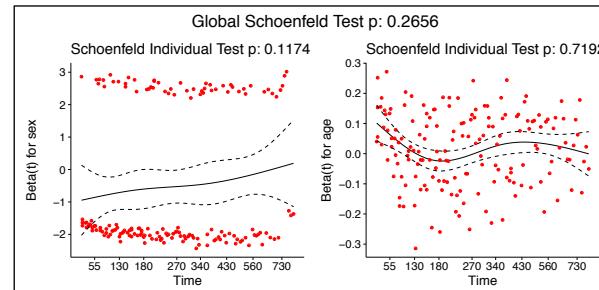
With lots of graphical parameters you have full control over look and feel of the survival plots; position and content of the legend; additional annotations like p-value, title, subtitle.

```
ggsurvplot(fit, data = lung,
           conf.int = TRUE,
           pval = TRUE,
           fun = "pct",
           risk.table = TRUE,
           size = 1,
           linetype = "strata",
           palette = c("#E7B800",
                      "#2E9FDF"),
           legend = "bottom",
           legend.title = "Sex",
           legend.labs = c("Male",
                         "Female"))
```

Diagnostics of Cox Model

The function `cox.zph()` from `survival` package may be used to test the proportional hazards assumption for a Cox regression model fit. The graphical verification of this assumption may be performed with the function `ggcoxzph()` from the `survminer` package. For each covariate it produces plots with scaled Schoenfeld residuals against the time.

```
library("survival")
fit <- coxph(Surv(time, status) ~ sex + age, data = lung)
ftest <- cox.zph(fit)
ftest
##          rho chisq      p
## sex     0.1236 2.452 0.117
## age    -0.0275 0.129 0.719
## GLOBAL    NA 2.651 0.266
library("survminer")
ggcoxzph(ftest)
```



The function `ggcoxdiagnostics()` plots different types of residuals as a function of time, linear predictor or observation id. The type of residual is selected with `type` argument. Possible values are "martingale", "deviance", "score", "schoenfeld", "dfbeta", "dfbetas", and "scaledsch".

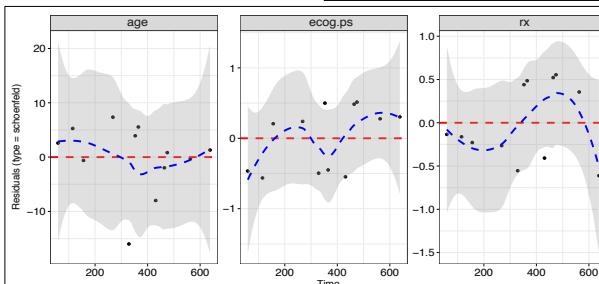
The `ox.scale` argument defines what shall be plotted on the OX axis. Possible values are "linear.predictions", "observation.id", "time".

Logical arguments `hline` and `sline` may be used to add horizontal line or smooth line to the plot.

```
library("survival")
library("survminer")
fit <- coxph(Surv(time, status) ~ sex + age, data = lung)
```

```
ggcoxdiagnostics(fit,
                  type = "deviance",
                  ox.scale = "linear.predictions")
```

```
ggcoxdiagnostics(fit,
                  type = "schoenfeld",
                  ox.scale = "time")
```



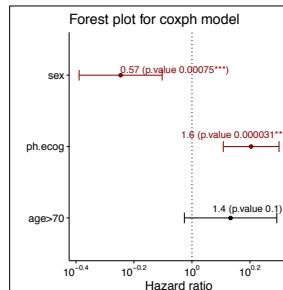
Summary of Cox Model

The function `ggforest()` from the `survminer` package creates a forest plot for a Cox regression model fit. Hazard ratio estimates along with confidence intervals and p-values are plotted for each variable.

```
library("survival")
library("survminer")
lung$age <- ifelse(lung$age > 70, ">70", "<= 70")
fit <- coxph(Surv(time, status) ~ sex + ph.ecog + age, data = lung)
```

```
## Call:
## coxph(formula = Surv(time, status) ~ sex+ph.ecog+age, data=lung)
##
##            coef exp(coef) se(coef)      z      p
## sex       -0.567   0.567   0.168 -3.37 0.00075
## ph.ecog    0.470   1.600   0.113  4.16 3.1e-05
## age>70    0.307   1.359   0.187  1.64 0.10175
##
## Likelihood ratio test=31.6 on
## n= 227, number of events= 164
```

```
ggforest(fit)
```



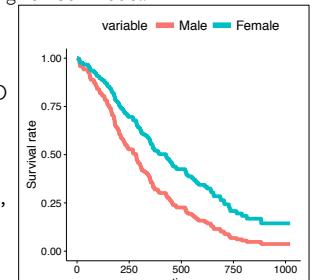
The function `ggcoadjustedcurves()` from the `survminer` package plots Adjusted Survival Curves for Cox Proportional Hazards Model. Adjusted Survival Curves show how a selected factor influences survival estimated from a Cox model.

Note that these curves differ from Kaplan Meier estimates since they present expected survival based on given Cox model.

```
library("survival")
library("survminer")
lung$sex <- ifelse(lung$sex == 1,
                   "Male", "Female")

fit <- coxph(Surv(time, status) ~ sex + ph.ecog + age,
             data = lung)

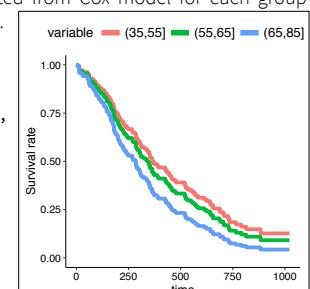
ggcoadjustedcurves(fit, data=lung,
                    variable=lung$sex)
```



Note that it is not necessary to include the grouping factor in the Cox model. Survival curves are estimated from Cox model for each group defined by the factor independently.

```
lung$age3 <- cut(lung$age,
                  c(35, 55, 65, 85))

ggcoadjustedcurves(fit, data=lung,
                    variable=lung$age3)
```



R Syntax Comparison :: CHEAT SHEET

Dollar sign syntax

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

SUMMARY STATISTICS:

one continuous variable:
`mean(mtcars$mpg)`

one categorical variable:
`table(mtcars$cyl)`

two categorical variables:
`table(mtcars$cyl, mtcars$am)`

one continuous, one categorical:
`mean(mtcars$mpg[mtcars$cyl==4])`
`mean(mtcars$mpg[mtcars$cyl==6])`
`mean(mtcars$mpg[mtcars$cyl==8])`

PLOTTING:

one continuous variable:
`hist(mtcars$disp)`

`boxplot(mtcars$disp)`

one categorical variable:
`barplot(table(mtcars$cyl))`

two continuous variables:
`plot(mtcars$disp, mtcars$mpg)`

two categorical variables:
`mosaicplot(table(mtcars$am, mtcars$cyl))`

one continuous, one categorical:
`histogram(mtcars$disp[mtcars$cyl==4])`
`histogram(mtcars$disp[mtcars$cyl==6])`
`histogram(mtcars$disp[mtcars$cyl==8])`

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

WRANGLING:

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

making a new variable:
`mtcars$efficient[mtcars$mpg>30] <- TRUE`
`mtcars$efficient[mtcars$mpg<30] <- FALSE`

Formula syntax

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

SUMMARY STATISTICS:

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

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

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

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

tilde

PLOTTING:

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

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

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

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

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

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

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

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

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

Tidyverse syntax

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

SUMMARY STATISTICS:

one continuous variable:
`mtcars %>% dplyr::summarize(mean(mpg))`

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

the pipe

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

one continuous, one categorical:
`mtcars %>% dplyr::group_by(cyl) %>%
dplyr::summarize(mean(mpg))`

PLOTTING:
one continuous variable:
`ggplot2::qplot(x=mpg, data=mtcars, geom = "histogram")`

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

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

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

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

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

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

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

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

R Syntax Comparison :: CHEAT SHEET

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

The three most prevalent R syntaxes are:

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

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

Internet research tip:

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

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



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

Even more ways to say the same thing

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

quickplot

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

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

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

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

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

ggplot

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

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

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

plus adds layers

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

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

ggformula

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

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

formulas in base plots

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

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

The teachR's :: CHEAT SHEET

Getting ready to teach some R? Use our cheat sheet to **prepare, teach and debrief**



DRAFT v0.1

Before the course (design)

Use these to prepare your lecture/course:

Who are your learners? (Persona Analysis)
(change according to requirements...[1])



The R novice

Background: some statistics, some programming

Prior knowledge: basic R course, base R syntax

Goals: understand tidy concepts,
expose to tidyverse practices

Special needs: First successes, mitigate fears, encourage learning

The R "false expert"

Background: working with R for some time, but doesn't keep-up

Prior knowledge: been using base R syntax, loops, and functions

Goals: strengthen tidyverse familiarity, apply dplyr workflow

Special needs: switch from obsolete methods to state-of-the-art R

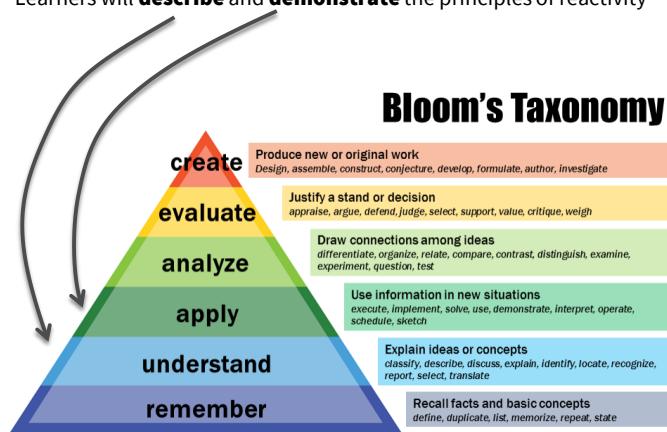
Define goals using Bloom's Taxonomy [2]

Design your classes to move your learners "up the pyramid"

Keep "realistic goals" for each persona



For example (R shiny - novice):
Learners will **describe** and **demonstrate** the principles of reactivity



Additional sources

[1] Dreyfus, Stuart E., and Hubert L. Dreyfus. *A five stage model of the mental activities involved in directed skill acquisition*. No. ORC-80-2. California Univ Berkeley Operations Research Center, 1980.

[2] Content downloaded from <https://cft.vanderbilt.edu/guides-sub-pages/blooms-taxonomy/>
(CC-BY-SA Vanderbilt University Center for Teaching)

After the course (learn/improve)

Make sure you make the most to improve your next lecture



Use feedback to understand what went well, and what you need to improve.



Measure the time each lecture takes you (or where did you get to), so that next time your time estimates will be better

Useful tips and tricks

Useful tips for preparations



Use github to upload course materials



RMarkdown for exercises

Recommended reading materials/references for R courses:

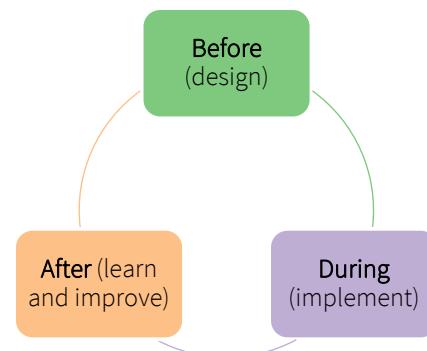
R for Data Science / Garrett Grolemund and Hadley Wickham
(r4ds.had.co.nz)

Advanced R / Hadley Wickham (adv-r.had.co.nz)

RStudio **Cheat sheets**:

<https://www.rstudio.com/resources/cheatsheets/>

Iterative work flow





Class Agnostic Time Series with tsbox :: CHEAT SHEET

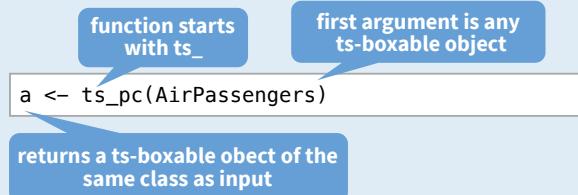
Basics

IDEA

tsbox provides a time series toolkit which:

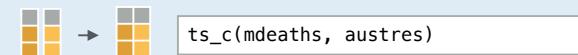
1. works identically with most time series **classes**
2. handles regular and irregular **frequencies**
3. **converts** between classes and frequencies

Most functions in tsbox have the same structure:

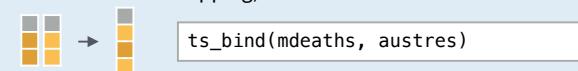


COMBINE TIME SERIES

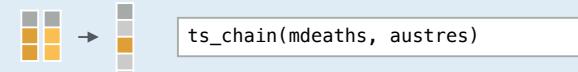
collect time series of **all classes** and **frequencies** as multiple time series



combine time series to a new, single time series (first series wins if overlapping)

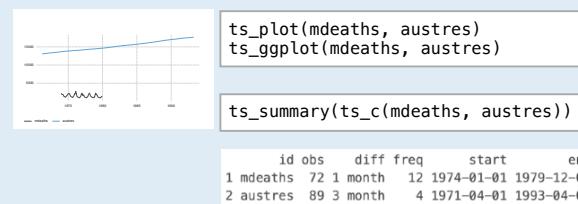


like ts_bind, but extra- and extrapolate, using growth rates



PLOT AND SUMMARIZE

Plot time series of **all classes** and **frequencies**



Helper Functions

Transform time series of **all classes** and **frequencies**

TRANSFORM



ts_trend(): Trend estimation based on loess

`ts_trend(fdeaths)`



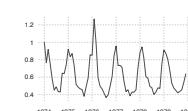
ts_pc(), **ts_pcy()**, **ts_pca()**, **ts_diff()**, **ts_difff()**: (annualized) Percentage change rates or differences to previous period, year

`ts_pc(fdeaths)`



ts_scale(): normalize mean and variance

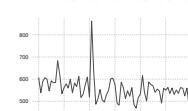
`ts_scale(fdeaths)`



ts_index(): Index, based on levels

ts_compound(): Index, based on growth rates

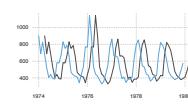
`ts_index(fdeaths, base = 1976)`



ts_seas(): seasonal adjustment using X-13

`ts_seas(fdeaths)`

SPAN AND FREQUENCY



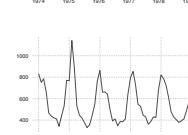
ts_lag(): Lag or lead of time series

`ts_lag(fdeaths, 4)`



ts_frequency(): convert to frequency

`ts_frequency(fdeaths, "year")`



ts_span(): filter time series for a time span.

`ts_span(fdeaths, "1976-01-01")`

`ts_span(fdeaths, "-5 year")`

Class Conversion

tsbox is built around a set of converters, which convert time series of the following **supported classes** to each other:

| converter function | ts-boxable class |
|---|------------------------|
| <code>ts_ts()</code> | ts, mts |
| <code>ts_data.frame()</code> , <code>ts_df()</code> | data.frame |
| <code>ts_data.table()</code> , <code>ts_dt()</code> | data.table |
| <code>ts_tbl()</code> | df_tbl, "tibble" |
| <code>ts_xts()</code> | xts |
| <code>ts_zoo()</code> | zoo |
| <code>ts_tibbletime()</code> | tibbletime |
| <code>ts_timeSeries()</code> | timeSeries |
| <code>ts_tsibble()</code> | tsibble |
| <code>ts_tslist()</code> | a list with ts objects |

Time Series in data frames

LONG STRUCTURE

Default structure to store multiple time series in long data frames (or data tables, or tibbles)

`ts_df(ts_c(fdeaths, mdeaths))`

| id | time | value |
|---------|------------|-------|
| fdeaths | 1974-01-01 | 901 |
| fdeaths | 1974-02-01 | 689 |
| fdeaths | 1974-03-01 | 827 |
| ... | ... | ... |

AUTO-DETECT COLUMN NAMES

tsbox auto-detects a `value`-, a `time`- and zero, one or several `id`-columns. Alternatively, the `time`- and the `value`-column can be explicitly named `time` and `value`.

`ts_default()`: standardize column names in data frames

RESHAPE

`ts_wide()`: convert default long structure to wide

`ts_long()`: convert wide structure to default long

USE WITH PIPE

tsbox plays well with tibbles and with `%>%`, so it can be easily integrated into a dplyr/pipe workflow

```
library(dplyr)
ts_c(fdeaths, mdeaths) %>%
  ts_tbl() %>%
  ts_trend() %>%
  ts_pc()
```

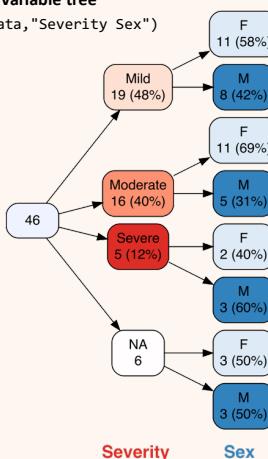
pass return value as first argument to the next function

vtree cheatsheet

Basics

Draw a basic variable tree

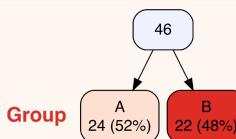
```
vtree(FakeData, "Severity Sex")
```



Severity Sex

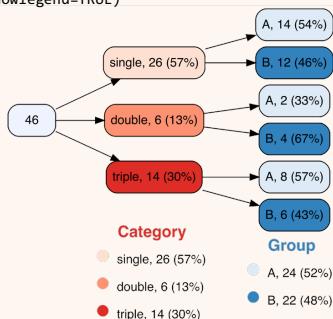
Draw the tree vertically

```
vtree(FakeData, "Group", horiz=FALSE)
```



Display a legend

```
vtree(FakeData, "Category Group", sameline=TRUE, showlegend=TRUE)
```



Variable specification

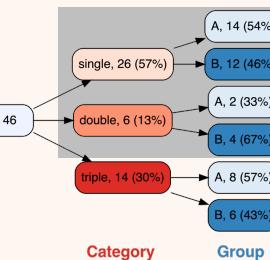
Modifier Effect

| Modifier | Effect |
|-----------------------------|----------------------------------|
| prefix <code>is.na</code> : | <code>is.na(variable)</code> |
| prefix <code>stem</code> : | all REDCap variables with stem |
| prefix <code>tri</code> : | trichotomize in each node |
| <code>variable=x</code> | dichotomize at <code>x</code> |
| <code>variable<x</code> | dichotomize below <code>x</code> |
| <code>variable>x</code> | dichotomize above <code>x</code> |

Pruning

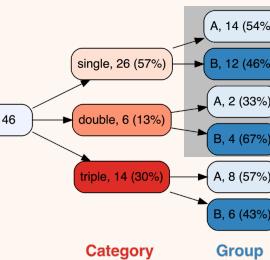
Prune single and double and their descendants

```
vtree(FakeData, "Category Group", sameline=TRUE, prune=list(Category=c("single", "double")))
```



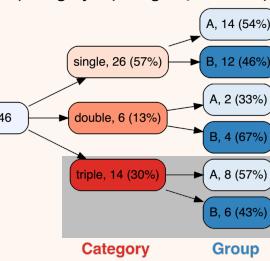
Prune nodes below single and double

```
vtree(FakeData, "Category Group", sameline=TRUE, prunebelow=list(Category=c("single", "double")))
```



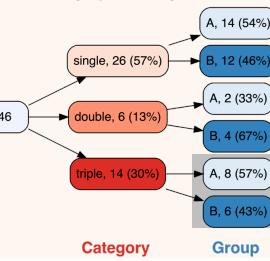
Only keep single and double and their descendants

```
vtree(FakeData, "Category Group", sameline=TRUE, keep=list(Category=c("single", "double")))
```



Only include descendants of single and double

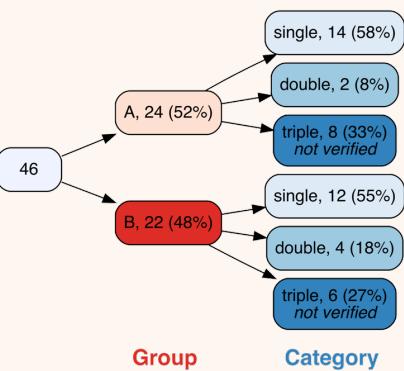
```
vtree(FakeData, "Category Group", sameline=TRUE, follow=list(Category=c("single", "double")))
```



Text

Add text to nodes

```
vtree(FakeData, "Group Category", sameline=TRUE, text=list(Category=c(triple="\n*not verified*")))
```



Group

Category

Labeling

Parameter

Change variable labels

```
labelvar=c(  
  Severity="New label for Severity")
```

Change node labels for a variable

```
labelnode=list(MyVar=  
  c(New="Old", New2="Old2"))
```

Change a specific node label

```
tlabelnode=list(  
  c(Group="A", Sex="F", label="girl"))
```

Font size (points) for variable names

```
varnamepointsize=15
```

Specify an optional label for the root node

```
title="All patients"
```

Show node labels?

```
shownodelabels=TRUE
```

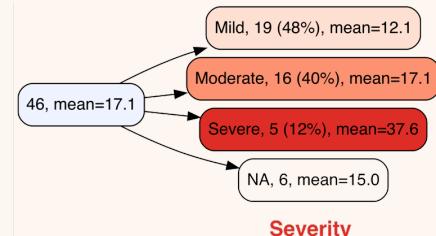
Show variable names?

```
showvarnames=TRUE
```

Summary information

Add summary statistics to nodes

```
vtree(FakeData, "Severity", sameline=TRUE, summary="Score , mean=%mean%")
```



Severity

Code Result

| | |
|-------------|---|
| %mean% | mean |
| %SD% | standard deviation |
| %min% | minimum |
| %max% | maximum |
| %px% | Xth percentile (e.g. p50 means the 50th percentile) |
| %median% | median, i.e. p50 |
| %IQR% | IQR, i.e. p25, p75 |
| %npct% | frequency and percentage |
| %pct% | just percentage |
| %list% | list of individual values, separated by commas |
| %listlines% | list of individual values, each on a separate line |
| %mv% | the number of missing values |

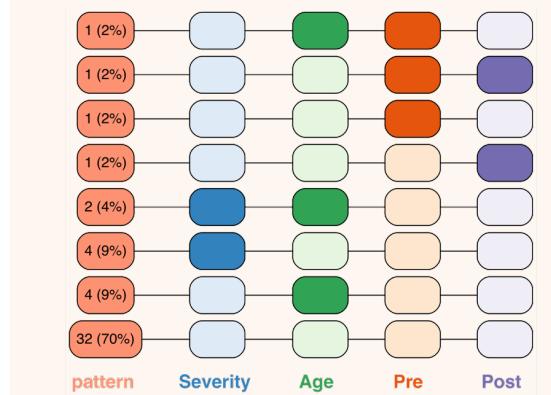
Code Summary information restricted to:

| | |
|------------|---------------------------|
| %noroot% | all nodes except the root |
| %leafonly% | leaf nodes |
| %var=v% | nodes of variable v |
| %node=n% | nodes named n |

Checking for missing values

Show missing value patterns; dark color = missing, light = not missing

```
vtree(FakeData, "Severity Age Pre Post", check.is.na=TRUE)
```



xplain Cheat Sheet

Important Links

- xplain package on CRAN <https://cran.r-project.org/web/packages/xplain/index.html>
- xplain web tutorial <http://www.zuckarelli.de/xplain/index.html>
- xplain cheat sheet http://www.zuckarelli.de/xplain/xplain_cheatsheet.pdf
- xplain on GitHub <https://github.com/jsugarelli/xplain>

Purpose & Application

- xplain allows to **write interpretation/explanation texts** for statistical functions in the form of XML files.
- The user of the functions can read these explanations **while working on his/her specific problems**.
- xplain explanations **can react to the user's results** and provide meaningful insights related to the user's problem.
- For this, the xplain **XML files can contain R code** and can **work with the return object** of the user's function call.

```
> xplain("lm(education ~ young + income + urban)")
> Your R^2 is 0.11 which is quite low. There is a serious
risk your model is misspecified. You should reconsider the
selection of variables included in your model.
```

xplain XML files

1 Any valid xplain XML must be enclosed in an `<xplain>` block. Multiple `<xplain>` blocks per XML file are possible.

2 A `<package>` block combines all functions from the same package.

3 Within a `<function>` block, explanations/interpretations for the function as such or for specific elements of the return object can be provided.

4 Packages explanations/ interpretations related to one element of the function's return object.

Main attributes: Overview

| | |
|--------------|---|
| name | Name of the element (package, function, result). |
| lang | Language (ISO code) of the explanation (e.g. "EN"). |
| level | Complexity level; integer number; cumulative, i.e. <code>level=1</code> explanations will also be presented when <code>level=2</code> or <code>level=3</code> are called. |

```
<xplain>
  1 <xplain>
    2 <package name = "stats">
      3 <function name = "lm">
        4 <title>This is about lm</title>
        5 <text>...</text>
        6 <result name = "coefficients">
          4 <title>...<title>
          5 <text>...</text>
        </result>
      </function>
    </package>
  </xplain>
</xplain>
```

Not case-sensitive

5 Structures explanations with headers.

6 The actual explanations/interpretations. Can include R code with references to the function's return object.

Attributes: Inheritance and necessity

- Elements **inherit attributes from higher-level** elements; e.g., if only one language, definition on `<xplain>` level suffices. Lower-level attributes overrule higher-level.
- name** attribute required for `<package>`, `<function>` and `<result>` elements.
- All levels shown, if no **level** is given to `xplain()`.

1
Direct call of
`xplain()`

call Call of the explained function as string

xml Path of the XML file providing the explanations

lang Language of the explanations to be shown (default means English)

level Complexity level of the explanations (cumulative! Default means "all")

Calling xplain()

2

Wrapper function with
`xplain.getcall()`

Example: lm

```
lm.xplain <- function(formula, data, subset, weights, na.action,
method = "qr", model = TRUE, x = FALSE, y = FALSE, qr = TRUE,
singular.ok = TRUE, contrasts = NULL, offset, ...) {
  call<-xplain.getcall("lm")
  xplain(call, xml="http://www.zuckarelli.de/example_lm.xml")}
```

Including R code

R code can be easily integrated into `<text></text>` elements:

```
<text> !%<% R code %!> </text>
      ↑           ↑
      R code delimiter tags
```

Access the explained function's (`<function name="...">`) return object:

- Access the full return object with `@`. Example: `summary(@)`.
- Access the current `<result name="...">` item of the return object with `##`. Example: `mean(##)`.

Using placeholders

```
<define name= "placeholder" > !%<% R code %!> </define>
      ↓           ↓
      Placeholder name delimiter tags
```

Example: `<define name="s">!%<% summary(@) %!> </define>`
`<text>And here is the summary !**s**! for your model</text>`

Iterating through (items of) the return object

- To apply a `<text>` element to a whole matrix, data frame, vector or list, use the **foreach** attribute.
- Value of foreach defines what is iterated over and (for 2D structures) in which sequence; `items` is for lists.
- \$ is a placeholder for the index of the current element. `"items"`
- Example** (shows all 1st column elements of the coefficient matrix):
`<text foreach="rows">!%<% @coefficients[$,1] %!> </text>`

foreach =
"rows"
"columns"
"rows, columns"
"columns, rows"
"items"