

Base R Cheat Sheet

Getting Help

Accessing the help files

?mean

Get help of a particular function.

help.search('weighted mean')

Search the help files for a word or phrase.

help(package = 'dplyr')

Find help for a package.

More about an object

str(iris)

Get a summary of an object's structure.

class(iris)

Find the class an object belongs to.

Using Packages

install.packages('dplyr')

Download and install a package from CRAN.

library(dplyr)

Load the package into the session, making all its functions available to use.

dplyr::select

Use a particular function from a package.

data(iris)

Load a built-in dataset into the environment.

Working Directory

getwd()

Find the current working directory (where inputs are found and outputs are sent).

setwd('C://file/path')

Change the current working directory.

Use projects in RStudio to set the working directory to the folder you are working in.

Vectors

Creating Vectors

| | | |
|-------------------|-------------|-----------------------------|
| c(2, 4, 6) | 2 4 6 | Join elements into a vector |
| 2:6 | 2 3 4 5 6 | An integer sequence |
| seq(2, 3, by=0.5) | 2.0 2.5 3.0 | A complex sequence |
| rep(1:2, times=3) | 1 2 1 2 1 2 | Repeat a vector |
| rep(1:2, each=3) | 1 1 1 2 2 2 | Repeat elements of a vector |

Vector Functions

sort(x)

Return x sorted.

rev(x)

Return x reversed.

table(x)

See counts of values.

unique(x)

See unique values.

Selecting Vector Elements

By Position

x[4]

The fourth element.

x[-4]

All but the fourth.

x[2:4]

Elements two to four.

x[!(2:4)]

All elements except two to four.

x[c(1, 5)]

Elements one and five.

By Value

x[x == 10]

Elements which are equal to 10.

x[x < 0]

All elements less than zero.

x[x %in% c(1, 2, 5)]

Elements in the set 1, 2, 5.

Named Vectors

x['apple']

Element with name 'apple'.

Programming

For Loop

```
for (variable in sequence){  
  Do something  
}
```

Example

```
for (i in 1:4){  
  j <- i + 10  
  print(j)  
}
```

While Loop

```
while (condition){  
  Do something  
}
```

Example

```
while (i < 5){  
  print(i)  
  i <- i + 1  
}
```

Functions

```
function_name <- function(var){  
  Do something  
  return(new_variable)  
}
```

Example

```
square <- function(x){  
  squared <- x*x  
  return(squared)  
}
```

Reading and Writing Data

Also see the **readr** package.

| Input | Output | Description |
|------------------------------|-------------------------------|--|
| df <- read.table('file.txt') | write.table(df, 'file.txt') | Read and write a delimited text file. |
| df <- read.csv('file.csv') | write.csv(df, 'file.csv') | Read and write a comma separated value file. This is a special case of read.table/write.table. |
| load('file.RData') | save(df, file = 'file.Rdata') | Read and write an R data file, a file type special for R. |

| Conditions | a == b | Are equal | a > b | Greater than | a >= b | Greater than or equal to | is.na(a) | Is missing |
|------------|--------|-----------|-------|--------------|--------|--------------------------|------------|------------|
| | a != b | Not equal | a < b | Less than | a <= b | Less than or equal to | is.null(a) | Is null |

Types

Converting between common data types in R. Can always go from a higher value in the table to a lower value.

| | | |
|--------------|---------------------------------|---|
| as.logical | TRUE, FALSE, TRUE | Boolean values (TRUE or FALSE). |
| as.numeric | 1, 0, 1 | Integers or floating point numbers. |
| as.character | '1', '0', '1' | Character strings. Generally preferred to factors. |
| as.factor | '1', '0', '1', levels: '1', '0' | Character strings with preset levels. Needed for some statistical models. |

Maths Functions

| | | | |
|--------------|---------------------------------|-------------|-------------------------|
| log(x) | Natural log. | sum(x) | Sum. |
| exp(x) | Exponential. | mean(x) | Mean. |
| max(x) | Largest element. | median(x) | Median. |
| min(x) | Smallest element. | quantile(x) | Percentage quantiles. |
| round(x, n) | Round to n decimal places. | rank(x) | Rank of elements. |
| signif(x, n) | Round to n significant figures. | var(x) | The variance. |
| cor(x, y) | Correlation. | sd(x) | The standard deviation. |

Variable Assignment

```
> a <- 'apple'  
> a  
[1] 'apple'
```

The Environment

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

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

Matrices

`m <- matrix(x, nrow = 3, ncol = 3)`
Create a matrix from x.

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

Lists

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

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

Also see the `dplyr` package.

Data Frames

`df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))`
A special case of a list where all elements are the same length.

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

Matrix subsetting

| | | | |
|-----------------------|--|---|------------------------------------|
| <code>df[, 2]</code> | | <code>nrow(df)</code> Number of rows. | <code>cbind</code> - Bind columns. |
| <code>df[2,]</code> | | <code>ncol(df)</code> Number of columns. | <code>rbind</code> - Bind rows. |
| <code>df[2, 2]</code> | | <code>dim(df)</code> Number of columns and rows. | |

Strings

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

Factors

| | |
|---------------------------------|--|
| <code>factor(x)</code> | Turn a vector into a factor. Can set the levels of the factor and the order. |
| <code>cut(x, breaks = 4)</code> | Turn a numeric vector into a factor by 'cutting' into sections. |

Statistics

| | | | |
|----------------------------------|--|------------------------------|--|
| <code>lm(y ~ x, data=df)</code> | Linear model. | <code>t.test(x, y)</code> | Test for a difference between proportions. |
| <code>glm(y ~ x, data=df)</code> | Generalised linear model. | <code>pairwise.t.test</code> | Perform a t-test for paired data. |
| <code>summary</code> | Get more detailed information out a model. | <code>aov</code> | Analysis of variance. |
| | | | |

Distributions

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

Plotting

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

Dates

See the `lubridate` package.

Data Import :: CHEAT SHEET

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

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

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

OTHER TYPES OF DATA

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

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

Save Data

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

Comma delimited file

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

File with arbitrary delimiter

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

CSV for excel

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

String to file

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

String vector to file, one element per line

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

Object to RDS file

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

Tab delimited files

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



Read Tabular Data

- These functions share the common arguments:

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

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

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

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

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

Comma Delimited Files

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

To make file.csv run:

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

Semi-colon Delimited Files

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

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

Files with Any Delimiter

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

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

Fixed Width Files

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

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

Tab Delimited Files

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

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

USEFUL ARGUMENTS

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

Example file

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

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

Skip lines

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

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

No header

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

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

Read in a subset

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

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

Provide header

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

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

Missing Values

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

Read Non-Tabular Data

Read a file into a single string

```
read_file(file, locale = default_locale())
```

Read each line into its own string

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

Read Apache style log files

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

Read a file into a raw vector

```
read_file_raw(file)
```

Read each line into a raw vector

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



Data types

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

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

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

age is an integer

sex is a character

1. Use **problems()** to diagnose problems.
`x <- read_csv("file.csv"); problems(x)`

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

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

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

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

Tibbles - an enhanced data frame

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

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

| # A tibble: 234 × 6 | manufacturer | model | displ | cyl | trans |
|---|--------------|------------|------------|-----|----------|
| 1 | audi | a4 | 1.80 | 4 | manual |
| 2 | audi | a4 | 1.80 | 4 | manual |
| 3 | audi | a4 | 2.00 | 4 | manual |
| 4 | audi | a4 | 2.00 | 4 | manual |
| 5 | audi | a4 | 2.00 | 4 | manual |
| 6 | audi | a4 | 2.00 | 4 | manual |
| 7 | audi | a4 | 2.00 | 4 | manual |
| 8 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 9 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 10 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 11 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 12 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 13 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 14 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 15 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 16 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 17 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 18 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 19 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| 20 | audi | a4 quattro | 1.80 | 4 | auto(l4) |
| ... | ... | ... | ... | ... | ... |
| 156 | 1999 | 6 | auto(l4) | | |
| 157 | 1999 | 6 | auto(l4) | | |
| 158 | 2008 | 6 | auto(l4) | | |
| 159 | 2008 | 8 | auto(s4) | | |
| 160 | 1999 | 4 | manual(m5) | | |
| 161 | 1999 | 4 | auto(m5) | | |
| 162 | 2008 | 4 | manual(m5) | | |
| 163 | 2008 | 4 | manual(m5) | | |
| 164 | 2008 | 4 | auto(l4) | | |
| 165 | 2008 | 4 | auto(l4) | | |
| 166 | 1999 | 4 | auto(l4) | | |
| [reached getOption("max.print") -- omitted 68 rows] | | | | | |

tibble display

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

data frame display

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

A large table to display

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

CONSTRUCT A TIBBLE IN TWO WAYS

tibble(...)
Construct by columns.
`tibble(x = 1:3, y = c("a", "b", "c"))`

tribble(...)
Construct by rows.
`tribble(~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.

Both make this tibble

Handle Missing Values

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

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

| x | x1 | x2 | x | 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.

| x | x1 | x2 | x | x1 | x2 |
|---|----|----|---|----|----|
| A | 1 | | A | 1 | |
| B | NA | | B | 2 | |
| C | NA | | C | 2 | |
| D | 3 | | D | 3 | |
| E | NA | | E | 2 | |

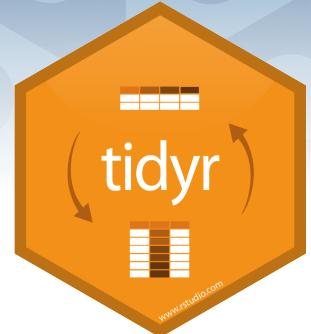
`replace_na(x, list(x2 = 2))`

Expand Tables

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)`

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.

| country | year | rate | country | year | cases | pop |
|---------|------|----------|---------|------|-------|-----|
| A | 1999 | 0.7K/19M | A | 1999 | 0.7K | 19M |
| A | 2000 | 2K/20M | A | 2000 | 2K | 20M |
| B | 1999 | 37K/172M | B | 1999 | 37K | 172 |
| B | 2000 | 80K/174M | B | 2000 | 80K | 174 |
| C | 1999 | 212K/1T | C | 1999 | 212K | 1T |
| C | 2000 | 213K/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_()**.

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

`separate_rows(table3, rate)`

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

Collapse cells across several columns to make a single column.

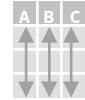
| country | century | year | country | year |
|---------|---------|------|---------|------|
| Afghan | 19 | 99 | Afghan | 1999 |
| Afghan | 20 | 0 | Afghan | 2000 |
| Brazil | 19 | 99 | Brazil | 1999 |
| Brazil | 20 | 0 | Brazil | 2000 |
| China | 19 | 99 | China | 1999 |
| China | 20 | 0 | China | 2000 |

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

Data Transformation with dplyr :: CHEAT SHEET



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



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



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



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

Summarise Cases

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



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

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

VARIATIONS

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

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

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

Group Cases

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



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

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

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

Manipulate Cases

EXTRACT CASES

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



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



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



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



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



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

Manipulate Variables

EXTRACT VARIABLES

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



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



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

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

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

MAKE NEW VARIABLES

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



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

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

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

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

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

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



Vector Functions

TO USE WITH MUTATE ()

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

vectorized function →

OFFSETS

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

CUMULATIVE AGGREGATES

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

RANKINGS

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

MATH

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

MISC

dplyr::case_when() - multi-case if_else()
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
 pmax() - element-wise max()
 pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

Summary Functions

TO USE WITH SUMMARISE ()

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

summary function →

COUNTS

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

LOCATION

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

LOGICALS

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

POSITION/ORDER

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

RANK

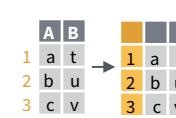
quantile() - nth quantile
min() - minimum value
max() - maximum value

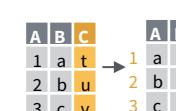
SPREAD

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

Row Names

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

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

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

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

Combine Tables

COMBINE VARIABLES

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

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

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

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

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

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

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

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

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

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

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

COMBINE CASES

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

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

| df A B C | bind_rows(..., .id = NULL) |
|---|--|
| x a t 1 x b u 2 x c v 3 z c v 3 z d w 4 | Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured) |
| | |

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

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

| A B C | union(x, y, ...) |
|----------------------------------|---|
| a t 1 b u 2 c v 3 d w 4 | Rows that appear in x or y. (Duplicates removed). union_all() retains duplicates. |
| | |

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

EXTRACT ROWS

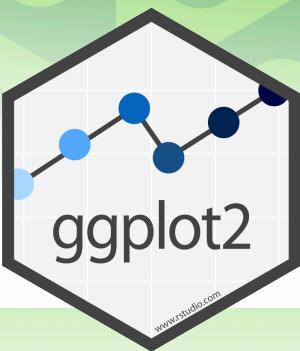
| X | y | = |
|----------------------------------|----------------------------------|---|
| A B C a t 1 b u 2 c v 3 | A B D a t 3 b u 2 d w 1 | |
| | | |

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

| A B C | semi_join(x, y, by = NULL, ...) |
|----------------|--|
| a t 1 b u 2 | Return rows of x that have a match in y. |
| | USEFUL TO SEE WHAT WILL BE JOINED. |

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

Data Visualization with ggplot2 :: CHEAT SHEET



Basics

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



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



Complete the template below to build a graph.

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

↑ required
Not required, sensible defaults supplied

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

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

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

last_plot() Returns the last plot

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

Geoms

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

GRAPHICAL PRIMITIVES

- a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
- a + geom_blank()**
(Useful for expanding limits)
- b + geom_curve(aes(yend = lat + 1, xend=long+1, curvature=z))** - x, yend, alpha, angle, color, curvature, linetype, size
- a + geom_path(lineend="butt", linejoin="round", linemitre=1)** - x, y, alpha, color, group, linetype, size
- a + geom_polygon(aes(group = group))** - x, y, alpha, color, fill, group, linetype, size
- b + geom_rect(aes(xmin = long, ymin=lat, xmax=long + 1, ymax = lat + 1))** - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size
- a + geom_ribbon(aes(ymin=unemploy - 900, ymax=unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

LINE SEGMENTS

- common aesthetics: x, y, alpha, color, linetype, size
- b + geom_abline(aes(intercept=0, slope=1))**
 - b + geom_hline(aes(yintercept = lat))**
 - b + geom_vline(aes(xintercept = long))**

- b + geom_segment(aes(yend=lat+1, xend=long+1))**
- b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

- c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
- c + geom_area(stat = "bin")** - x, y, alpha, color, fill, linetype, size
- c + geom_density(kernel = "gaussian")** - x, y, alpha, color, fill, group, linetype, size, weight
- c + geom_dotplot()** - x, y, alpha, color, fill
- c + geom_freqpoly()** - x, y, alpha, color, group, linetype, size
- c + geom_histogram(binwidth = 5)** - x, y, alpha, color, fill, linetype, size, weight
- c2 + geom_qq(aes(sample = hwy))** - x, y, alpha, color, fill, linetype, size, weight

discrete

- d <- ggplot(mpg, aes(f1))
- d + geom_bar()** - x, alpha, color, fill, linetype, size, weight

TWO VARIABLES

continuous x , continuous y

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

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

- e + geom_point()** - x, y, alpha, color, fill, shape, size, stroke

- e + geom_quantile()** - x, y, alpha, color, group, linetype, size, weight

- e + geom_rug(sides = "bl")** - x, y, alpha, color, linetype, size

- e + geom_smooth(method = lm)** - x, y, alpha, color, fill, group, linetype, size, weight

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

discrete x , continuous y

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

- f + geom_col()** - x, y, alpha, color, fill, group, linetype, size

- f + geom_boxplot()** - x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

- f + geom_dotplot(binaxis = "y", stackdir = "center")** - x, y, alpha, color, fill, group

- f + geom_violin(scale = "area")** - x, y, alpha, color, fill, group, linetype, size, weight

discrete x , discrete y

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

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

THREE VARIABLES

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

- l + geom_contour(aes(z = z))** - x, y, z, alpha, colour, group, linetype, size, weight

continuous bivariate distribution

- h <- ggplot(diamonds, aes(carat, price))
- h + geom_bin2d(binwidth = c(0.25, 500))** - x, y, alpha, color, fill, linetype, size, weight

- h + geom_density2d()** - x, y, alpha, colour, group, linetype, size

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

continuous function

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

- i + geom_area()** - x, y, alpha, color, fill, linetype, size

- i + geom_line()** - x, y, alpha, color, group, linetype, size

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

visualizing error

- df <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))

- j + geom_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size

- j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width (also **geom_errorbarh()**)

- j + geom_linerange()** - x, ymin, ymax, alpha, color, group, linetype, size

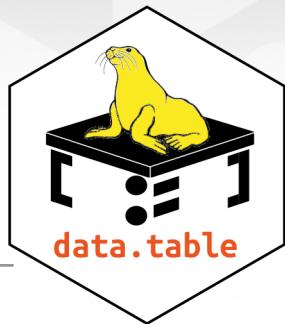
- j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

maps

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

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

Data Transformation with 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**.

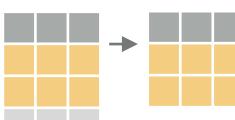
data.tables are also data frames – functions that work with data frames therefore also work with data.tables.

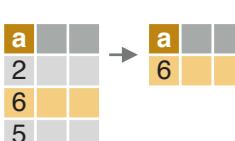
Create a data.table

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.

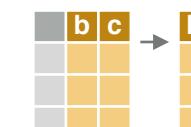
LOGICAL OPERATORS TO USE IN i

| | | | | | |
|---|----|----------|------|---|-----------|
| < | <= | 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.

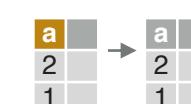
SUMMARIZE

 **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.

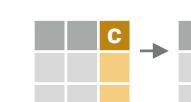
COMPUTE COLUMNS*

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

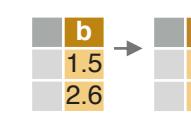
 **dt[a == 1, c := 1 + 2]** – compute a column based on an expression but only for a subset of rows.

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

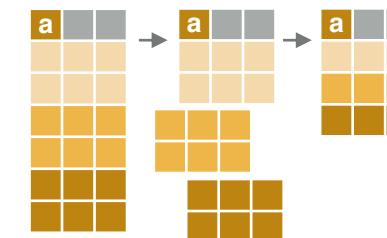
DELETE COLUMN

 **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.

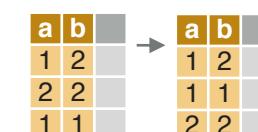
dt[, .SD[N], by = a] – extract last row of groups.

Chaining

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

Functions for data.tables

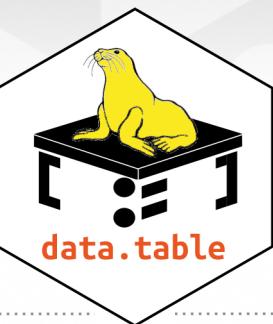
REORDER



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)`”.



UNIQUE ROWS

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

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

`uniqueN(dt, by = c("a", "b"))` – count the number of unique rows based on columns specified in “by”.

RENAME COLUMNS

| a | b |
|---|---|
| x | y |

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

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

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

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

| a | b | c |
|---|---|---|
| 1 | c | 7 |
| 2 | a | 5 |
| 3 | b | 6 |

| x | y | z |
|---|---|---|
| 3 | b | 4 |
| 2 | c | 5 |
| 1 | a | 8 |

| a | b | c | x |
|---|---|---|----|
| 3 | b | 4 | 3 |
| 1 | c | 5 | 2 |
| 2 | a | 8 | NA |

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

ROLLING JOIN

| a | id | date |
|---|----|------------|
| 1 | A | 01-01-2010 |
| 2 | A | 01-01-2012 |
| 3 | A | 01-01-2014 |
| 1 | B | 01-01-2010 |
| 2 | B | 01-01-2012 |

| b | id | date |
|---|----|------------|
| 1 | A | 01-01-2013 |
| 1 | B | 01-01-2013 |

| a | id | date | b |
|---|----|------------|---|
| 2 | A | 01-01-2013 | 1 |
| 1 | B | 01-01-2013 | 1 |

`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 |
|---|---|
| | |

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

| a | b |
|---|---|
| | |

| x | y |
|---|---|
| | |

| a | b | x | y |
|---|---|---|---|
| | | | |
| | | | |

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

Apply function to cols.

APPLY A FUNCTION TO MULTIPLE COLUMNS

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

`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 |
|---|---|
| 1 | 1 |
| 2 | 2 |
| 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.

RESHAPE TO WIDE FORMAT

| id | y | a | b |
|----|---|---|---|
| A | x | 1 | 3 |
| A | z | 2 | 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`
`id ~ y`
`value.var`
A data.table.
Formula with a LHS: ID columns containing IDs for multiple entries. And a RHS: columns with values to spread in column headers.
Columns containing values to fill into cells.

| id | a | x | a | z | b | x | b | z |
|----|---|---|---|---|---|---|---|---|
| A | 1 | 2 | 3 | 4 | | | | |
| B | 1 | 2 | 3 | 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`
`id.vars`
`measure.vars`
`variable.name,`
`value.name`
A data.table.
ID columns with IDs for multiple entries.
Columns containing values to fill into cells (often in pattern form).
Names of new columns for variables and values derived from old headers.

Sequential rows

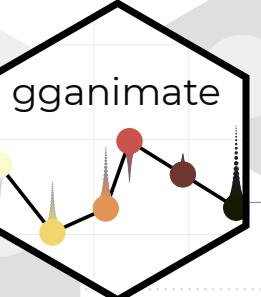
ROW IDS

| a | b |
|---|---|
| 1 | a |
| 2 | a |
| 3 | b |

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

| a | b |
|---|---|
|---|---|

Animate ggplots with gganimate :: CHEAT SHEET



Core Concepts

gganimate 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 gganimate's default settings .

Starting Plots

```
library(tidyverse)
library(gganimate)

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_filters()

```
a + transition_filter(transition_length = 3,
                      filter_length = 1,
                      cut == "Ideal",
                      Deep = depth >= 60)
```

`transition_length` and `filter_length` work the same as `transition/state_length()` in `transition_states()`...

... 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 = 3, state_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...

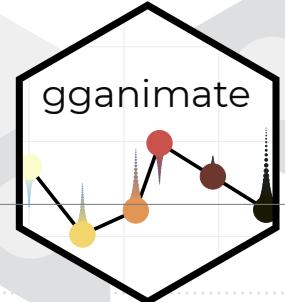
... and we're cycling between **seven** views. Seven is the number of steps in the transition, so the view is changing when the points are static, and visa versa. Views are determined by what data is in the current frame.

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 gganimate :: CHEAT SHEET



enter/exit_*

Every enter_*() function has a corresponding exit_*() function, and visa versa.

enter/exit_fade()

```
anim_a + enter_fade()
```

When new points need to be added, they will start transparent and become opaque.

enter_grow()/exit_shrink()

```
anim_a + exit_shrink()
```

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 gifski 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

gganimate'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 |
|--------------------------------------|---|--------------------------|
| transitioning | TRUE if the current frame is an transition frame, FALSE otherwise | states, layers, filter |
| previous_state/layer | Last shown state/layer | states, layers |
| next_state/layer | State/layer that will been shown next | states, layers |
| closest_state/layer | State/layer that current frame is closest to (if between states/layers, either next or closest). | states, layers |
| previous/closest/_filter/_expression | Similar to their state/layer analogs. *_filter variables return the name of the filter, *_expression variables return the condition. | filter |
| frame_time | Time of current frame | time, components, events |
| frame_along | Current frame's value for the dimension we're transitioning over | reveal |
| nlayers | Number of layers (total, not just currently shown) | layer |

h2o:: CHEAT SHEET

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.

| | | | |
|------------------|--------------|----------|------------------|
| x[j] ## column J | x[i] | <- value | Value Assignment |
| x[i, j] | x[i, j, ...] | <- value | |
| x[[i]] | x[[i]] | <- value | |
| x\$name | x\$i | <- value | |

Selection

Value Assignment

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:

h2o.asfactor, as.factor: Factor.

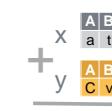
h2o.as_date, as.Date: Date.

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.

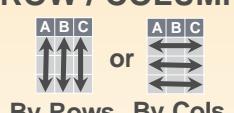
h2o.all: Given a set of logical vectors, determine if all of the values are true.

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.

GRID SEARCH

h2o.grid: Efficient method to build multiple models with different hyperparameters.

h2o.getGrid: Get a grid object from H2O distributed K/V store.

MODEL SCORING

h2o.predict: Obtain predictions from various fitted H2O model objects.

h2o.scoreHistory: Get Model Score History.

MODEL METRICS

h2o.make_metrics: Given predicted values (target for regression, class-1 probabilities, or binomial or per-class probabilities for multinomial), compute a model metrics object.

GENERAL MODEL HELPER

h2o.performance: Evaluate the predictive performance of a Supervised Learning Regression or Classification Model via various metrics. Set **xval = TRUE** for retrieving the training cross-validation metrics.

REGRESSION MODEL HELPER

h2o.mse: Display the mean squared error calculated from "Predicted Responses" and "Actual (Reference) Responses". Set **xval = TRUE** for retrieving the cross-validation MSE.

CLASSIFICATION MODEL HELPERS

h2o.accuracy: Get Model Accuracy metric.

h2o.auc: Retrieve the AUC (area under ROC curve). Set **xval = TRUE** for retrieving the cross-validation AUC.

h2o.confusionMatrix: Display prediction errors for classification data ("Predicted" vs "Reference : Real Values").

h2o.hit_ratio_table: Retrieve the Hit Ratios. Set **xval = TRUE** for retrieving the cross-validation Hit Ratio.

CLUSTERING MODEL HELPER

h2o.betweenss: Get the between cluster Sum of Squares.

h2o.centers: Retrieve the Model Centers.

PREDICTOR VARIABLE IMPORTANCE

h2o.varimp: Retrieve the variable importance

h2o.varimp_plot: Plot Variable Importances.

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.

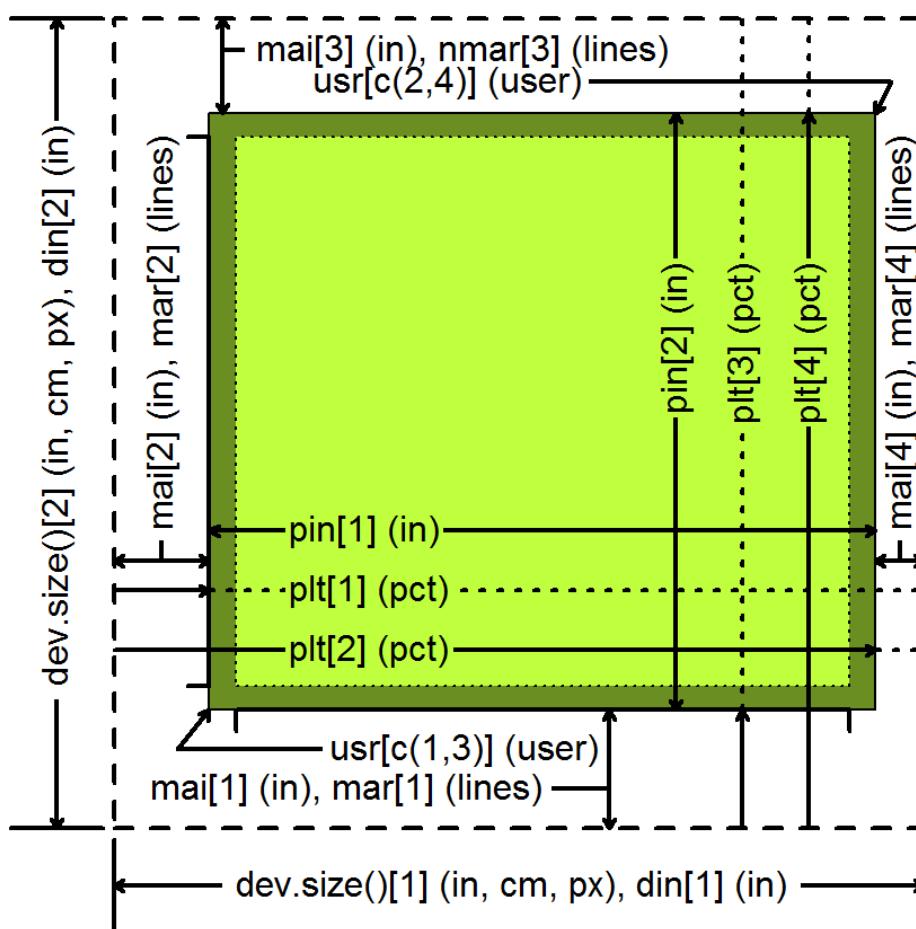


How Big is Your Graph?

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 **par** 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

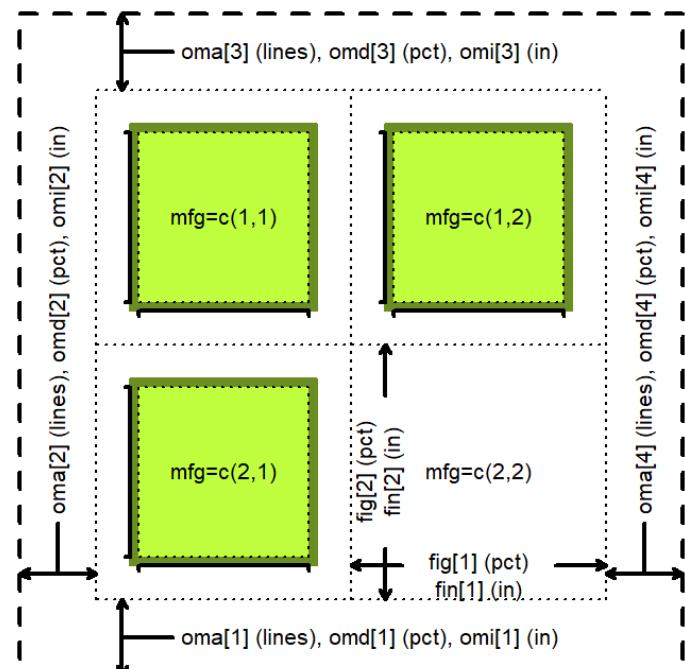
If you display multiple plots within a single graphics window (e.g., with the `mfrow` or `mfcol` arguments of `par` or with the `layout` function), then the `fig` and `fin` arguments will tell you the size of the current subplot window in percent or inches, respectively.

`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 ("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.

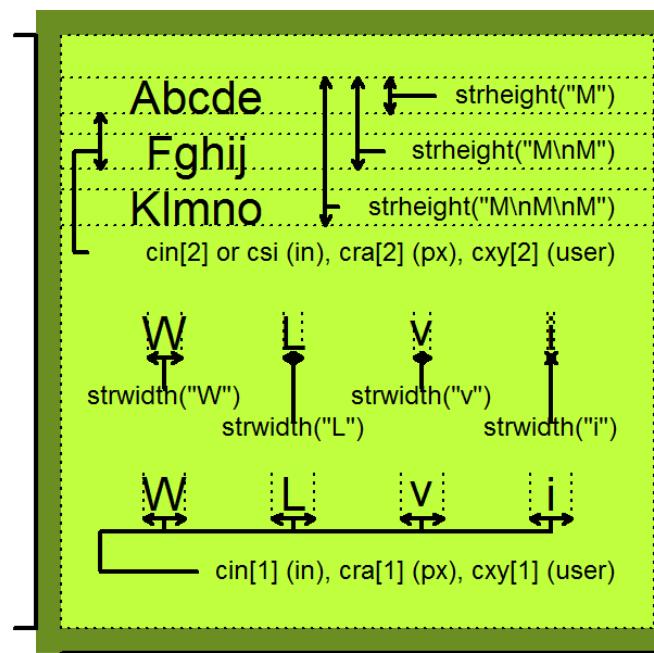
`par("cin")` (r.o.) (width, height) in inches
`par("csi")` (r.o.) height in inches
`par("cra")` (r.o.) (width, height) in pixels
`par("cxy")` (r.o.) (width, height) in xy coordinates

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

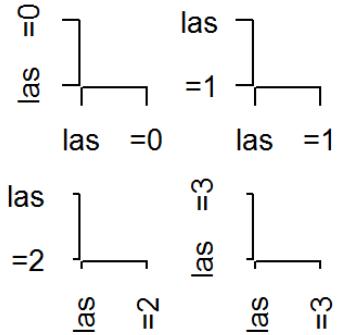
Fixing this takes a bit of trial and error.

1. Specify a larger/smaller value for the `pointsize` argument when you open your graphics device.
2. 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.
3. Use the `cex` argument to increase or decrease the relative size of your fonts.

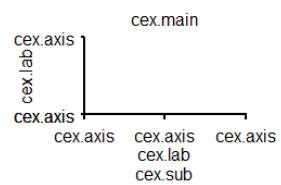
If your axes don't fit

There are several possible solutions.

1. You can assign wider margins using the `mar` or `mai` argument in `par`.
2. You can change the orientation of the axis labels with `las`. Choose among
 - a. `las=0` both axis labels parallel
 - b. `las=1` both axis labels horizontal
 - c. `las=2` both axis labels perpendicular
 - d. `las=3` both axis labels vertical.



3. change the relative size of the font
 - a. `cex.axis` for the tick mark labels.
 - b. `cex.lab` for `xlab` and `ylab`.
 - c. `cex.main` for the main title
 - d. `cex.sub` for the subtitle.



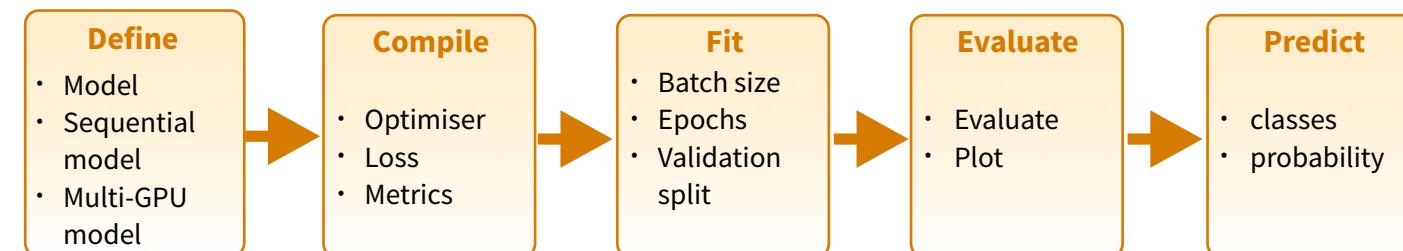
Deep Learning with Keras :: CHEAT SHEET



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>

The “Hello, World!”
of deep learning

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 on input activity



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



`layer_flatten()` Flattens an input

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`'.

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`

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

| | |
|---|--|
|  | <code>layer_conv_1d()</code> 1D, e.g. temporal convolution |
|  | <code>layer_conv_2d_transpose()</code> Transposed 2D (deconvolution) |
|  | <code>layer_conv_2d()</code> 2D, e.g. spatial convolution over images |
|  | <code>layer_conv_3d_transpose()</code> Transposed 3D (deconvolution) <code>layer_conv_3d()</code> 3D, e.g. spatial convolution over volumes |
|  | <code>layer_conv_lstm_2d()</code> Convolutional LSTM |
|  | <code>layer_separable_conv_2d()</code> Depthwise separable 2D |
|  | <code>layer_upsampling_1d()</code> <code>layer_upsampling_2d()</code> <code>layer_upsampling_3d()</code> Upsampling layer |
|  | <code>layer_zero_padding_1d()</code> <code>layer_zero_padding_2d()</code> <code>layer_zero_padding_3d()</code> Zero-padding layer |
|  | <code>layer_cropping_1d()</code> <code>layer_cropping_2d()</code> <code>layer_cropping_3d()</code> Cropping layer |

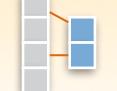
POOLING LAYERS

| | |
|---|---|
|  | <code>layer_max_pooling_1d()</code> <code>layer_max_pooling_2d()</code> <code>layer_max_pooling_3d()</code> Maximum pooling for 1D to 3D |
|  | <code>layer_average_pooling_1d()</code> <code>layer_average_pooling_2d()</code> <code>layer_average_pooling_3d()</code> Average pooling for 1D to 3D |
|  | <code>layer_global_max_pooling_1d()</code> <code>layer_global_max_pooling_2d()</code> <code>layer_global_max_pooling_3d()</code> Global maximum pooling |
|  | <code>layer_global_average_pooling_1d()</code> <code>layer_global_average_pooling_2d()</code> <code>layer_global_average_pooling_3d()</code> Global average pooling |

ACTIVATION LAYERS

| | |
|---|---|
|  | <code>layer_activation()</code> object, activation Apply an activation function to an output |
|  | <code>layer_activation_leaky_relu()</code> Leaky version of a rectified linear unit |
|  | <code>layer_activation_parametric_relu()</code> Parametric rectified linear unit |
|  | <code>layer_activation_thresholded_relu()</code> Thresholded rectified linear unit |
|  | <code>layer_activation_elu()</code> Exponential linear unit |

DROPOUT LAYERS

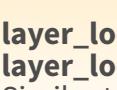
| | |
|---|---|
|  | <code>layer_dropout()</code> Applies dropout to the input |
|  | <code>layer_spatial_dropout_1d()</code> <code>layer_spatial_dropout_2d()</code> <code>layer_spatial_dropout_3d()</code> Spatial 1D to 3D version of dropout |

RECURRENT LAYERS

| | |
|---|---|
|  | <code>layer_simple_rnn()</code> Fully-connected RNN where the output is to be fed back to input |
|  | <code>layer_gru()</code> Gated recurrent unit - Cho et al |
|  | <code>layer_cudnn_gru()</code> Fast GRU implementation backed by CuDNN |

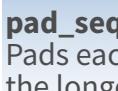
| | |
|---|--|
|  | <code>layer_lstm()</code> Long-Short Term Memory unit - Hochreiter 1997 |
|  | <code>layer_cudnn_lstm()</code> Fast LSTM implementation backed by CuDNN |

LOCALLY CONNECTED LAYERS

| | |
|---|--|
|  | <code>layer_locally_connected_1d()</code> <code>layer_locally_connected_2d()</code> Similar to convolution, but weights are not shared, i.e. different filters for each patch |
|---|--|

Preprocessing

SEQUENCE PREPROCESSING

| | |
|---|---|
|  | <code>pad_sequences()</code> Pads each sequence to the same length (length of the longest sequence) |
|  | <code>skipgrams()</code> Generates skipgram word pairs |
|  | <code>make_sampling_table()</code> Generates word rank-based probabilistic sampling table |

TEXT PREPROCESSING

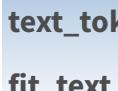
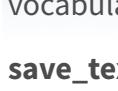
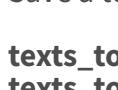
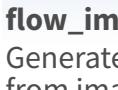
| | |
|---|---|
|  | <code>text_tokenizer()</code> Text tokenization utility |
|  | <code>fit_text_tokenizer()</code> Update tokenizer internal vocabulary |
|  | <code>save_text_tokenizer(); load_text_tokenizer()</code> Save a text tokenizer to an external file |
|  | <code>texts_to_sequences(); texts_to_sequences_generator()</code> Transforms each text in texts to sequence of integers |
|  | <code>texts_to_matrix(); sequences_to_matrix()</code> Convert a list of sequences into a matrix |
|  | <code>text_one_hot()</code> One-hot encode text to word indices |
|  | <code>text_hashing_trick()</code> Converts a text to a sequence of indexes in a fixed-size hashing space |
|  | <code>text_to_word_sequence()</code> Convert text to a sequence of words (or tokens) |

IMAGE PREPROCESSING

| | |
|---|---|
|  | <code>image_load()</code> Loads an image into PIL format. |
|  | <code>flow_images_from_data()</code> <code>flow_images_from_directory()</code> Generates batches of augmented/normalized data from images and labels, or a directory |
|  | <code>image_data_generator()</code> Generate minibatches of image data with real-time data augmentation. |
|  | <code>fit_image_data_generator()</code> Fit image data generator internal statistics to some sample data |
|  | <code>generator_next()</code> Retrieve the next item |
|  | <code>image_to_array(); image_array_resize(); image_array_save()</code> 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

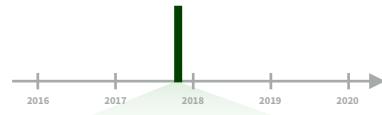
A callback is a set of functions to be applied at given stages of the training procedure. You can use callbacks to get a view on internal states and statistics of the model during training.

`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(20170131)

July 4th, 2000

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

4th of July '99

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

2001: Q3

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

2:01

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

2017.5

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



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

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

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

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

2017-11-28

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

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

12:00:00

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

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

GET AND SET COMPONENTS

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

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

2018-01-31 11:59:59

date(x) Date component. date(dt)

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

hour(x) Hour. hour(dt)

2018-01-31 11:59:59

minute(x) Minutes. minute(dt)

2018-01-31 11:59:59

second(x) Seconds. second(dt)

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

Round Date-times



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

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

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

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

Tip: use a date with day > 12

Stamp Date-times

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

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

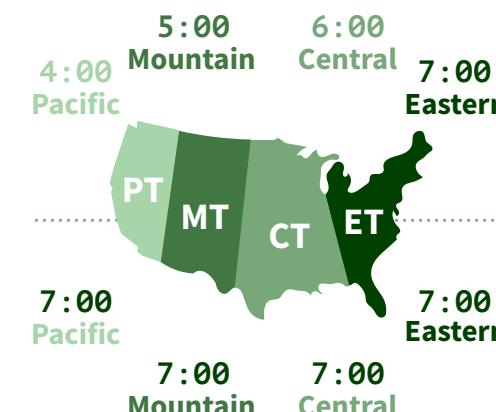
2. Apply the template to dates
`sf(ymd("2010-04-05"))`
`## [1] "Created Monday, Apr 05, 2010 00:00"`

Time Zones

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

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

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



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

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



Math with Date-times

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

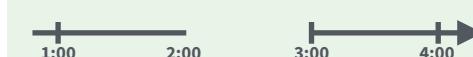
A normal day

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



The start of daylight savings (spring forward)

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



The end of daylight savings (fall back)

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



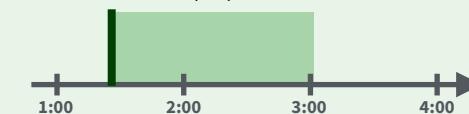
Leap years and leap seconds

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

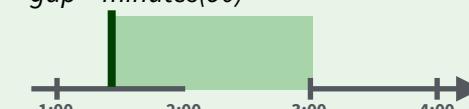


Periods track changes in clock times, which ignore time line irregularities.

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



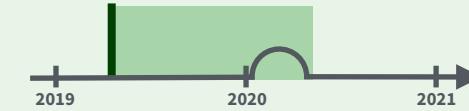
```
gap + minutes(90)
```



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

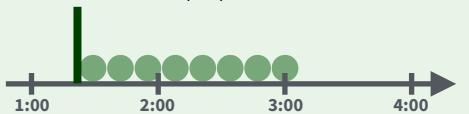


```
leap + years(1)
```



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

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



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



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



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



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

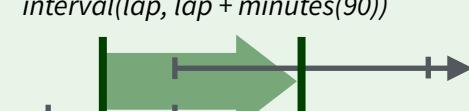
```
interval(nor, nor + minutes(90))
```



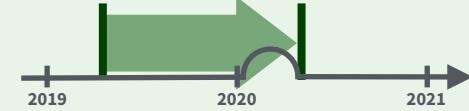
```
interval(gap, gap + minutes(90))
```



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



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



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

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

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

```
jan31 <- ymd(20180131)
```

```
jan31 + months(1)
```

```
## NA
```

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

```
jan31 %m+% months(1)
```

```
## "2018-02-28"
```

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

```
add_with_rollback(jan31, months(1), roll_to_first = TRUE)
```

```
## "2018-03-01"
```

PERIODS

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

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

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

```
"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
```



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())

Machine Learning with R



Introduction

mlr offers a unified interface for the basic building blocks of machine learning: tasks, learners, hyperparameters, etc.

Tasks contain a description of a task (classification, regression, clustering, etc.) and a data set.

Learners specify a machine learning algorithm (GLM, SVM, xgboost, etc.) and its parameters.

Hyperparameters are learner settings that can be specified directly or tuned. A **parameter set** lists the possible hyperparameters for a given learner.

Wrapped Models are learners that have been trained on a task and can be used to make predictions.

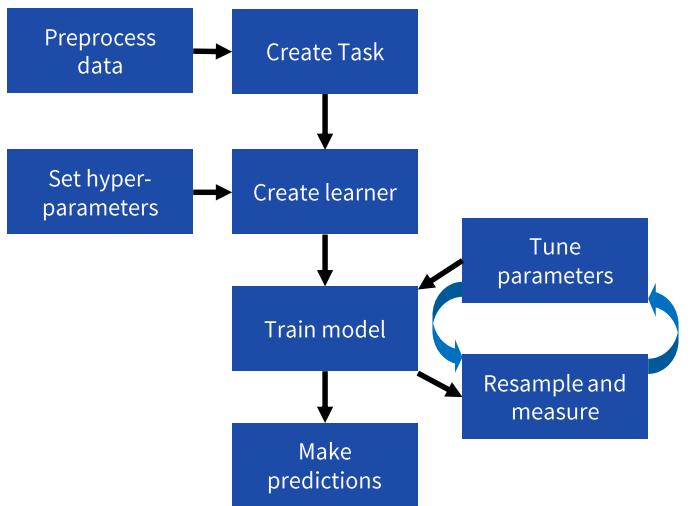
Predictions are the results of applying a model to either new data or the original training data.

Measures control how learner performance is evaluated, e.g. RMSE, LogLoss, AUC, etc.

Resampling estimates generalization performance by separating training data from test data. Common strategies include holdout and cross-validation.

Links: [Tutorial](#) | [CRAN](#) | [Github](#)

mlr workflow



Setup

Preprocessing data

`createDummyFeatures(obj=, target=, method=, cols=)`
Creates (0,1) flags for each non-numeric variable excluding `target`. Can be applied to entire dataset or only specific `cols`

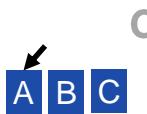
`normalizeFeatures(obj=, target=, method=, cols=, range=, on.constant=)`
Normalizes numerical features according to specified `method`:

- "center" (subtract mean)
- "scale" (divide by std. deviation)
- "standardize" (center and scale)
- "range" (linear scale to given range, default `range=c(0,1)`)

`mergeSmallFactorLevels(task=, cols=, min.perc=)`
Combine infrequent factor levels into a single merged level

`summarizeColumns(obj=)` where `obj` is a data.frame or task.
Provides type, NA, and distributional data about each column

See also `capLargeValues` `dropFeatures` `removeConstantFeatures` `summarizeLevels`

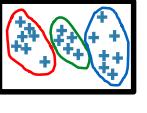


Creating a task

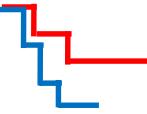
`makeClassifTask(data=, target=)`
Classification of a target variable, with optional positive class `positive`

`makeRegrTask(data=, target=)`
Regression on a target variable

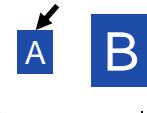
`makeMultilabelTask(data=, target=)`
Classification where the target can belong to more than one class per observation



`makeClusterTask(data=)`
Unsupervised clustering on a data set



`makeSurvTask(data=, target= c("time", "event"))`
Survival analysis with a survival time column and an event column



`makeCostSensTask(data=, costs=)`
Cost-sensitive classification where each observation-cost pair has a specified cost

Other arguments that can be passed to a `task`:

- `weights`= Weighting vector to apply to observations
- `blocking`= Factor vector where each level indicates a block of observations that will not be split up in resampling

Making a learner

`makeLearner(cl=, predict.type=, ..., par.vals=)`
Choose an algorithm class to perform the task and determine what that algorithm will predict

- `cl`=name of algorithm, e.g. `"classif.xgboost"` `"regr.randomForest"` `"cluster.kmeans"`
- `predict.type="response"` returns a prediction type that matches the source data; `"prob"` returns a predicted probability for classification problems only; `"se"` returns the standard error of the prediction for regression problems only. Only certain learners can return `"prob"` and `"se"`
- `par.vals`= takes a list of hyperparameters and passes them to the learner; parameters can also be passed directly (...)

You can make multiple learners at once with `makeLearners()`

mlr has integrated over 170 different learning algorithms

- Full list: `View(listLearners())` shows all learners
- Available learners for a task: `View(listLearners(task))`
- Filtered list: `View(listLearners("classif", properties=c("prob", "factors")))` shows all classification learners `"classif"` which can predict probabilities `"prob"` and handle factor inputs `"factors"`
- See also `getLearnerProperties()`

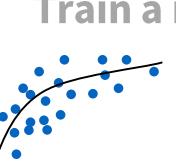
Training & Testing

Setting hyperparameters

`setHyperPars(learner=, ...)`
Set the hyperparameters (settings) for each learner, if you don't want to use the defaults. You can also specify hyperparameters in the `makeLearner()` call

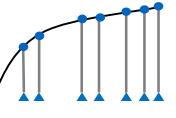


`getParamSet(learner=)`
Show the possible universe of parameters for your learner; can take a learner directly, or a text string such as `"classif.qda"`



Train a model and predict

`train(learner=, task=)`
Train a model (`WrappedModel`) by applying a learner to a task. By default, the model will train on all observations. The underlying model can be extracted with `getLearnerModel()`



`predict(object=, task=, newdata=)`
Use a trained model to make predictions on a task or dataset. The resulting `pred` object can be viewed with `View(pred)` or accessed by `as.data.frame(pred)`



Measuring performance

`performance(pred=, measures=)`
Calculate performance of predictions according to one or more of several measures (use `listMeasures()` for full list):

- `classif` `acc` `auc` `bac` `ber` `brier[,scaled]` `f1` `fdr` `fnr` `fpr` `gmean` `multiclass[,au1]` `aunp` `aunu` `brier` `npv` `ppv` `qsr` `ssr` `tn` `tnr` `tp` `tpr` `wkappa`
- `regr` `rsq` `expvar` `kendalltau` `mae` `mape` `medae` `medse` `mse` `msle` `rae` `rmse` `rmsle` `rrse` `rsq` `sae` `spearmanrho` `sse`
- `cluster` `db` `dunn` `G1` `G2` `silhouette`
- `multilabel` `multilabel[,f1]` `subset01` `.tpr` `.ppv` `.acc` `.hamloss`
- `costsens` `mcp` `meancosts`
- `surv` `cindex`
- `other` `featperc` `timeboth` `timelpredict` `timetrain`

For detailed performance data on classification tasks, use:

- `calculateConfusionMatrix(pred=)`
- `calculateROCMeasures(pred=)`

Resampling a learner

`makeResampleDesc(method=, ..., stratify=)`

`method` must be one of the following:

- "CV" (cross-validation, for number of folds use `iters=`)
 - "LOO" (leave-one-out cross-validation, for folds use `iters=`)
 - "RepCV" (repeated cross-validation, for number of repetitions use `reps=`, for folds use `folds=`)
 - "Subsample" (aka Monte-Carlo cross-validation, for iterations use `iters=`, for train % use `split=`)
 - "Bootstrap" (out-of-bag bootstrap, uses `iters=`)
 - "Holdout" (for train % use `split=`)
- `stratify` keeps target proportions consistent across samples.

`makeResampleInstance(desc=, task=)` can reduce noise by ensuring the resampling is done identically every time.

`resample(learner=, task=, resampling=, measures=)`
Train and test model according to specified resampling strategy.

mlr includes several pre-specified resample descriptions: `cv2` (2-fold cross-validation), `cv3`, `cv5`, `cv10`, `hout` (holdout with split 2/3 for training, 1/3 for testing). Convenience functions also exist to `resample()` with a specific strategy: `crossval()`, `repCV()`, `holdout()`, `subsample()`, `bootstrap00B()`, `bootstrapB632()`, `bootstrapB632plus()`

Refining Performance

Tuning hyperparameters

Set search space using `makeParamSet(make<type>Param())`

- `makeNumericParam(id=, lower=, upper=, trafo=)`
 - `makeIntegerParam(id=, lower=, upper=, trafo=)`
 - `makeIntegerVectorParam(id=, len=, lower=, upper=, trafo=)`
 - `makeDiscreteParam(id=, values=c(...))` (can also be used to test discrete values of numeric or integer parameters)
- `trafo` transforms the parameter output using a specified function, e.g. `lower=-2, upper=2, trafo=function(x) 10^x` would test values between 0.01 and 100, scaled exponentially
- Other acceptable parameter types include `Logical` `LogicalVector` `CharacterVector` `DiscreteVector`

Set a search algorithm with `makeTuneControl<type>()`

- `Grid(resolution=10)` Grid of all possible points
- `Random(maxit=100)` Randomly sample search space
- `MBO(budget=)` Use Bayesian model-based optimization
- `Irace(n.instances=)` Iterated racing process
- Other types: `CMAES`, `Design`, `GenSA`

Tune using `tuneParams(learner=, task=, resampling=, measures=, par.set=, control=)`

Quickstart

Prepare data for training and testing

```

library(mlbench)
data(Soybean)
soy = createDummyFeatures(Soybean, target="Class")
tsk = makeClassifTask(data=soy, target="Class")
ho = makeResampleInstance("Holdout", tsk)
tsk.train = subsetTask(tsk, ho$train.ind[1])
tsk.test = subsetTask(tsk, ho$test.ind[1])

```

Convert the factor inputs in the Soybean dataset into (0,1) dummy features which can be used by the XGboost algorithm. Create a task to predict the "Class" column. Create a train set with 2/3 of data and a test set with the remaining 1/3 (default).

Create learner and evaluate performance

```

lrn = makeLearner("classif.xgboost", nrounds=10)
cv = makeResampleDesc("CV", iters=5)
res = resample(lrn, tsk.train, cv, acc)

```

Create an XGboost learner which will build 10 trees. Then test performance using 5-fold cross-validation. Accuracy should be between 0.90-0.92.

Tune hyperparameters and retrain model

```

ps = makeParamSet(makeNumericParam("eta", 0, 1),
  makeNumericParam("lambda", 0, 200),
  makeIntegerParam("max_depth", 1, 20))
tc = makeTuneControlMBO(budget=100)
tr = tuneParams(lrn, tsk.train, cv5, acc, ps, tc)
lrn = setHyperPars(lrn, par.vals=tr$x)

```

Tune hyperparameters `eta`, `lambda`, and `max_depth` by defining a search space and using Model Based Optimization (MBO) to control the search. Then perform 100 rounds of 5-fold cross-validation, improving accuracy to ~0.93. Update the XGboost learner with the tuned hyperparameters.

```

mdl = train(lrn, tsk.train)
prd = predict(mdl, tsk.test)
calculateConfusionMatrix(prd)
mdl = train(lrn, tsk)

```

Train the model on the train set and make predictions on the test set. Show performance as a confusion matrix. Finally, re-train model on the full set to use on new data. You are now ready to go out into the real world and make 93% accurate predictions!

Legend for functions (not all parameters shown):

`function(required_parameters, optional_parameters=)`

Configuration

mlr's default settings can be changed using `configureMlr()`:

- `show.info` Whether to show verbose output by default when training, tuning, resampling, etc. (`TRUE`)
- `on.learner.error` How to handle a learner error. `"stop"` halts execution, `"warn"` returns NAs and displays a warning, `"quiet"` returns NAs with no warning (`"stop"`)
- `on.learner.warning` How to handle a learner warning. `"warn"` displays a warning, `"quiet"` suppresses it (`"warn"`)
- `on.par.without.desc` How to handle a parameter with no description. `"stop"`, `"warn"`, `"quiet"` (`"stop"`)
- `on.par.out.of.bounds` How to handle a parameter with an out-of-bounds value. `"stop"`, `"warn"`, `"quiet"` (`"stop"`)
- `on.measure.not.applicable` How to handle a measure not applicable to a learner. `"stop"`, `"warn"`, `"quiet"` (`"stop"`)
- `show.learner.output` Whether to show learner output to the console during training (`TRUE`)
- `on.error.dump` Whether to create an error dump for crashed learners if `on.learner.error` is not set to `"stop"` (`TRUE`)

Use `getMlrOptions()` to see current settings

Parallelization

mlr works with the `parallelMap` package to take advantage of multicore and cluster computing for faster operations. mlr automatically detects which operations are able to run in parallel.

To begin parallel operation use:

- ```
parallelStart(mode=, cpus=, level=)
```
- `mode` determines how the parallelization is performed:
    - `"local"` no parallelization applied, simply uses `mapply`
    - `"multicore"` multicore execution on a single machine, uses `parallel::mclapply`. Not available in Windows.
    - `"socket"` multicore execution in socket mode
    - `"mpi"` Snow MPI cluster on one or multiple machines using `parallel::makeCluster` and `parallel::clusterMap`
    - `"BatchJobs"` Batch queuing HPC clusters using `BatchJobs::batchMap`
  - `cpus` determines how many logical cores will be used
  - `level` controls parallelization: `"mlr.benchmark"`, `"mlr.resample"`, `"mlr.selectFeatures"`, `"mlr.tuneParams"`, `"mlr.ensemble"`

To end parallelization, use `parallelStop()`

# Imputation

`impute(obj=, target=, cols=, dummy.cols=, dummy.type=)`  
Applies specified logic to data frame or task containing NAs and returns an imputation description which can be used on new data

- `obj`=data frame or task on which to perform imputation
- `target`=specify target variable which will not be imputed
- `cols`=column names and logic for imputation\*
- `dummy.cols`=column names to create a NA (T/F) column\*
- `dummy.type`=set to `"numeric"` to use (0,1) instead of (T/F)

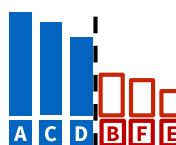
\*Can also use `classes` and `dummy.classes` in place of `cols`

Imputation logic is passed to `cols` or `classes` via a list, e.g.: `cols=list(V1=imputeMean())` where `V1` is the column to which to apply the imputation, and `imputeMean()` is the imputation method. Available imputation methods include:  
`imputeConst(const=)` `imputeMedian()` `imputeMode()` `imputeMin(multiplier=)` `imputeMax(multiplier=)` `imputeNormal(mean=, sd=)` `imputeHist(breaks=, use.mids=)` `imputeLearner(learner=, features=)` `impute` returns a list containing the imputed dataset or task as well as an imputation description that can be used to reapply the same imputation to new data using `reimpute`

`reimpute(obj=, desc=)` Imputes missing values on a task or dataset (`obj`) using a description (`desc`) created by `impute`

# Feature Extraction

## Feature filtering



`filterFeatures(task=, method=, perc=, abs=, threshold=)`  
Uses a learner-agnostic feature evaluation method to rank feature importance, then includes only features in the top n percent (`perc=`), top n (`abs=`), or which meet a set performance threshold (`threshold=`).

Outputs a task with features that failed the test omitted. `method` defaults to `"randomForestSRC.rfsrc"`, but can be set to:  
`"anova.test"` `"carscore"` `"cforest.importance"`  
`"chi.squared"` `"gain.ratio"` `"information.gain"`  
`"kruskal.test"` `"linear.correlation"` `"mrmr"` `"oneR"`  
`"permutation.importance"` `"randomForest.importance"`  
`"randomForestSRC.rfsrc"` `"randomForestSRC.var.select"`  
`"rank.correlation"` `"relief"`  
`"symmetrical.uncertainty"` `"univariate.model.score"`  
`"variance"`

## Feature selection



`selectFeatures(learner=, task=, resampling=, measures=, control=)`  
Uses a feature selection algorithm (`control`) to resample and build a model repeatedly using different feature sets each time in order to find the best set.

Available controls include:

- `makeFeatSelControlExhaustive(max.features=)` Try every combination of features up to optional `max.features`
- `makeFeatSelControlRandom(maxit=, prob=, max.features=)` Randomly sample features with probability `prob` (default 0.5) until `maxit` (default 100) iterations; return the best one found
- `makeFeatSelControlSequential(method=, maxit=, max.features=, alpha=, beta=)` Perform an iterative search using a `method` from the following: `"sfs"` forward search, `"sbs"` backward search, `"sfbs"` floating forward search, `"sfbs"` floating backward search. `alpha` indicates minimum improvement required to add a feature; `beta` indicates minimum required to remove a feature
- `makeFeatSelControlGA(maxit=, max.features=, mu=, lambda=, crossover.rate=, mutation.rate=)` Genetic algorithm trains on random feature vectors, then uses crossover on the best performers to produce 'offspring', repeated over generations. `mu` is size of parent population, `lambda` is size of children population, `crossover.rate` is probability of choosing a bit from first parent, `mutation.rate` is probability of flipping a bit (on or off)

`selectFeatures` returns a `FeatSelResult` object which contains optimal features and an optimization path. To apply feature selection result (`fsr`) to your task (`tsk`), use:  
`tsk = subsetTask(tsk, features=fsr$x)`

# Benchmarking

`benchmark(learners=, tasks=, resamplings=, measures=)`  
Allows easy comparison of multiple learners on a single task, a single learner on multiple tasks, or multiple learners on multiple tasks. Returns a benchmark result object.

Benchmark results can be accessed with a variety of functions beginning with `getBMR<object>.AggrPerformance`  
`FeatSelResults` `FilteredFeatures` `LearnerIds`  
`LeanerShortNames` `Learners` `MeasureIds` `Measures`  
`Models` `Performances` `Predictions` `TaskDescs` `TaskIds`  
`TuneResults`

mlr contains several toy tasks which are useful for benchmarking:  
`agri.task` `bc.task` `bh.task` `costiris.task` `iris.task`  
`lung.task` `mtcars.task` `pid.task` `sonar.task`  
`wpbc.task` `yeast.task`

# Visualization

## Performance

`generateThreshVsPerfData(obj=, measures=)` Measure performance at different probability cutoffs to determine optimal decision threshold for binary classification problems

- `plotThreshVsPerf(obj=)` Plot visual representation of threshold curve(s) from `ThreshVsPerfData`
- `plotROCCurves(obj=)` Plot receiver operating characteristic (ROC) curve from `ThreshVsPerfData`. Must set `measures=list(fpr, tpr)`

## Residuals

- `plotResiduals(obj=)` Plots residuals for `Prediction` or `BenchmarkResult`

## Learning curve

`generateLearningCurveData(learners=, task=, resampling=, percs=, measures=)` Measure performance of learner(s) trained on different percentages of task data

- `plotLearningCurve(obj=)` Plot curve showing learner performance vs. proportion of data used, uses `LearningCurveData`

## Feature importance

`generateFilterValuesData(task=, method=)` Get feature importance rankings using specified filter method

- `plotFilterValues(obj=)` Plot bar chart of feature importance based on filter method using `FilterValuesData`

## Hyperparameter tuning

`generateHyperParsEffectData(tune.result=)` Get the impact of different hyperparameter settings on model performance

- `plotHyperParsEffect(hyperpars.effec.t.data=, x=, y=, z=)` Create a plot showing hyperparameter impact on performance using `HyperParsEffectData`

See also:

- `plotOptPath(op=)` Display details of optimization process. Takes `<obj>$opt.path`, where `<obj>` is an object of class `tuneResult` or `featSelResult`
- `plotTuneMultiCritResult(res=)` Show pareto front for results of tuning to multiple performance measures

## Partial dependence

`generatePartialDependenceData(obj=, input=)` Get partial dependence of model (`obj`) prediction over each feature of data (`input`)

- `plotPartialDependence(obj=)` Plots partial dependence of model using `PartialDependenceData`

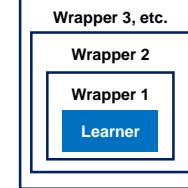
## Benchmarking

`plotBMRBoxplots(bmr=)` Distribution of performances  
`plotBMRSummary(bmr=)` Scatterplot of avg. performances  
`plotBMRanksAsBarChart(bmr=)` Rank learners in bar plot

## Other

- `generateCritDifferencesData(bmr=, measure=, p.value=, test=)` Perform critical-differences test using either the Bonferroni-Dunn ("bd") or "Nemenyi" test
- `plotCritDifferences(obj=)`
- `generateCalibrationData(obj=)` Evaluate calibration of probability predictions vs. true incidence
- `plotCalibration(obj=)`

# Wrappers



**Wrappers** fuse a learner with additional functionality. mlr treats a learner with wrappers as a single learner, and hyperparameters of wrappers can be tuned jointly with underlying model parameters. Models trained with wrappers will apply them to new data.

## Preprocessing and imputation

`makeDummyFeaturesWrapper(learner=)`  
`makeImputeWrapper(learner=, classes=, cols=)`  
`makePreprocWrapper(learner=, train=, predict=)`  
`makePreprocWrapperCaret(learner=, ...)`  
`makeRemoveConstantFeaturesWrapper(learner=)`

## Class imbalance

`makeOverBaggingWrapper(learner=)`  
`makeSMOTEWrapper(learner=)`  
`makeUndersampleWrapper(learner=)`  
`makeWeightedClassesWrapper(learner=)`

## Cost-sensitive learning

`makeCostSensClassifWrapper(learner=)`  
`makeCostSensRegrWrapper(learner=)`  
`makeCostSensWeightedPairsWrapper(learner=)`

## Multilabel classification

`makeMultilabelBinaryRelevanceWrapper(learner=)`  
`makeMultilabelClassifierChainsWrapper(learner=)`  
`makeMultilabelDBRWrapper(learner=)`  
`makeMultilabelNestedStackingWrapper(learner=)`  
`makeMultilabelStackingWrapper(learner=)`

## Other

`makeBaggingWrapper(learner=)`  
`makeConstantClassWrapper(learner=)`  
`makeDownsampleWrapper(learner=, dw.perc=)`  
`makeFeatSelWrapper(learner=, resampling=, control=)`  
`makeFilterWrapper(learner=, fw.perc=, fw.abs=, fw.threshold=)`  
`makeMultiClassWrapper(learner=)`  
`makeTuneWrapper(learner=, resampling=, par.set=, control=)`

## Nested Resampling

mlr supports **nested resampling** for complex operations such as tuning and feature selection through wrappers. In order to get a good estimate of generalization performance and avoid data leakage, both an outer (for tuning/feature selection) and an inner (for the base model) resampling process are advised.

- Outer resampling can be specified in `resample` or `benchmark`
- Inner resampling can be specified in `makeTuneWrapper`, `makeFeatSelWrapper`, etc.

## Ensembles

`makeStackedLearner(base.learners=, super.learner=, method=)` Combines multiple learners to create an ensemble

- `base.learners`=learners to use for initial predictions
- `super.learner`=learner to use for final prediction
- `method`=how to combine base learner predictions:
  - `"average"` simple average of all base learners
  - `"stack.nocv", "stack.cv"` train super learner on results of base learners, with or without cross-validation
  - `"hill.climb"` search for optimal weighted average
  - `"compress"` with a neural network for faster performance

# Intro stats with mosaic

(lattice version)

## 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)
```

## Loading packages

```
library(mosaic)
```

## Arithmetic operations

|                           |                             |
|---------------------------|-----------------------------|
| <code>+ - * /</code>      | basic operations            |
| <code>^</code>            | exponentiation              |
| <code>( )</code>          | grouping                    |
| <code>sqrt(x)</code>      | square root                 |
| <code>abs(x)</code>       | absolute value              |
| <code>log10(x)</code>     | logarithm, base 10          |
| <code>log(x)</code>       | natural logarithm, base $e$ |
| <code>exp(x)</code>       | exponential function $e^x$  |
| <code>factorial(k)</code> | $k! = k(k-1) \dots 1$       |

## Logical operators

|                                   |                                                               |
|-----------------------------------|---------------------------------------------------------------|
| <code>==</code>                   | is equal to (note double equal sign)                          |
| <code>!=</code>                   | is not equal to                                               |
| <code>&lt;</code>                 | is less than                                                  |
| <code>&lt;=</code>                | is less than or equal to                                      |
| <code>&gt;</code>                 | is greater than                                               |
| <code>&gt;=</code>                | is greater than or equal to                                   |
| <code>&amp;</code>                | <code>A &amp; B</code> is TRUE if both A and B are TRUE       |
| <code> </code>                    | <code>A   B</code> is TRUE if one or both of A and B are TRUE |
| <code>%in%</code>                 | includes; for example                                         |
| <code>"C" %in% c("A", "B")</code> | 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, groups = w,
 data = mydata)
```

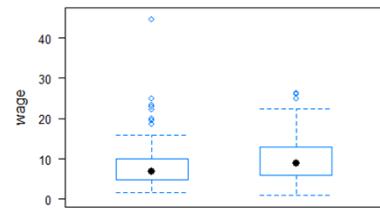
`y` : y-axis variable (*optional*)

`x` : x-axis variable (*required*)

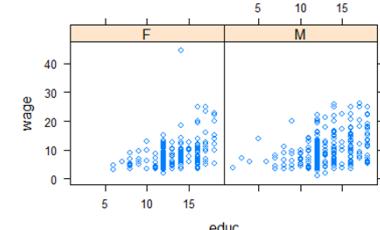
`z` : panel-by variable (*optional*)

`w` : color-by variable (*optional*)

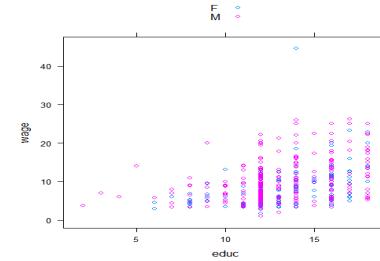
```
bwplot(wage ~ sex, data = CPS85)
```



```
xyplot(wage ~ educ | sex,
 data = CPS85)
```



```
xyplot(wage ~ educ,
 groups = sex, data = CPS85,
 auto.key = TRUE)
```



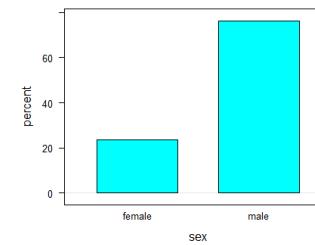
## One categorical variable

Counts by category

```
tally(~ sex, data = HELPrc)
```

Percentages by category

```
tally(~ sex, format =
 "percent", data = HELPrc)
bargraph(~ sex, type =
 "percent", data = HELPrc)
```



Tests and confidence intervals

Exact test

```
result1 <-
binom.test(~ (homeless ==
"homeless"), data = HELPrc)
```

Approximate test (large samples)

```
result2 <-
prop.test(~ (homeless ==
"homeless"), data = HELPrc)
```

Extract confidence intervals and p-values

```
confint(result1)
pval(result2)
```

## 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 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

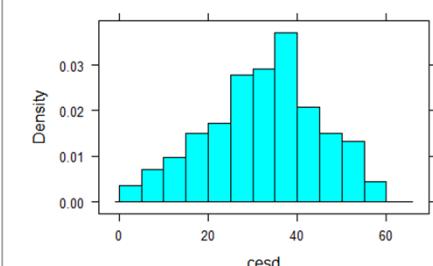
```
favstats(cesd ~ sex,
 data = HELPrc)
```

Quantiles

```
quantile(~ cesd, data = HELPrc,
 prob = c(0.25, 0.5, 0.8))
```

Histogram

```
histogram(~ cesd, width = 5,
 center = 2.5, data = HELPrc)
```



Normal probability plot

```
qqmath(~ cesd, dist = "qnorm",
 data = HELPrc)
```

Density plot

```
densityplot(~ cesd, data =
 HELPrc)
```

Dot plot

```
dotPlot(~ cesd, data = HELPrc)
```

One-sample t-test

```
result <- t.test(~ cesd, mu =
 34, data = HELPrc)
```

Extract confidence intervals and p-values

```
confint(result)
pval(result)
```

## Two categorical variables

Contingency table with margins

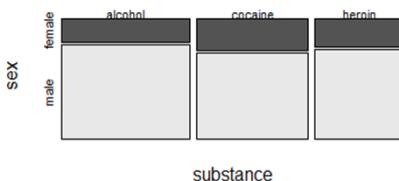
```
tally(~ substance + sex,
 margins = TRUE,
 data = HELPrc)
```

Percentages by column

```
tally(~ sex | substance,
 format = "percent",
 data = HELPrc)
```

Mosaic plot

```
mosaicplot(~ substance + sex,
 color = TRUE, data = HELPrc)
```



Chi-square test

```
xchisq.test(~ substance + sex,
 data = HELPrc,
 correct = FALSE)
```

## 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)
```

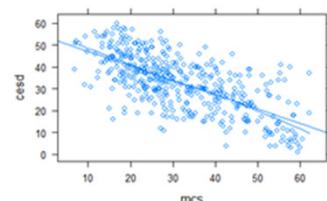
## Two quantitative variables

Correlation coefficient

```
cor(cesd ~ mcs, data = HELPrc)
```

Scatterplot with regression line and smooth

```
xyplot(cesd ~ mcs,
 type = c("p", "r", "smooth"),
 data = HELPrc)
```



Simple linear regression

```
cesdmodel <- lm(cesd ~ mcs,
 data = HELPrc)
```

```
msummary(cesdmodel)
```

Prediction

```
lmfunction <- makeFun(cesdmodel)
lmfunction(mcs = 35)
```

Extract useful quantities

```
anova(cesdmodel)
```

```
coef(cesdmodel)
```

```
confint(cesdmodel)
```

```
rsquared(cesdmodel)
```

Diagnostics; plot residuals

```
histogram(~resid(cesdmodel),
 density = TRUE)
```

```
qqmath(~resid(cesdmodel))
```

Diagnostics; plot residuals vs. fitted

```
xyplot(resid(cesdmodel) ~
 fitted(cesdmodel),
 type = c("p", "smooth", "r"))
```

## Categorical response, quantitative predictor

Logistic regression

```
logit_mod <-
 glm(homeless ~ age + female,
 family = binomial, data = HELPrc)
```

```
msummary(logitmod)
```

Odds ratios and confidence intervals

```
exp(coef(logit_mod))
```

```
exp(confint(logit_mod))
```

## Data management

From dplyr package

Drop or reorder variables

```
select()
```

Create new variables from existing ones

```
mutate()
```

Retain specific rows from data

```
filter()
```

Sort data rows

```
arrange()
```

Compute summary statistics by group

```
group_by()
```

```
summarize()
```

Merge data tables

```
left_join()
```

```
inner_join()
```

## Importing data

Import file from computer or URL

```
MustangPrice <-
 read.file("C:/MustangPrice.csv")
NOTE: R uses forward slashes!
Dome <-
 read.file("http://www.mosaic-
web.org/go/datasets/Dome.csv")
```

## Randomization and simulation

Fix random number sequence

```
set.seed(42)
```

Tossing 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)
```

## Quantitative response, categorical predictor

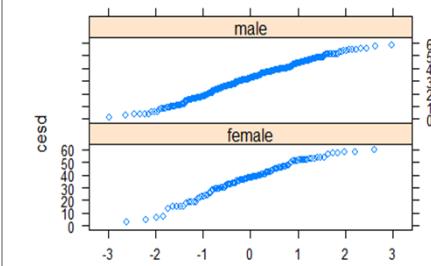
Two-level predictor: two-sample  $t$  test

Numeric summaries

```
favstats(~cesd | sex,
 data = HELPrc)
```

Comparative normal probability plot

```
qqmath(~cesd | sex, data = HELPrc,
 layout = c(1, 2)) # also bwplot
```



Dotplot for smaller samples

```
xyplot(sex ~ length, alpha = 0.6,
 cex = 1.4, data = KidsFeet)
```

Two-sample  $t$ -test and confidence interval

```
result <- t.test(cesd ~ sex,
 var.equal = FALSE, data = HELPrc)
```

```
confint(result)
```

More than two levels: Analysis of variance

Numeric summaries

```
favstats(cesd ~ substance,
 data = HELPrc)
```

Graphic summaries

```
bwplot(cesd ~ substance, pch = "|",
 data = HELPrc)
```

Fitt and summarize model

```
modsubstance <- lm(cesd ~ substance,
 data = HELPrc)
```

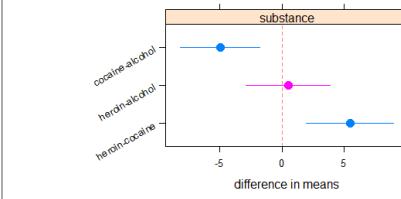
```
anova(modsubstance)
```

Which differences are significant?

```
pairwise <- TukeyHSD(modsubstance)
```

```
mplot(pairwise)
```

95% family-wise confidence level



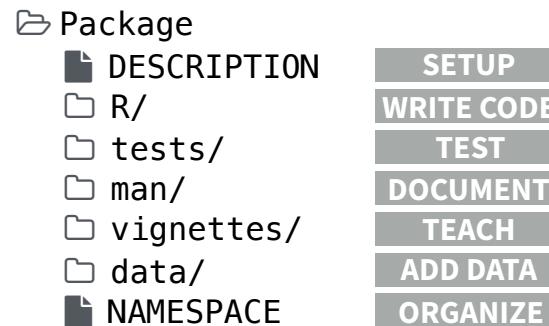
# Package Development: : CHEAT SHEET



## Package Structure

A package is a convention for organizing files into directories.

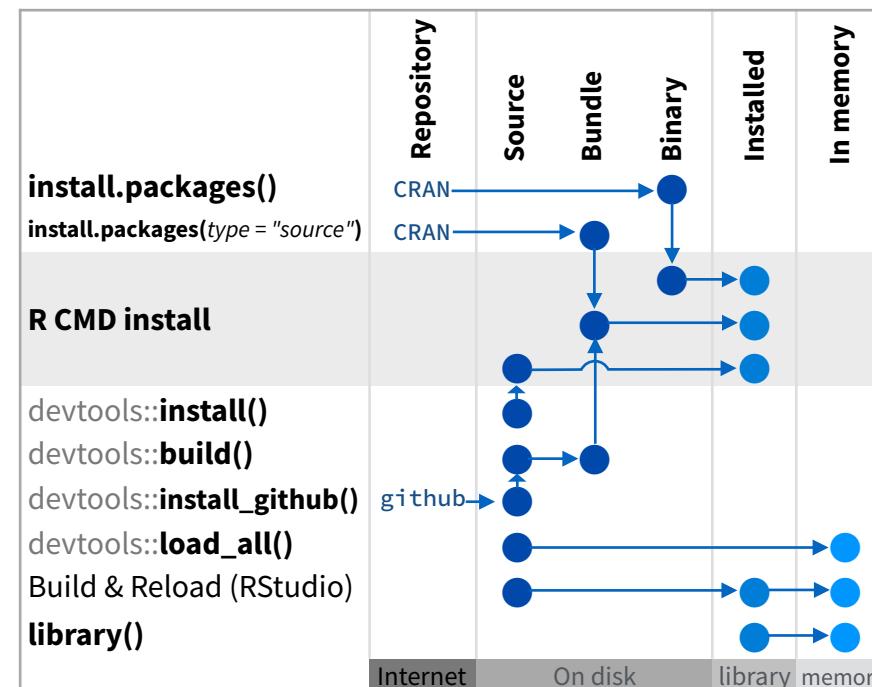
This sheet shows how to work with the 7 most common parts of an R package:



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                  | MIT                                            | GPL-2                                                                                          |
|----------------------|------------------------------------------------|------------------------------------------------------------------------------------------------|
| No strings attached. | MIT license applies to your code if re-shared. | GPL-2 license applies to your code, <i>and all code anyone bundles with it</i> , if re-shared. |

Package: mypackage  
Title: Title of Package  
Version: 0.1.0  
Authors@R: person("Hadley", "Wickham", email = "hadley@me.com", 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}{<br>left \tab centered \tab right \cr<br>cell \tab cell \tab cell \cr} |

### 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 | <b>@seealso</b> |      |
| @concepts        | @keywords      | @format         |      |
| @describeln      | <b>@param</b>  | @source         | data |
| <b>@examples</b> | @rdname        | @include        |      |
| <b>@export</b>   | <b>@return</b> | @slot           | S4   |
| @family          | @section       | @field          | RC   |

## 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](http://rmarkdown.rstudio.com))

```

```

```
title: "Vignette Title"
author: "Vignette Author"
date: "`r Sys.Date()`"
output: rmarkdown::html_vignette
vignette: >
 \%VignetteIndexEntry{Vignette Title}
 \%VignetteEngine{knitr::rmarkdown}
 \usepackage[utf8]{inputenc}
```

```

```

## 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](http://r-pkgs.had.co.nz/release.html)

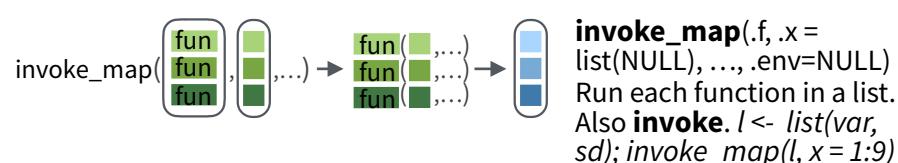
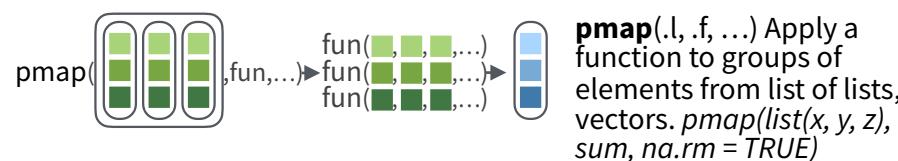
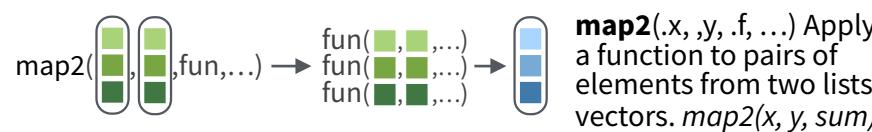
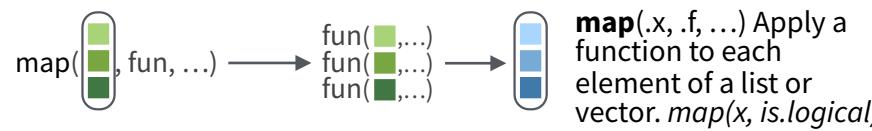


# Apply functions with purrr :: CHEAT SHEET



## Apply Functions

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



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

### OUTPUT

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

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

### SHORTCUTS - within a purrr function:

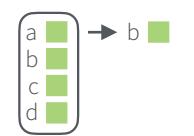
"**name**" becomes `function(x) x[["name"]]`, e.g. `map(l, "a")` extracts `a` from each element of `l`

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

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

## Work with Lists

### FILTER LISTS



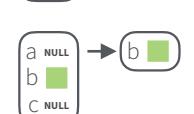
**pluck(.x, ..., .default=NULL)** Select an element by name or index, `pluck(x, "b")`, or its attribute with `attr_getter`. `pluck(x, "b", attr_getter("n"))`



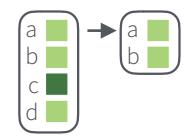
**keep(.x, .p, ...)** Select elements that pass a logical test. `keep(x, is.na)`



**discard(.x, .p, ...)** Select elements that do not pass a logical test. `discard(x, is.na)`



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



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

### SUMMARISE LISTS



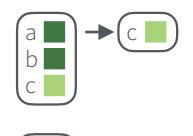
**every(.x, .p, ...)** Do all elements pass a test? `every(x, is.character)`



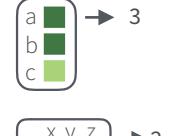
**some(.x, .p, ...)** Do some elements pass a test? `some(x, is.character)`



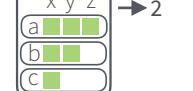
**has\_element(.x, .y)** Does a list contain an element? `has_element(x, "foo")`



**detect(.x, .f, ..., .right=FALSE, .p)** Find first element to pass. `detect(x, is.character)`

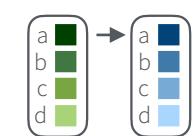


**detect\_index(.x, .f, ..., .right = FALSE, .p)** Find index of first element to pass. `detect_index(x, is.character)`

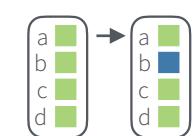


**vec\_depth(x)** Return depth (number of levels of indexes). `vec_depth(x)`

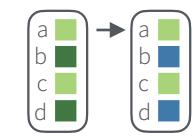
### TRANSFORM LISTS



**modify(.x, .f, ...)** Apply function to each element. Also **map**, **map\_chr**, **map\_dbl**, **map\_dfc**, **map\_dfr**, **map\_int**, **map\_lgl**. `modify(x, ~.+2)`



**modify\_at(.x, .at, .f, ...)** Apply function to elements by name or index. Also **map\_at**. `modify_at(x, "b", ~.+2)`



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

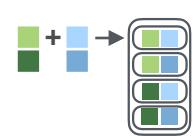


**modify\_depth(.x, .depth, .f, ...)** Apply function to each element at a given level of a list. `modify_depth(x, 1, ~.+2)`

### WORK WITH LISTS



**array\_tree(array, margin = NULL)** Turn array into list. Also **array\_branch**. `array_tree(x, margin = 3)`



**cross2(.x, .y, .filter = NULL)** All combinations of .x and .y. Also **cross**, **cross3**, **cross\_df**. `cross2(1:3, 4:6)`

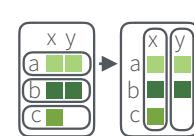


**set\_names(x, nm = x)** Set the names of a vector/list directly or with a function. `set_names(x, c("p", "q", "r"))`  
`set_names(x, tolower)`

### RESHAPE LISTS

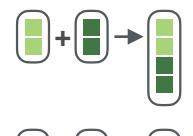


**flatten(.x)** Remove a level of indexes from a list. Also **flatten\_chr**, **flatten\_dbl**, **flatten\_dfc**, **flatten\_dfr**, **flatten\_int**, **flatten\_lgl**. `flatten(x)`

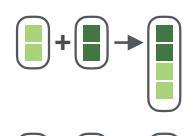


**transpose(.l, .names = NULL)** Transposes the index order in a multi-level list. `transpose(x)`

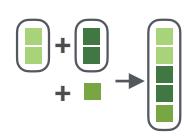
### JOIN (TO) LISTS



**append(x, values, after = length(x))** Add to end of list. `append(x, list(d = 1))`

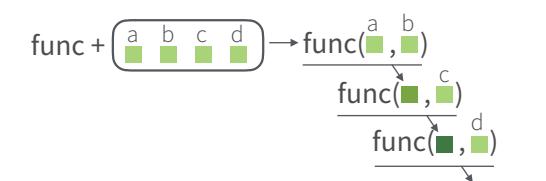


**prepend(x, values, before = 1)** Add to start of list. `prepend(x, list(d = 1))`

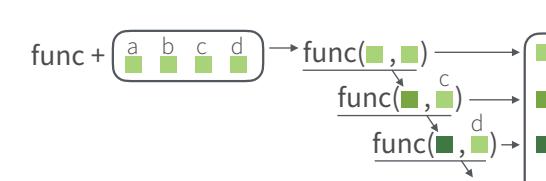


**splice(...)** Combine objects into a list, storing S3 objects as sub-lists. `splice(x, y, "foo")`

## Reduce Lists



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



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

## Modify function behavior

**compose()** Compose multiple functions.

**lift()** Change the type of input a function takes. Also **lift\_dl**, **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 whenever an error occurs (instead of error).

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

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



# Nested Data

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

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

n\_iris\$data[[1]]

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

n\_iris

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

n\_iris\$data[[2]]

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

n\_iris\$data[[3]]

Use a nested data frame to:

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

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

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

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

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

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

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

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

n\_iris %>% unnest()

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

Unnests a nested data frame.

# List Column Workflow

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

## 1 Make a list column

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

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

## 2 Work with list columns

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

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

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

## 3 Simplify the list column

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

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

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

## 1. MAKE A LIST COLUMN

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

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

Makes list column when needed

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

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

Saves list input as list columns

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

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

Converts multi-level list to tibble with list cols

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

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

Returns list col when result returns list.

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

## 2. WORK WITH LIST COLUMNS

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

`map2(<tibble [50x4]>, <S3: lm>, fun, ...)`

`fun(<tibble [50x`

# quanteda Cheat Sheet

Quantitative Analysis of Textual Data

## 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
1789-Washington 625 1538 23 1789 Washington George
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
1789-Washington 1 1 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         |



# 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.

```
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 ```
7
8 ```{python, echo = FALSE}
9 import seaborn as sns
10 fmri = sns.load_dataset("fmri")
11 ```
12
13 ```{r}
14 f1 <- subset(py$fmri, region == "parietal")
15
16
17 ```{python}
18 import matplotlib as mpl
19 sns.lmplot("timepoint","signal", data=r.f1)
20 mpl.pyplot.show()
21 ```
```

R Console: R Markdown:

## 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))`

```
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
```

## 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$`



## Python in the IDE

Requires reticulate plus RStudio v1.2 or higher.

Syntax highlighting for Python scripts and chunks

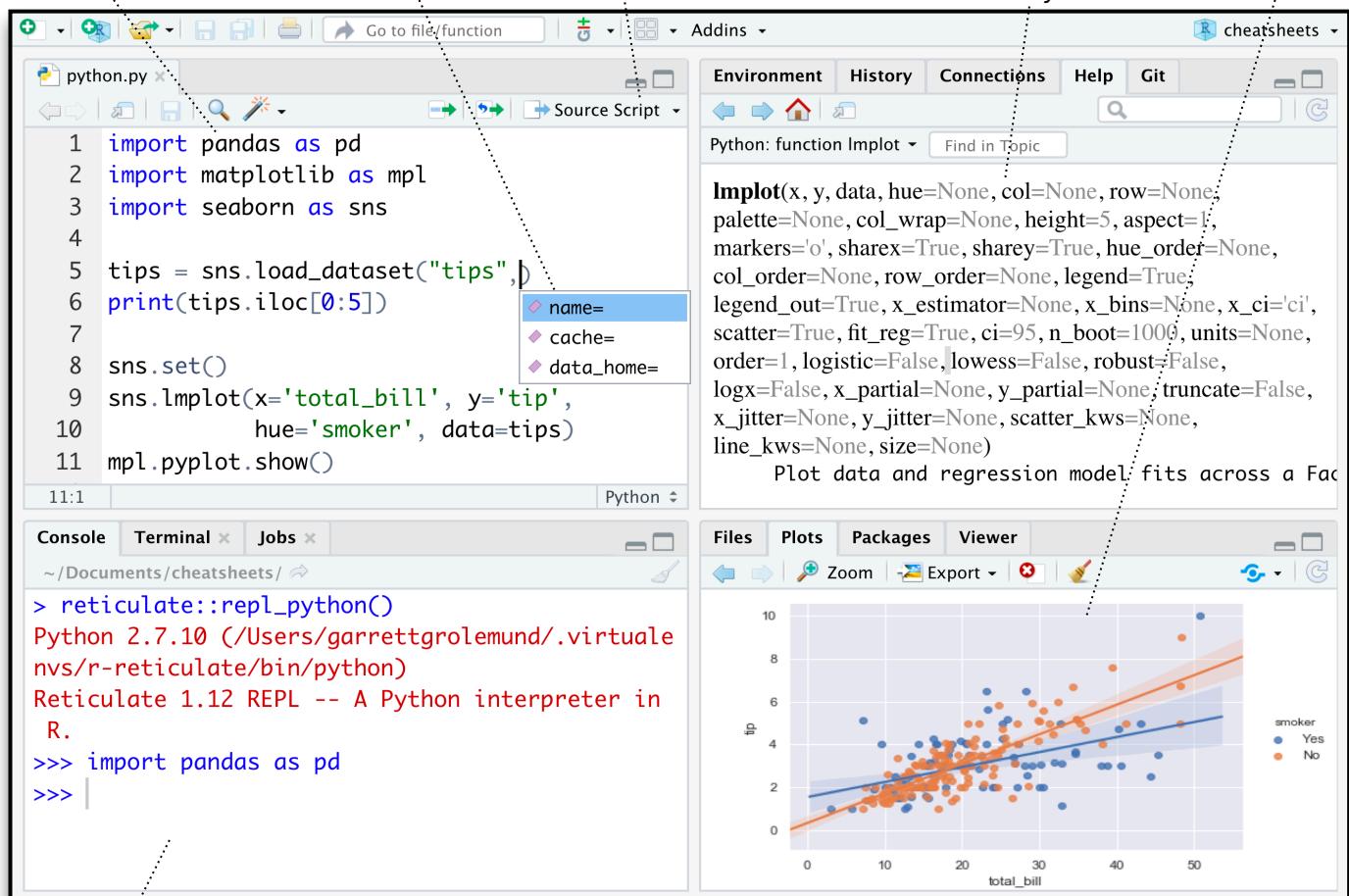
Tab completion for Python functions and objects (and Python modules imported in R scripts)

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.

matplotlib plots display in plots pane.

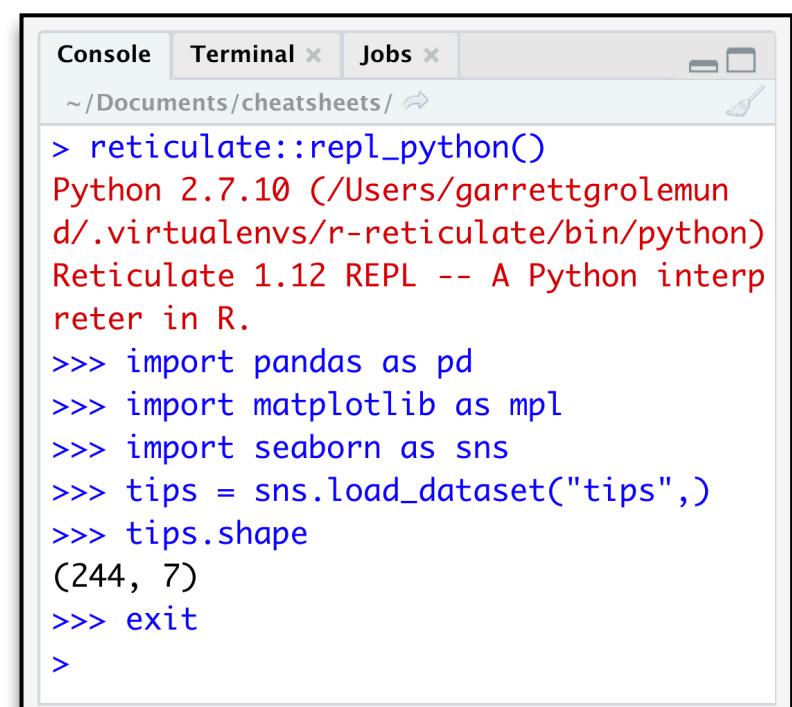


A Python REPL opens in the console when you run Python code with a keyboard shortcut. Type **exit** to close.

## 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**).
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")**
- **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")**

- **virtualenv\_list()** List all available virtualenvs. Also **virtualenv\_root().virtualenv\_list()**
- **conda\_list(conda = "auto")** List all available conda envs. Also **conda\_binary()** and **conda\_version().conda\_list()**

### 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?

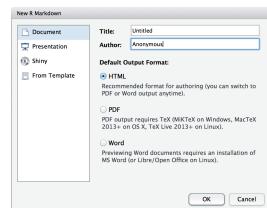


**.Rmd files** • An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.

**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

## render

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

**input** - file to render  
**output\_format**

**output\_options** - List of render options (as in YAML)

**output\_file**  
**output\_dir**

**params** - list of params to use

**envir** - environment to evaluate code chunks in

**encoding** - of input file

## Embed code with knitr syntax

### INLINE CODE

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

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

### CODE CHUNKS

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

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

```

### GLOBAL OPTIONS

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

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

```

### IMPORTANT CHUNK OPTIONS

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

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

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

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

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

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

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

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

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

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

**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, aniopts, autodep, background, cache.comments, cache.lazy, cache.rebuild, cache.vars, dev, dev.args, dpi, engine.opts, engine.path, fig.asp, fig.env, fig.ext, fig.keep, fig.lp, fig.path, fig.pos, fig.process, fig.retina, fig.scap, fig.show, fig.showtext, fig.subcap, interval, out.extra, out.height, out.width, prompt, purl, ref.label, render, size, split, tidy.opts

## .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.)

```

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

```

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

**R Markdown** ► Knit to HTML  
Knit to PDF  
Knit to Word  
Knit with Parameters...

## Interactive Documents

Turn your report into an interactive Shiny document in 4 steps

1. Add runtime: shiny to the YAML header.
2. Call Shiny input functions to embed input objects.
3. Call Shiny render functions to embed reactive output.
4. Render with rmarkdown::run or click Run Document in RStudio IDE

```

output: html_document
runtime: shiny

```{r, echo = FALSE}
numericInput("n",
  "How many cars?", 5)
renderTable({
  head(cars, input$n)
})```

```

| How many cars? | |
|----------------|-------|
| speed | dist |
| 1 4.00 | 2.00 |
| 2 4.00 | 10.00 |
| 3 7.00 | 4.00 |
| 4 7.00 | 22.00 |
| 5 8.00 | 16.00 |

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

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



Pandoc's Markdown

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

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

```
equation block:
```

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

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

```
equation block:
```

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

```
block quote
```

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

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

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

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

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

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

```
\textbf{Text ignored in HTML}
<em>HTML ignored in pdfs</em>
```

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

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

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

```
* item 2
```

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

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

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

```
continues after
```

```
2. an interruption
```

```
Term 1
```

```
Definition 1
```

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

- slide bullet 1
- slide bullet 2

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

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

```
***
```

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

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

Set render options with YAML

When you render, R Markdown

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



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

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

output value

creates

| | |
|-----------------------|----------------------------------|
| html_document | html |
| pdf_document | pdf (requires Tex) |
| word_document | Microsoft Word (.docx) |
| odt_document | OpenDocument Text |
| rtf_document | Rich Text Format |
| md_document | Markdown |
| github_document | Github compatible markdown |
| ioslides_presentation | ioslides HTML slides |
| slidy_presentation | slidy HTML slides |
| beamer_presentation | Beamer pdf slides (requires Tex) |

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

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

html tabs

Use tablet css class to place sub-headers into tabs

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



Create a Reusable Template

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

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

template.yaml (see below)

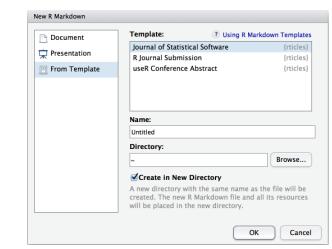
skeleton.Rmd (contents of the template)

any supporting files

3. **Install the package**

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

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



sub-option

description

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

Table Suggestions

Several functions format R data into tables

| Table with kable | | |
|--------------------------------|------|-------|
| <code>eruptions waiting</code> | | |
| 1 | 3.60 | 79.00 |
| 3.600 | 79 | |
| 2 | 1.80 | 54.00 |
| 1.800 | 54 | |
| 3 | 3.33 | 74.00 |
| 3.333 | 74 | |
| 4 | 2.28 | 62.00 |
| 2.283 | 62 | |

| eruptionswaiting | | |
|------------------|-------|----|
| 1 | 3.600 | 79 |
| 2 | 1.800 | 54 |
| 3 | 3.333 | 74 |
| 4 | 2.283 | 62 |

| Table with xtable | | |
|-------------------|---------|-------|
| eruptions | waiting | |
| 1 | 3.60 | 79.00 |
| 3.600 | 79 | |
| 2 | 1.80 | 54.00 |
| 1.800 | 54 | |
| 3 | 3.33 | 74.00 |
| 3.333 | 74 | |
| 4 | 2.28 | 62.00 |
| 2.283 | 62 | |

RStudio IDE :: CHEAT SHEET

Documents and Apps



Check spelling Render output Choose output format Choose output location Insert code chunk

Jump to previous chunk Jump to next chunk Run selected lines Publish to server Show file outline

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

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

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

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

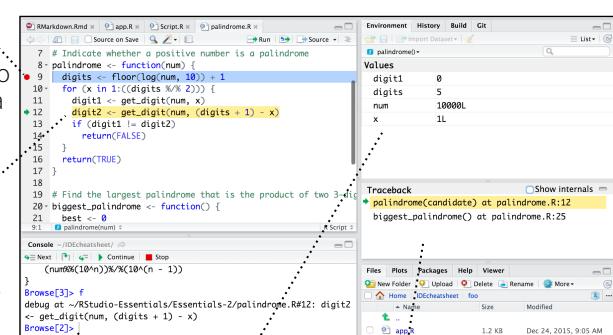
Debug Mode

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

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

Highlighted line shows where execution has paused

Run commands in environment where execution has paused

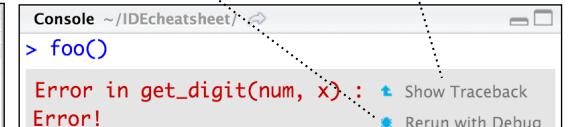


Examine variables in executing environment

Select function in traceback to debug

Launch debugger mode from origin of error

Open traceback to examine the functions that R called before the error occurred



Step through code one line at a time

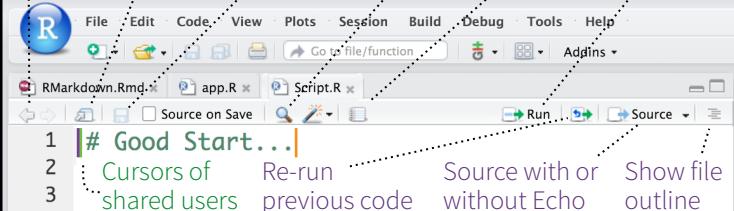
Step into and out of functions to run

Resume execution mode

Quit debug

Write Code

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



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

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

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

Syntax highlighting based on your file's extension

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

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

Jump to function in file

Change file type

Working Directory

Maximize, minimize panes

Drag pane boundaries

R Support

Import data with wizard

History of past commands to run/copy

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

Load workspace Save workspace Delete all saved objects

Search inside environment

Choose environment to display from list of parent environments

Display objects as list or grid

Data Values Functions

150 obs. of 5 variables

a 1

foo function (x)

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

Copy... Move... Export... Set As Working Directory Go To Working Directory

Create folder Upload file Delete file Rename file Change directory

Path to displayed directory

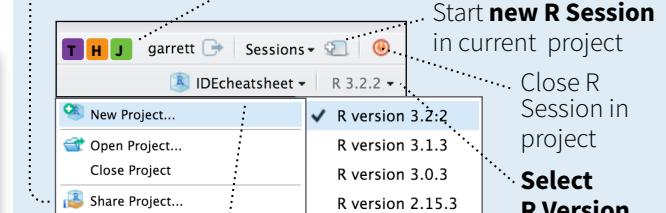
Maximize, minimize panes

Drag pane boundaries

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

Pro Features

Share Project Active shared with Collaborators



Start **new R Session** in current project

Close R Session in project

Select R Version

PROJECT SYSTEM **File > New Project**

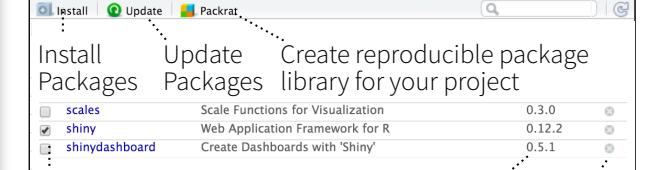
RStudio saves the call history, workspace, and working directory associated with a project. It reloads each when you re-open a project.



Open in recent plots Export plot Delete plot Delete all plots

RStudio opens plots in a dedicated Plots pane

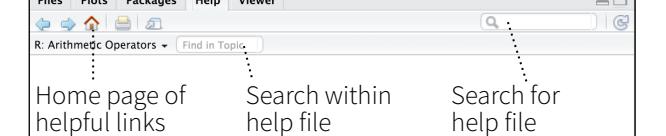
GUI Package manager lists every installed package



Click to load package with **library()**. Unclick to detach package with **detach()**

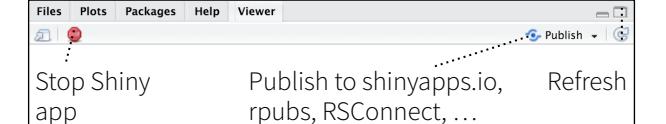
Package version installed Delete from library

RStudio opens documentation in a dedicated Help pane



Home page of helpful links Search within help file Search for help file

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



Stop Shiny app Publish to shinyapps.io, rpubs, RSConnect, ... Refresh

RStudio opens documentation in a dedicated Help pane



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

Filter rows by value or value range Sort by values Search for value



1 LAYOUT

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

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

2 RUN CODE

Search command history

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

Interrupt current command

Clear console

Quit Session (desktop only)

Restart R Session

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

Source the current file

Source with echo

Windows/Linux **Mac**

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

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

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

Ctrl+Enter **Cmd+Enter**

Alt+Enter Option+Enter

Ctrl+Alt+E Cmd+Option+E

Ctrl+Alt+F Cmd+Option+F

Ctrl+Alt+G Cmd+Option+G

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

Ctrl+Shift+Enter Cmd+Shift+Enter

4 WRITE CODE

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

Select

Select Word

Select to Line Start

Select to Line End

Select Page Up/Down

Select to Start/End

Delete Word Left

Delete Word Right

Delete to Line End

Delete to Line Start

Indent

Outdent

Yank line up to cursor

Yank line after cursor

Insert yanked text

Insert <->

Insert %>%

Show help for function

Show source code

New document

New document (Chrome)

Open document

Save document

Close document

Close document (Chrome)

Close all documents

Extract function

Extract variable

Reindent lines

(Un)Comment lines

Reflow Comment

Reformat Selection

Select within braces

Show Diagnostics

Transpose Letters

Move Lines Up/Down

Copy Lines Up/Down

Add New Cursor Above

Add New Cursor Below

Move Active Cursor Up

Move Active Cursor Down

Find and Replace

Use Selection for Find

Replace and Find

Windows /Linux

Tab or Ctrl+Space

↑/↓

Enter, Tab, or →

Esc

Ctrl+Z

Ctrl+Shift+Z

Ctrl+X

Cmd+C

Cmd+V

Ctrl+A

Ctrl+D

Shift+[Arrow]

Option+Shift+ ←/→

Alt+Shift+←

Alt+Shift+→

Shift+PageUp/Down

Shift+Alt+↑/↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Ctrl+Shift+M

F1

F2

Cmd+Shift+N

Ctrl+Alt+Shift+N

Ctrl+O

Ctrl+S

Cmd+S

Ctrl+W

Cmd+Option+W

Ctrl+Shift+W

Ctrl+Alt+X

Cmd+Option+X

Ctrl+Alt+V

Cmd+Option+V

Ctrl+I

Ctrl+Shift+C

Ctrl+Shift+/

Ctrl+Shift+A

Ctrl+Shift+E

Ctrl+Shift+E

Cmd+Shift+Opt+P

Ctrl+T

Alt+↑/↓

Shift+Alt+↑/↓

Cmd+Option+↑/↓

Ctrl+Option+Up

Ctrl+Option+Down

Ctrl+Option+Shift+Up

Ctrl+Opt+Shift+Down

Ctrl+Opt+Shift+J

Mac

Tab or Cmd+Space

↑/↓

Enter, Tab, or →

Esc

Cmd+Z

Cmd+Shift+Z

Cmd+X

Cmd+C

Cmd+V

Cmd+A

Cmd+D

Shift+[Arrow]

Option+Shift+ ←/→

Cmd+Shift+←

Cmd+Shift+→

Shift+PageUp/Down

Shift+Alt+↑/↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Cmd+Shift+M

F1

F2

Cmd+Shift+N

Cmd+Shift+Opt+N

Cmd+O

Cmd+S

Cmd+W

Cmd+Option+W

Cmd+Shift+W

Cmd+Alt+X

Cmd+Option+X

Cmd+I

Cmd+Shift+C

Cmd+Shift+/

Cmd+Shift+A

Cmd+Shift+E

Cmd+Shift+E

Cmd+Shift+Opt+P

Cmd+T

Option+↑/↓

Shift+Alt+↑/↓

Cmd+Option+↑/↓

Ctrl+Option+Up

Ctrl+Option+Down

Ctrl+Option+Shift+Up

Ctrl+Opt+Shift+Down

Ctrl+Opt+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/

<h2

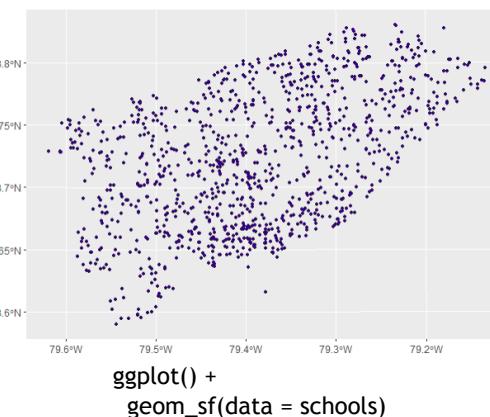
Spatial manipulation with sf: : CHEAT SHEET



The sf package provides a set of tools for working with geospatial vectors, i.e. points, lines, polygons, etc.

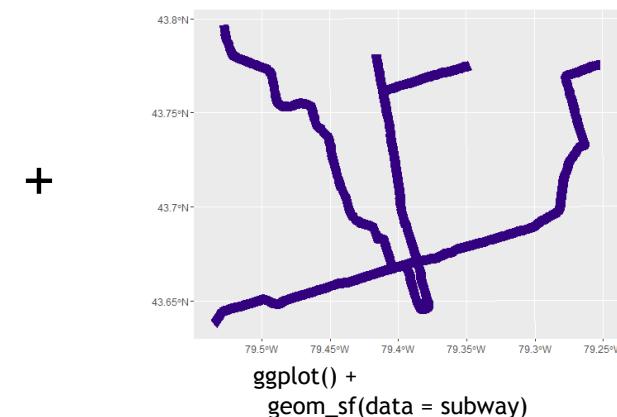
Geometric confirmation

- 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



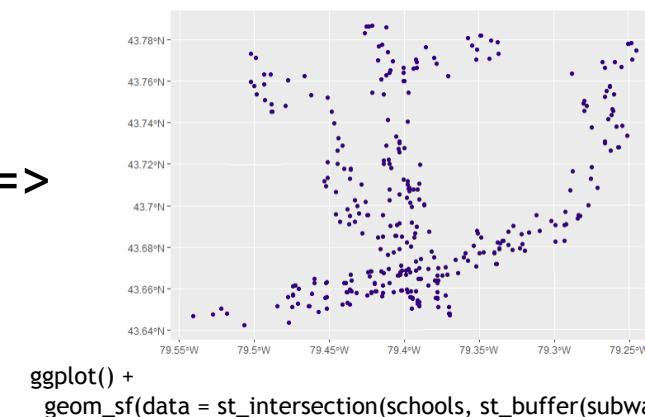
Geometric operations

- 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

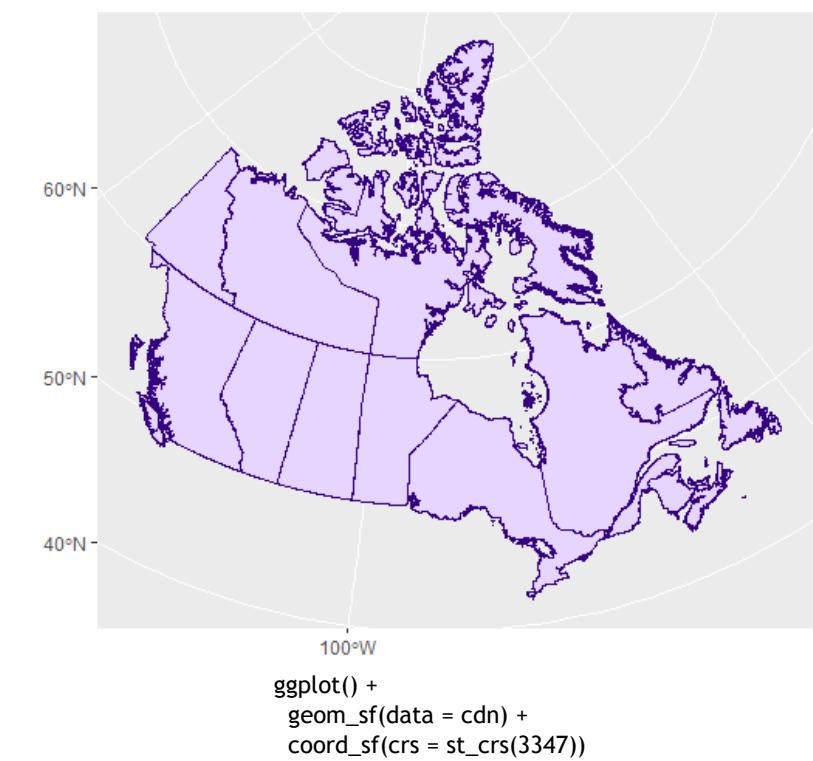
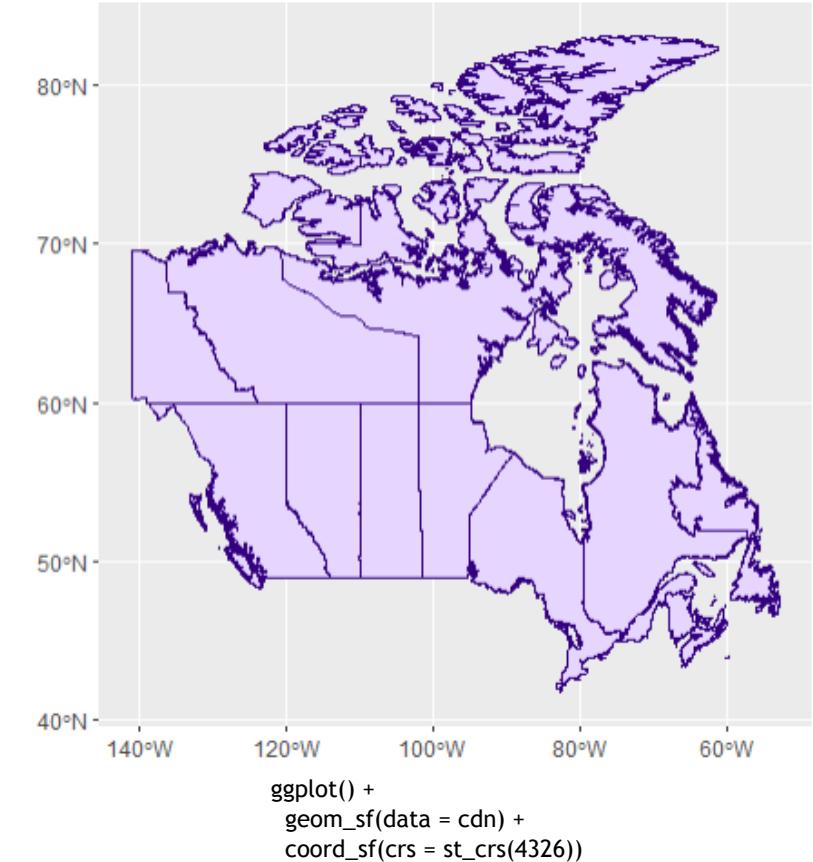
- ↳ `st_contains(x, y, ...)` Identifies if x is within y (i.e. point within polygon)
- ↳ `st_crop(x, y, ..., xmin, ymin, xmax, ymax)` Creates geometry of x that intersects a specified rectangle
- ↳ `st_difference(x, y)` Creates geometry from x that does not intersect with y
- ↳ `st_intersection(x, y)` 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

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

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





Shiny :: CHEAT SHEET

Basics

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



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

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(
  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** - 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
at www.rstudio.com/products/shiny-server/



Building an App

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

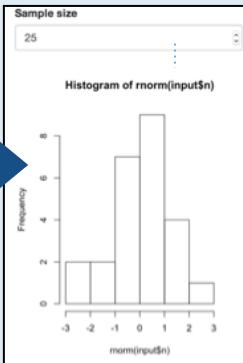
Add inputs to the UI with `*Input()` functions

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

```
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)
```



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.

● ● ● **app-name**

- app.R** (optional) defines objects available to both ui.R and server.R
- global.R** (optional) used in showcase mode
- DESCRIPTION** (optional) data, scripts, etc.
- README** (optional) directory of files to share with web browsers (images, CSS, .js, etc.) Must be named "**www**"
- <other files>

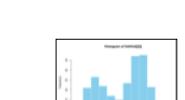
The directory name is the name of the app
 (optional) defines objects available to both ui.R and server.R
 (optional) used in showcase mode
 (optional) data, scripts, etc.
 (optional) directory of files to share with web browsers (images, CSS, .js, etc.) Must be named "www"

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)`



`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)`

works with

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

`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**.

Action

Link

- Choice 1
 Choice 2
 Choice 3
 Check me

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()`)

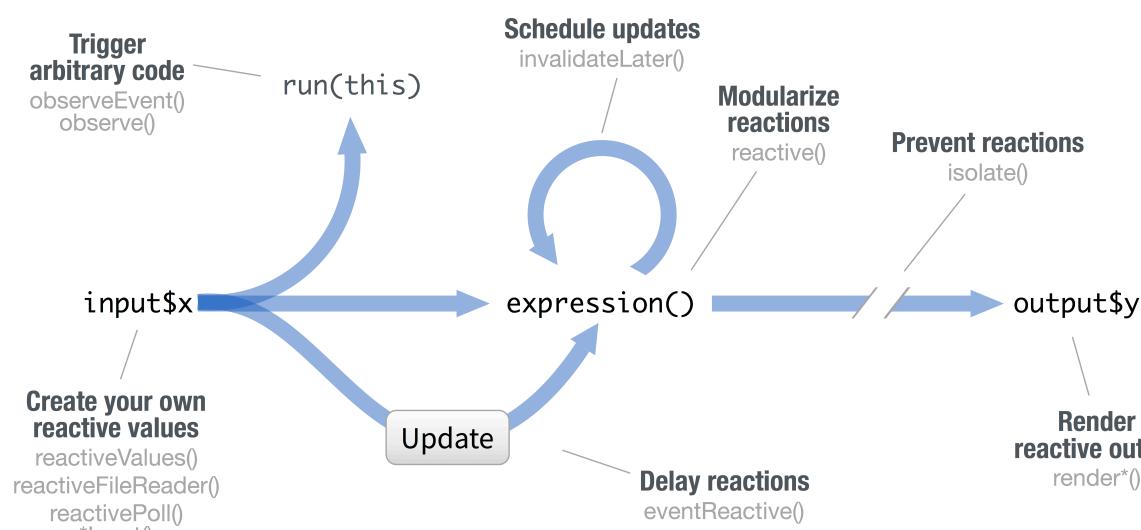
sliderInput(inputId, label, min, max, value, step, round, format, locale, ticks, animate, width, sep, pre, post)

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
}

reactiveValues() creates a list of reactive values whose values you can set.
```

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.

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>**

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>
```



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\$ol | tags\$sub |
| tags\$article | tags\$details | tags\$hrt | tags\$optgroup | tags\$summary |
| tags\$aside | tags\$dfn | tags\$HTML | tags\$option | tags\$sup |
| tags\$audio | tags\$div | tags\$si | tags\$output | tags\$table |
| tags\$b | tags\$dl | tags\$iframe | tags\$p | tags\$tbody |
| tags\$base | tags\$dt | tags\$img | tags\$param | tags\$td |
| tags\$bdi | tags\$em | tags\$input | tags\$pre | tags\$textarea |
| tags\$bdo | tags\$embed | tags\$ins | tags\$progress | tags\$tfoot |
| tags\$blockquote | tags\$eventsouce | tags\$kbd | tags\$q | tags\$th |
| tags\$body | tags\$fieldset | tags\$keygen | tags\$ruby | tags\$thead |
| tags\$br | tags\$figcaption | tags\$label | tags\$rp | tags\$time |
| tags\$button | tags\$figure | tags\$legend | tags\$rt | tags\$title |
| tags\$canvas | tags\$footer | tags\$li | tags\$ss | tags\$tr |
| tags\$caption | tags\$form | tags\$link | tags\$amp | tags\$track |
| tags\$cite | tags\$h1 | tags\$mark | tags\$script | tags\$u |
| tags\$code | tags\$h2 | tags\$map | tags\$section | tags\$ul |
| tags\$col | tags\$h3 | tags\$menu | tags\$select | tags\$var |
| tags\$colgroup | tags\$h4 | tags\$meta | tags\$small | tags\$video |
| tags\$command | tags\$h5 | tags\$meter | tags\$source | tags\$wbr |

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

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



To include a CSS file, use **includeCSS()**, or

1. Place the file in the **www** subdirectory
2. Link to it with

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

To include JavaScript, use **includeScript()** or

1. Place the file in the **www** subdirectory
2. Link to it with

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

IMAGES To include an image

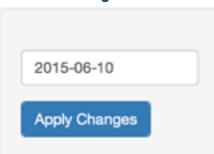
1. Place the file in the **www** subdirectory
2. Link to it with **img(src=<file name>")**

Layouts



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

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



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

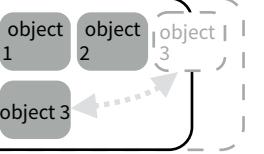
Organize panels and elements into a layout with a layout function. Add elements as arguments of the layout functions.

fluidRow()



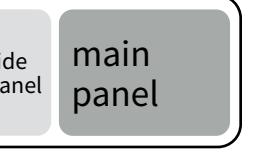
```
ui <- fluidPage(
  fluidRow(column(width = 4),
    column(width = 2, offset = 3)),
  fluidRow(column(width = 12)))
```

flowLayout()



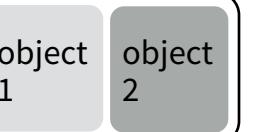
```
ui <- fluidPage(
  flowLayout(object1,
  object2,
  object3))
```

sidebarLayout()



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

splitLayout()

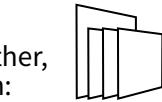


```
ui <- fluidPage(
  splitLayout(object1,
  object2))
```

verticalLayout()



```
ui <- fluidPage(
  verticalLayout(object1,
  object2,
  object3))
```



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

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



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



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



Data Science in Spark with sparklyr :: CHEAT SHEET

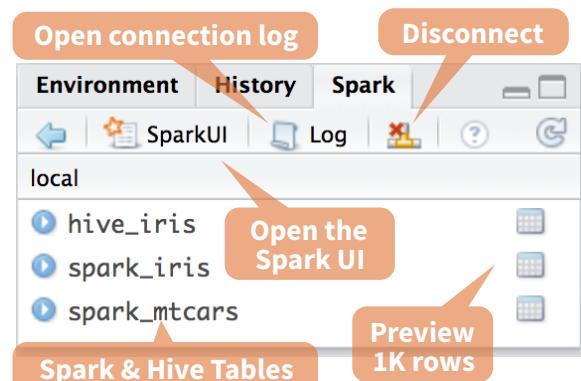


Intro

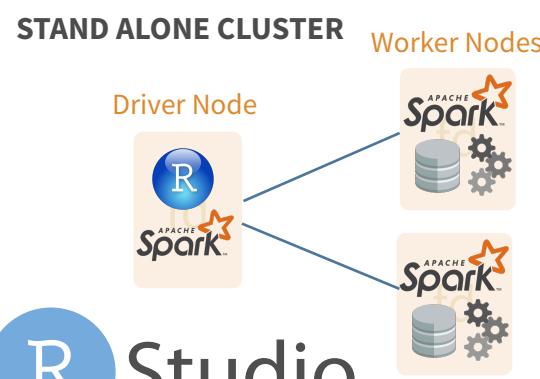
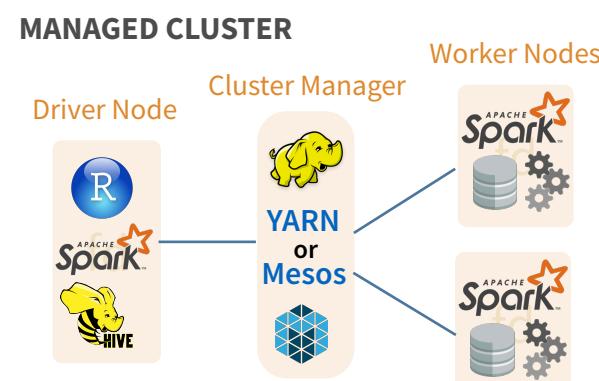
sparklyr is an R interface for Apache Spark™, it provides a complete **dplyr** backend and the option to query directly using **Spark SQL** statement. With sparklyr, you can orchestrate distributed machine learning using either **Spark's MLLib** or **H2O** Sparkling Water.

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

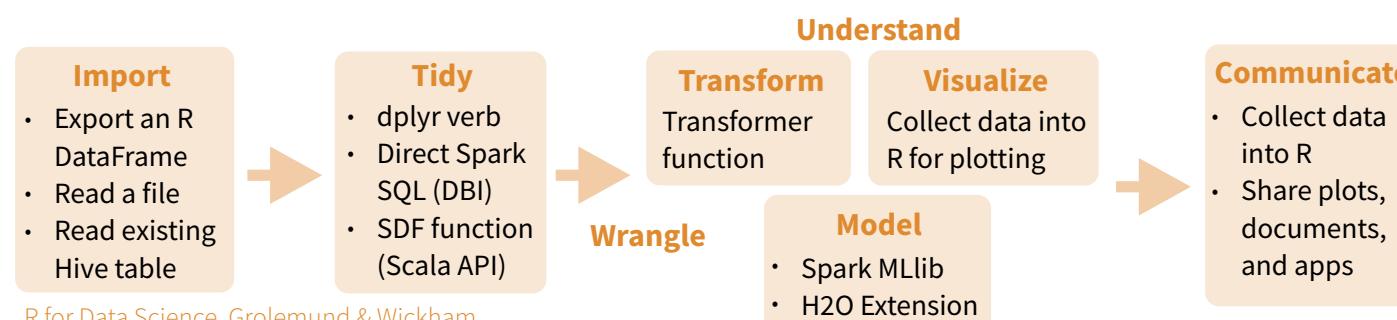


Cluster Deployment



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 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])`

ON A YARN MANAGED CLUSTER

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

ON A SPARK STANDALONE CLUSTER

1. Install RStudio Server or RStudio Pro on one of the existing nodes or a server in the same LAN
2. Install a local version of Spark:
`spark_install (version = "2.0.1")`
3. Open a connection
`spark_connect(master="spark://host:port", version = "2.0.1", spark_home = spark_home_dir())`

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"))
```

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

Create reference to Spark table

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

```
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

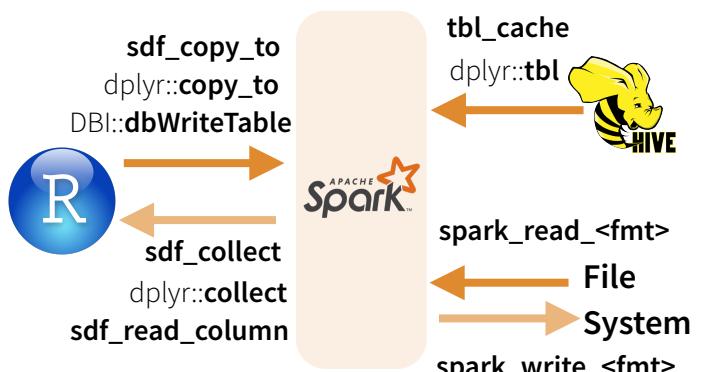
Arguments that apply to all functions: `x, path`

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_jobj()` 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())` Same options for: `ml_gradient_boosted_trees`

`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())`
Same options for: `ml_logistic_regression`

`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")`

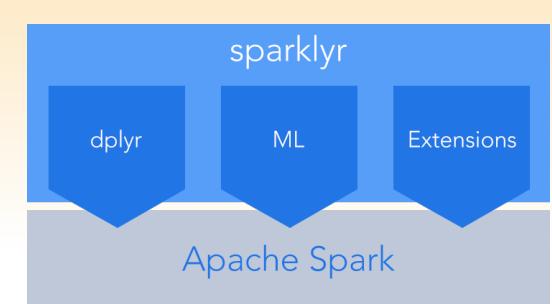
`ml_classification_eval(predicted_tbl_spark, label, predicted_lbl, metric = "f1")`

`ml_tree_feature_importance(sc, model)`

sparklyr

is an R interface
for

Apache Spark



String manipulation with stringr :: CHEAT SHEET



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

Detect Matches

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

Subset Strings

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

Manage Lengths

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

Mutate Strings

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

Join and Split

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

Order Strings

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

Helpers

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

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

R Syntax Comparison :: CHEAT SHEET

Dollar sign syntax

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

SUMMARY STATISTICS:

one continuous variable:

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

one categorical variable:

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

two categorical variables:

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

one continuous, one categorical:

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

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

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

PLOTTING:

one continuous variable:

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

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

one categorical variable:

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

two continuous variables:

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

two categorical variables:

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

one continuous, one categorical:

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

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

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

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

WRANGLING:

subsetting:

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

making a new variable:

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

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

Formula syntax

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

SUMMARY STATISTICS:

one continuous variable:

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

one categorical variable:

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

two categorical variables:

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

one continuous, one categorical:

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

tilde

PLOTTING:

one continuous variable:

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

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

one categorical variable:

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

two continuous variables:

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

two categorical variables:

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

one continuous, one categorical:

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

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

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

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

Tidyverse syntax

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

SUMMARY STATISTICS:

one continuous variable:

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

one categorical variable:

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

two categorical variables:

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

one continuous, one categorical:

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

the pipe

PLOTTING:

one continuous variable:

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

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

one categorical variable:

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

two continuous variables:

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

two categorical variables:

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

one continuous, one categorical:

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

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

WRANGLING:

subsetting:

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

making a new variable:

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

R Syntax Comparison :: CHEAT SHEET

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

The three most prevalent R syntaxes are:

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

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

Internet research tip:

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

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



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

Even more ways to say the same thing

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

quickplot

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

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

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

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

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

ggplot

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

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

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

plus adds layers

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

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

ggformula

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

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

formulas in base plots

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

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

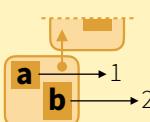
Tidy evaluation with rlang :: CHEAT SHEET



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

abs (1)

Constant - a bare value (i.e. an atomic vector of length 1). `is_bare_atomic(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
q
a + b, [a b]

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

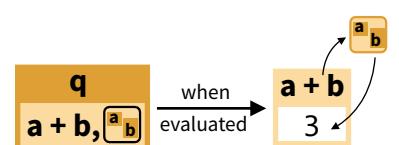
Expression Vector - a list of pieces of quoted code created by base R's `expression` and `parse` functions. Not to be confused with **expression**.

R Studio

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.

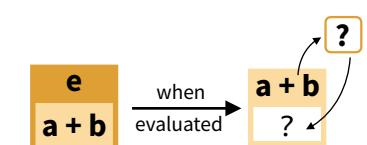
`rlang::quo(expr)` Quote contents as a quosure. Also **quos** to quote multiple expressions. `a <- 1; b <- 2; q <- quo(a + b); qs <- quos(a, b)`

`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.

• • •

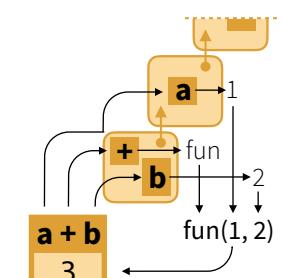
`rlang::parse_expr(x)` Convert a string to an expression. Also **parse_exprs**, **sym**, **parse_quo**, **parse_quos**. `e <- parse_expr("a+b")`

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)`

Evaluation



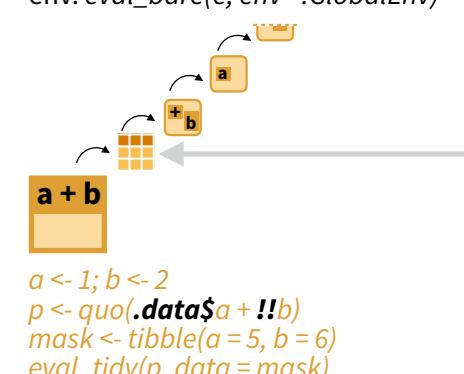
To evaluate an expression, R :

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.

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

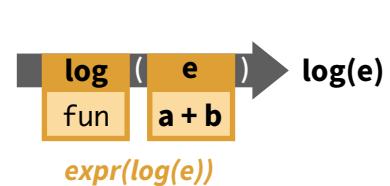
```
call2("log", x = 4, base = 2)
call2("log", !!!args)
exec("log", x = 4, base = 2)
exec("log", !!!args)
```



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

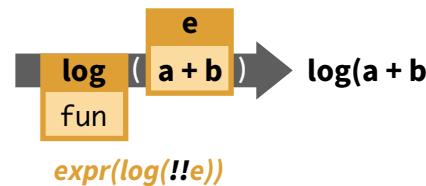
QUOTATION

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



QUASIQUOTATION

Quoting *some* parts of an expression while evaluating and then inserting the results of others (**unquoting** others).
e <- expr(a + b)

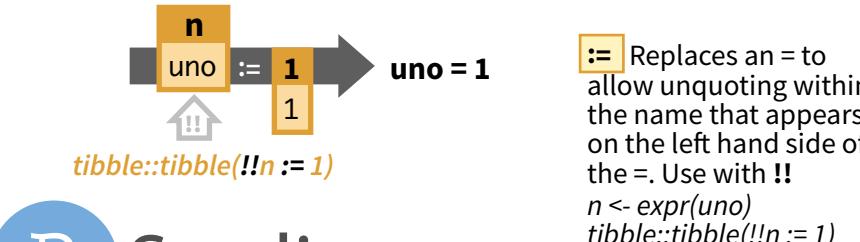
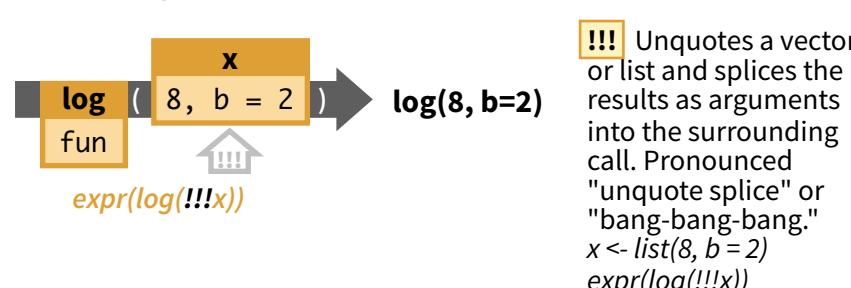
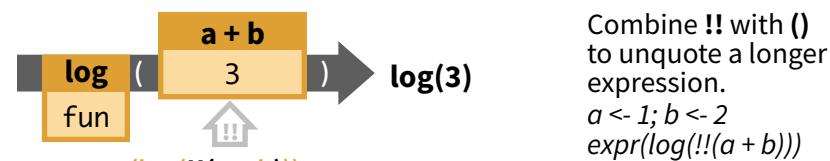
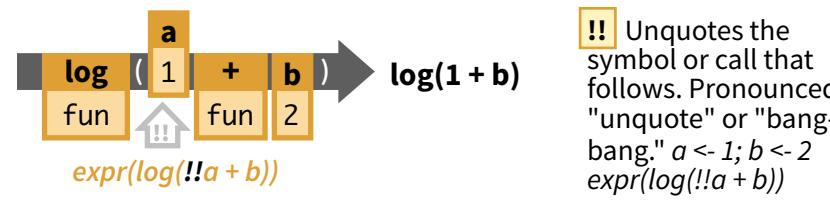


rlang provides !!, !!!, and := for doing quasiquotation.

!!, !!!, 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)) 2
}
```

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)
  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 quo to use
3. **Evaluate** the new expression/ quo 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
}
```

1
2

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
  mutate_y <- function(df) {
    dplyr::mutate(df, y = .data$a + 1)
  }
```

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