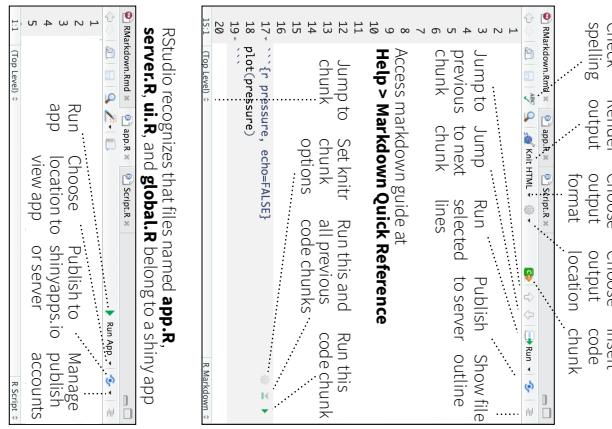


RStudio IDE :: CHEAT SHEET

Documents and Apps

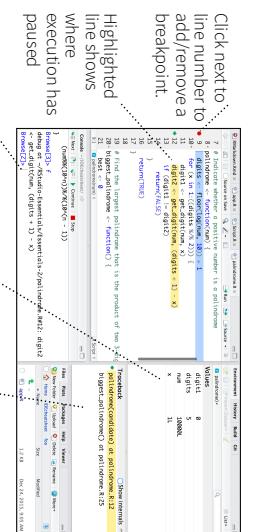


Write Code



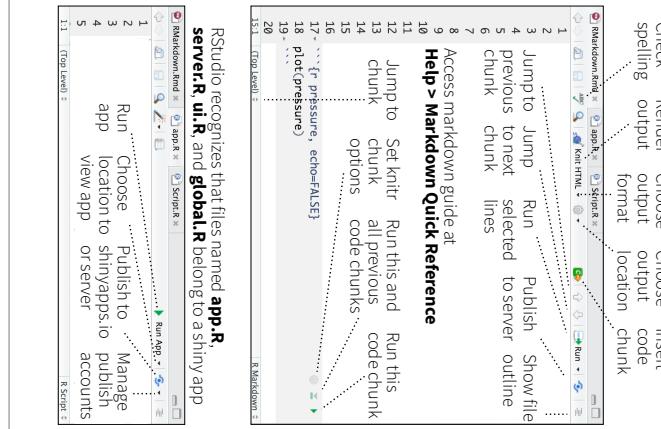
Debug Mode

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

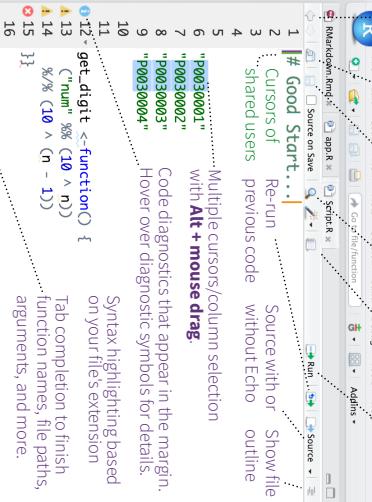


Version Control

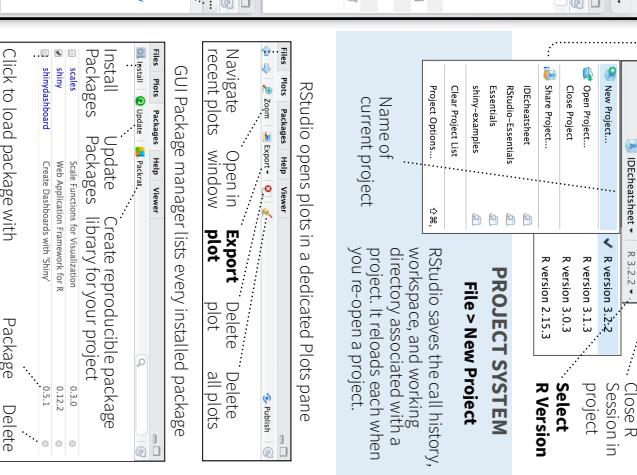
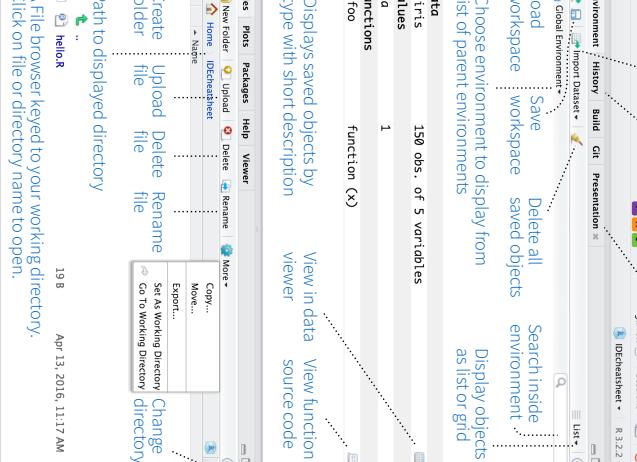
Turn on at **Tools > Project Options > Git/SVN**



R Support



Pro Features





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

2 RUN CODE

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

Interrupt current command

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

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 %>%

Types

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

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

Maths Functions

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

Variable Assignment

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

The Environment

<code>ls()</code>	List all variables in the environment.
<code>rm(x)</code>	Remove x from the environment.
<code>rm(list = ls())</code>	Remove all variables from the environment.

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

Matrices	
<code>m <- matrix(x, nrow = 3, ncol = 3)</code>	Create a matrix from x.
<code>t(m)</code>	Transpose
<code>m %*% n</code>	Matrix Multiplication
<code>solve(m, n)</code>	Find x in: $m \cdot x = n$
<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.

Lists

```
l <- list(x = 1:5, y = c('a', 'b'))
```

A list is a collection of elements which can be of different types.

```
l[[2]] l[[1]] l$x l['y']
```

<code>l[[2]]</code>	New list with only the first element.
<code>l[[1]]</code>	Element named x.

Also see the `dplyr` package.

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

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

List subsetting

	x	y
1	a	df\$x
2	b	df[[2]]
3	c	

Linear model.

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

Generalised linear model.

summary

Get more detailed information about a model.

pairwise.t.test

Perform a t-test for paired data.

Distributions

Random Variates

Density Function

Cumulative Distribution

Quantile

`rnorm`

`dnorm`

`pnorm`

`qnorm`

Normal

Poisson

Binomial

Uniform

`rpois`

`dpois`

`ppois`

`dpois`

`rbinom`

`dbinom`

`pbinom`

`qbinom`

`rnorm`

`dnorm`

`pnorm`

`qnorm`

`runif`

`dunif`

`unif`

`pnunif`

`qunif`

Strings

Also see the `stringr` package.

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

Join multiple vectors together.

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

Join elements of a vector together.

`grep(pattern, x)`

Find regular expression matches in x.

`gsub(pattern, replace, x)`

Replace matches in x with a string.

`toupper(x)`

Convert to uppercase.

`tolower(x)`

Convert to lowercase.

`nchar(x)`

Number of characters in a string.

factor(x)

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

cut(x, breaks = 4)

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

Statistics

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

Test for a difference between means.

`t.test(x, y)`

Test for a difference between proportions.

`prop.test`

Test for a difference between two proportions.

`pairwise.t.test`

Perform a t-test for paired data.

`aov`

Analysis of variance.

Distributions

`glm`

Generalised linear model.

`summary`

Get more detailed information about a model.

`head(df)`

See the first 6 rows.

`tail(df)`

See the last 6 rows.

`View(df)`

See the full data frame.

`head(df)`

See the first 6 rows.

`nrow(df)`

Number of rows.

`cbind`

Bind columns.

`rbind`

Bind rows.

`ncol(df)`

Number of columns.

`df[2, 1]`

Number of columns and rows.

`df[2, 2]`

Number of columns and rows.

Also see the `ggplot2` package.

Dates

See the `lubridate` package.

Advanced R

Cheat Sheet

Created by: Arianne Colton and Sean Chen

Environment Basics

- Environment – **Data structure** (with two components below) that powers lexical scoping

Create environment: `env1<-new.env()`

- Named list ("Bag of names") – each name points to an object stored elsewhere in memory.

If an object has no names pointing to it, it gets automatically deleted by the garbage collector.

- Access with: `ls('env1')`
- Parent environment** – used to implement lexical scoping. If a name is not found in an environment, then R will look in its parent (and so on).
- Access with: `parent.env('env1')`

Four special environments

- Empty environment – ultimate ancestor of all environments

- Parent: none
- Access with: `emptyenv()`

- Base environment - environment of the base package

- Parent: empty environment
- Access with: `baseenv()`

- Global environment – the interactive workspace that you normally work in

- Parent: environment of last attached package
- Access with: `globalenv()`

- Current environment – environment that R is currently working in (may be any of the above and others)

- Parent: empty environment
- Access with: `environment()`

Search Path

Search path – mechanism to look up objects, particularly functions.

- Access with: `search()` – lists all parents of the global environment (see Figure 1)
- Access any environment on the search path:
`as.environment('package:base')`

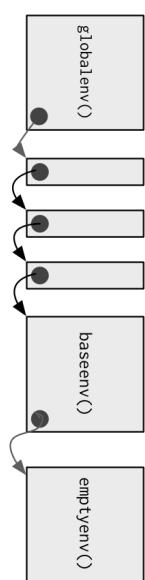


Figure 1 – The Search Path

- Mechanism : always start the search from global environment, then inside the latest attached package environment.
- New package loading with `library()`/`require()` : new package is attached right after global environment. (See Figure 2)
- Name conflict in two different package : functions with the same name, latest package function will get called.

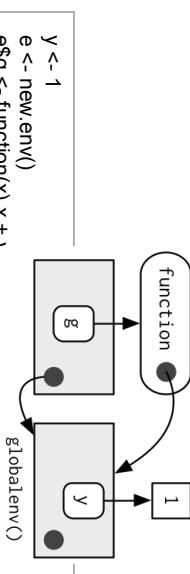
search():

```
'.GlobalEnv' ... 'Autoloads' 'package:base'
```

```
'GlobalEnv' 'package:reshape2' ... 'Autoloads' 'package:base'
```

NOTE: Autoloads : special environment used for saving memory by only loading package objects (like big datasets) when needed

Figure 2 – Package Attachment



Example (for enclosing and binding environment):

- Enclosing environment - an environment where the function is created. It determines how function finds value.
 - Enclosing environment never changes, even if the function is moved to a different environment.
- Binding environment - all environments that the function has a binding to. It determines how we find the function.
 - Access with: `environment('func1')`

Environments

Function Environments

- Enclosing environment - an environment where the function is created. It determines how function finds value.

- Execution environment - new created environments to host a function call execution.
 - Two parents :
 - Enclosing environment of the function
 - Calling environment of the function
 - Execution environment is thrown away once the function has completed.

- Calling environment - environments where the function was called.

- Access with: `parent.frame('func1')`
- Dynamic scoping :
 - About : look up variables in the calling environment rather than in the enclosing environment
- Usage : most useful for developing functions that aid interactive data analysis

Warning: If `<-` doesn't find an existing variable, it will create one in the global environment.

Data Structures

	Homogeneous	Heterogeneous
1d	Atomic vector	List
2d	Matrix	Data frame
nd	Array	

Note: R has no 0-dimensional or scalar types. Individual numbers or strings are actually vectors of length one, NOT scalars.

Human readable description of any R data structure :

```
str(variable)
```

Every Object has a mode and a class

1. **Mode**: represents how an object is stored in memory

- 'type' of the object from R's point of view

2. **Class**: represents the object's abstract type

- 'type' of the object from R's object-oriented programming point of view

- Access with: `class()`

	<code>typeof()</code>	<code>class()</code>
strings or vector of strings	character	character
numbers or vector of numbers	numeric	numeric
list	list	list
data.frame	list	data.frame

Factors

1. Factors are built on top of integer vectors using two attributes :

```
class(x) -> 'factor'  
levels(x) # defines the set of allowed values
```

2. Useful when you know the possible values a variable may take, even if you don't see all values in a given dataset.

Warning on Factor Usage:

1. Factors look and often behave like character vectors, they are actually integers. Be careful when treating them like strings.
2. Most data loading functions automatically convert character vectors to factors. (Use argument `stringAsFactors = FALSE` to suppress this behavior)

Object Oriented (OO) Field Guide

Object Oriented Systems

R has three object oriented systems :

1. **S3** is a very casual system. It has no formal definition of classes. It implements generic function OO.

- **Generic-function OO** - a special type of function called a generic function decides which method to call.

```
Example: drawRect(canvas, "blue")
```

```
Language: R
```

- **Message-passing OO** - messages (methods) are sent to objects and the object determines which function to call.

Example:	<code>canvas.drawRect('blue')</code>
Language:	Java, C++, and C#

2. **S4** works similarly to S3, but is more formal.

- Two major differences to S3 :
 - **Formal class definitions** - describe the representation and inheritance for each class, and has special helper functions for defining generics and methods.
 - **Multiple dispatch** - generic functions can pick methods based on the class of any number of arguments, not just one.

3. **Reference classes** are very different from S3 and S4:

- **Implements message-passing OO** - methods belong to classes, not functions.
- **Notation** - \$ is used to separate objects and methods, so method calls look like `canvas$drawRect('blue')`.

Base Type (C Structure)

R base types - the internal C-level types that underlie the above OO systems.

- **Includes** : atomic vectors, list, functions, environments, etc.
- **Useful operation** : Determine if an object is a base type (Not S3, S4 or RC) `is.object(x)` returns FALSE

About S3:

- R's first and simplest OO system
- Only OO system used in the base and stats package

- Methods belong to functions, not to objects or classes.
- Methods have a generic method for the class.

Notation :

- `generic.class()`

<code>mean.Date()</code>	Date method for the generic - <code>mean()</code>
--------------------------	---

Useful 'Generic' Operations

- Get all methods that belong to the 'mean' generic:
- List all generics that have a method for the 'Date' class :
 - `methods(class = 'Date')`

4. **S3 objects** are usually built on top of lists, or atomic vectors with attributes.

- Factor and data frame are S3 class
- Useful operations:

Check if object is an S3 object	<code>is.object(x) & !isS4(x)</code> or <code>pryr::objtype()</code>
Check if object inherits from a specific class	<code>inherits(x, 'classname')</code>
Determine class of any object	<code>class(x)</code>

Functions

Function Basics

Functions – objects in their own right

- All R functions have three parts:

<code>body()</code>	code inside the function
<code>formals()</code>	list of arguments which controls how you can call the function
<code>environment()</code>	"map" of the location of the function's variables (see "Enclosing Environment")

Every operation is a function call

- `+`, `for`, `if`, `[`, `$`, `{`
- `x + y` is the same as ``+`(x, y)`

Note: the backtick ````, lets you refer to functions or variables that have otherwise reserved or illegal names.

Lexical Scoping

What is Lexical Scoping?

- Looks up value of a symbol. (see "Enclosing Environment")
- `findGlobals()` - lists all the external dependencies of a function

```
f <- function() x + 1
```

`codetools::findGlobals(f)`

`> '+' X'`

`environment(f) <- emptyEnv()`

`f()`
error in f(): could not find function "`+`"

- R relies on lexical scoping to find everything, even the `+` operator.

Return Values

What are Replacement Functions?

- Act like they modify their arguments in place, and have the special name `xxx <-`
- Actually create a modified copy. Can use `pryr::address()` to find the memory address of the underlying object

```
'second<- <- function(x, value) {
  x[2] <- value
}
```

`x`
`X <- 1:10`

`second(x) <- 5L`

Function Arguments

Arguments – passed by reference and copied on modify

- Arguments are matched first by exact name (perfect matching), then by prefix matching, and finally by position.
- Check if an argument was supplied : `missing()`
- Lazy evaluation – since `x` is not used `stop("This is an error!")` never get evaluated.

```
i <- function(a, b) {
  missing(a) >- # return true or false
}

f <- function(x) {
  10
}
f(stop("This is an error!")) >- 10
```

- Force evaluation

```
f <- function(x) {
  force(x)
  10
}
```

- Default arguments evaluation

```
f <- function(x = ls()) {
  a <- 1
  x
}
```

- Useful way of providing a default value in case the output of another function is `NULL`:

```
'%>+%' <- function(a, b) paste0(a, b)
'new' %>+%' string'
```

Influx Functions

What are Influx Functions?

- Function name comes in between its arguments, like `+ or -`
- All user-created infix functions must start and end with `%>`.
- Useful way of providing a default value in case the output of another function is `NULL`:

```
'%>+%' <- function(a, b) if (is.null(a)) a else b
function _that_might_return_null() %>%> default value
```

Primitive Functions

What are Primitive Functions?

- Call C code directly with `.Primitive()` and contain no R code
- `formals()`, `body()` and `environment()` are all `NULL`
- Only found in base package
- More efficient since they operate at a low level

`print(sum) :`
`> function (... , na.rm = FALSE) .Primitive('sum')`

Subsetting

Subsetting returns a copy of the original data, NOT copy-on modified

Simplifying vs. Preserving Subsetting

1. Simplifying subsetting

- Returns the **simpliest** possible data structure that can represent the output

2. Preserving subsetting

- Keeps the structure of the output the **same** as the input.
- When you use `drop = FALSE`, it's preserving

Simplifying*	Preserving
<code>Vector</code>	<code>x[[1]]</code>
<code>List</code>	<code>x[[1]]</code>
<code>Factor</code>	<code>x[1:4, drop = T]</code>
<code>Array</code>	<code>x[1,] or x[, 1]</code>
<code>Data frame</code>	<code>x[, 1] or x[[1]]</code>
	<code>x[1]</code>

Simplifying behavior varies slightly between different data types:

1. Atomic Vector

- `x[[1]]` is the same as `x[1]`

2. List

- `[]` always returns a `list`
- Use `[[]]` to get `list` contents, this returns a single value piece out of a `list`

3. Factor

- Drops any unused levels but it remains a factor class

4. Matrix or Array

- If any of the dimensions has length 1, that dimension is dropped

5. Data Frame

- If `output` is a single column, it returns a vector instead of a data frame

Data Frame Subsetting

Data Frame – possesses the **characteristics of both lists and matrices**. If you subset with a single vector, they behave like lists; if you subset with two vectors, they behave like matrices

1. Subset with a single vector

- Behave like lists

```
df1[c('col1', 'col2')]
```

2. Subset with two vectors

- Behave like matrices

```
df1[, c('col1', 'col2')]
```

The results are the same in the above examples, however, results are different if subsetting with only one column. (see below)

1. Behave like matrices

```
str(df[, 'col1']) -> int[1:3]
```

2. Behave like lists

```
str(df1['col1']) -> 'data.frame'
```

• Result: the result remains a data frame of 1 column

\$ Subsetting Operator

1. About Subsetting Operator

- Useful shorthand for `[[]]` combined with character subsetting

```
x$y is equivalent to x[[y], exact = FALSE]
```

2. Difference vs. `[[`

- `$` does partial matching, `[[` does not

```
x <- list(xbc = 1)
x$a -> 1 # since "exact = FALSE"
x[[a]] -> # would be an error
```

3. Common mistake with `$`

- Using it when you have the name of a column stored in a variable

```
var <- 'cyl'
x$var
# doesn't work, translated to x[['var']]
# Instead use x[[var]]
```

Examples

1. Lookup tables (character subsetting)

```
x <- c('m', 'f', 'u', 'r', 'f', 'm', 'm')
lookup <- c(m = 'Male', f = 'Female', u = NA)
> m f u f f m n
> 'Male' 'Female' NA 'Female' 'Female' 'Male' 'Male'
uname(lookup[x])
```

2. Matching and merging by hand (integer subsetting)

Lookup table which has multiple columns of information:

```
grades <- c(1, 2, 2, 3, 1)
info <- data.frame(
  grade = 3:1,
  desc = c('Excellent', 'Good', 'Poor'),
  fail = c(F, F, T)
)
First Method
id <- match(grades, info$grade)
info[id,]
Second Method
rownames(info) <- info$grade
info[as.character(grades),]
```

3. Expanding aggregated counts (integer subsetting)

Problem: a data frame where identical rows have been collapsed into one and a count column has been added

Solution: `rep()` and integer subsetting make it easy to uncollapse the data by subsetting with a repeated row index: `rep(x, y)` `rep` replicates the values in `x`, `y` times.

4. Removing columns from data frames (character subsetting)

There are two ways to remove columns from a data frame:

Set individual columns to NULL	df1\$col3 <- NULL
Subset to return only columns you want	df1[c('col1', 'col2')]

5. Selecting rows based on a condition (logical subsetting)

This is the most commonly used technique for extracting rows out of a data frame.

```
df1[df1$cold == 5 & df1$cold2 == 4, ]
```

Subsetting continued

Boolean Algebra vs. Sets (Logical and Integer Subsetting)

- Using **integer subsetting** is more effective when:
 - You want to find the first (or last) TRUE.
 - You have very few TRUES and very many FALSEs; a set representation may be faster and require less storage.
- which()** - conversion from boolean representation to integer representation

```
which(c(T, F, T F)) -> 1 3
```
- Integer representation length : is always <= boolean representation length
- Common mistakes :
 - Use **x[which(y)]** instead of **x[y]**
 - x[which(y)]** is not equivalent to **x[y]**

Recommendation:

Avoid switching from logical to integer subsetting unless you want, for example, the first or last TRUE value

Debugging, Condition Handling and Defensive Programming

Debugging Methods

1. traceback() or RStudio's error inspector

- Lists the sequence of calls that lead to the error

2. browser() or RStudio's breakpoints tool

- Opens an interactive debug session at an arbitrary location in the code

3. options(error = browser) or RStudio's "Rerun with Debug" tool

- Opens an interactive debug session where the error occurred

Error Options:

options(error = recover)

- Difference vs. 'browser': can enter environment of any of the calls in the stack

options(error = dump_and_quit)

- Equivalent to 'recover' for non-interactive mode

options(error = fast.dump.rda)

- Creates **fast.dump.rda** in the current working directory

In batch R process :

```
dump_and_quit <- function(){  
  # Save debugging info to file  
  last.dump.rda  
  dump.frames(to.file = TRUE)  
}  
  
# Quit R with error status  
q(status = 1)  
}  
options(error = dump_and_quit)
```

Condition Handling of Expected Errors

1. Communicating potential problems to users:

- stop()**
 - Action : raise fatal error and force all execution to terminate

- warning()**
 - Action : generate warnings to display potential problems

- message()**
 - Action : generate messages to give informative output
 - Example usage : when some of elements of a vectorized input are invalid

2. Handling conditions programmatically:

- try()**
 - Action : gives you the ability to continue execution even when an error occurs

- tryCatch()**
 - Action : lets you specify handler functions that control what happens when a condition is signaled

```
result = tryCatch(code,  
  error = function(c) "error",  
  warning = function(c) "warning",  
  message = function(c) "message")  
)
```

Use **conditionMessage(c)** or **c\$message** to extract the message associated with the original error.

Defensive Programming

Basic principle : "fail fast", to raise an error as soon as something goes wrong

- stopifnot()** or use 'assertthat' package - check inputs are correct
- Avoid subset(), transform() and with()** - these are non-standard evaluation, when they fail, often fail with uninformative error messages.
- Avoid lapply() and supply()** - functions that can return different types of output.
 - Recommendation : Whenever subsetting a data frame in a function, you should always use **drop = FALSE**

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 **tidy**.

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
- **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 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())
```

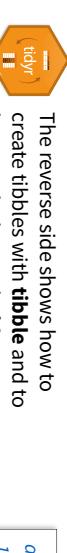
### Read Apache style log files

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

## Read Tabular Data

- These functions share the common arguments:

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



## Comma Delimited Files

**read\_csv(file.csv")**

To make file.csv run:

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

## Semi-colon Delimited Files

**read\_csv2("file2.csv")**

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

## 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|b|c\n1|2|3\n4|5|NA", path = "file.tsv")
```

## USEFUL ARGUMENTS

| Example file          |        |         |
|-----------------------|--------|---------|
| a,b,c                 | 1,2,3  | 4,5,NA  |
| A B C                 | 1 2 3  | 4 5 NA  |
| <b>No header</b>      |        |         |
| A B C                 | 1 2 3  | 4 5 NA  |
| <b>Provide header</b> |        |         |
| A B C                 | 1 2 3  | 4 5 NA  |
| <b>Missing Values</b> |        |         |
| A B C                 | NA 2 3 | NA 5 NA |

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()**

## 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:
col(
age = col_integer(),
sex = col_character(),
earn = col_double()
)
```

**age** is an **integer**  
**sex** is a **character**  
**earn** is a **double (numeric)**



## 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,

- 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: 734 × 6   |
| manufacture <dbl>     |
| model <dbl>           |
| displ <dbl>           |
| hp <dbl>              |
| drat <dbl>            |
| wt <dbl>              |
| qsec <dbl>            |
| vs <dbl>              |
| am <dbl>              |
| cyl <dbl>             |
| gear <dbl>            |
| carb <dbl>            |
| mpg <dbl>             |
| ... # omitted 68 rows |

|                       |
|-----------------------|
| # A tibble: 734 × 6   |
| manufacture <dbl>     |
| model <dbl>           |
| displ <dbl>           |
| hp <dbl>              |
| drat <dbl>            |
| wt <dbl>              |
| qsec <dbl>            |
| vs <dbl>              |
| am <dbl>              |
| cyl <dbl>             |
| gear <dbl>            |
| carb <dbl>            |
| mpg <dbl>             |
| ... # omitted 68 rows |

Each variable is in its own **column**, gathering the column values into a single value column.

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

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

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

gather(table4a, `1999`, `2000`)

key = "year", value = "cases")

spread(table2, type, count)

separate(table3, rate,

into = c("cases", "pop"))

.....

table3

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

.....

table4

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

.....

table5

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

.....

table6

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

.....

table7

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

.....

table8

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

.....

table9

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

.....

table10

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

.....

table11

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

.....

table12

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

.....

table13

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

.....

table14

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

.....

table15

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

.....

table16

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

.....

table17

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

.....

table18

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

.....

table19

|         |         |      |
|---------|---------|------|
| country | century | year |
| Afghan  | 19      | 99   |
| Afghan  |         |      |

# Data Wrangling

with `dplyr` and `tidyverse`

Cheat Sheet



## Syntax - Helpful conventions for wrangling

### `dplyr::tbl_df(iris)`

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

| Source: local data frame [150 x 5] |              |             |              |             |
|------------------------------------|--------------|-------------|--------------|-------------|
|                                    | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width |
| 1                                  | 5.1          | 3.5         | 1.4          | 0.2         |
| 2                                  | 4.9          | 3.0         | 1.4          | 0.2         |
| 3                                  | 4.7          | 3.2         | 1.3          | 0.2         |
| 4                                  | 4.6          | 3.1         | 1.5          | 0.2         |
| 5                                  | 5.0          | 3.6         | 1.4          | 0.2         |
| ..                                 | ..           | ..          | ..           | ..          |
| 98                                 | 5.1          | 3.3         | 1.7          | 0.5         |
| 99                                 | 5.1          | 3.7         | 1.5          | 0.4         |
| 100                                | 5.0          | 3.0         | 1.3          | 0.4         |

Variables not shown: Petal.Width (dbl), Species (fctr)

### `dplyr::glimpse(iris)`

Information dense summary of `tbl` data.

### `utils::View(iris)`

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

### `dplyr::tbl`

Extract rows that meet logical criteria.

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

`dplyr::distinct(iris)`  
Remove duplicate rows.

`dplyr::sample_frac(iris, 0.5, replace = TRUE)`  
Randomly select fraction of rows.

`dplyr::sample_n(iris, 10, replace = TRUE)`  
Randomly select n rows.

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

`dplyr::top_n(storms, 2, date)`  
Select and order top n entries (by group if grouped data).

`dplyr::%>%`  
Passes object on left hand side as first argument (or argument) of function on righthand side.

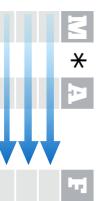
`x %>% f(y)` is the same as `f(x, y)`  
`y %>% f(x, , z)` is the same as `f(x, y, z)`

"Piping" with `%>%` makes code more readable, e.g.  
`iris %>% group_by(Species) %>% summarise(avg = mean(Sepal.Width)) %>% arrange(avg)`

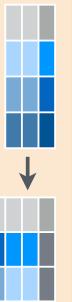
## Tidy Data - A foundation for wrangling in R



Tidy data complements R's **vectorized operations**. R will automatically preserve observations as you manipulate variables.



## Reshaping Data - Change the layout of a data set



`dplyr::data_frame(a = 1:3, b = 4:6)`  
Combine vectors into data frame (optimized).

`dplyr::arrange(mtcars, mpg)`  
Order rows by values of a column (low to high).

`dplyr::arrange(mtcars, desc(mpg))`  
Order rows by values of a column (high to low).

`dplyr::rename(tb, y = year)`  
Rename the columns of a data frame.

`dplyr::gather(cases, "year", "n", 2:4)`  
Gather columns into rows.

`dplyr::spread(pollution, size, amount)`  
Spread rows into columns.

`dplyr::select(iris, Sepal.Length, Petal.Length, Species)`  
Select columns by name or helper function.

### Helper functions for `select` - ?`select`

`select(iris, contains("x"))`  
Select columns whose name contains a character string.

`select(iris, ends_with("Length"))`  
Select columns whose name ends with a character string.

`select(iris, everything())`  
Select every column.

`select(iris, matches("x:t"))`  
Select columns whose name matches a regular expression.

`select(iris, num_range("x", 1:5))`  
Select columns named x1,x2,x3,x4,x5.

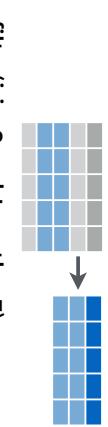
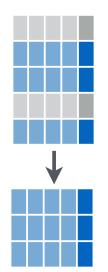
`select(iris, one_of(c("Species", "Genus")))`  
Select columns whose names are in a group of names.

`select(iris, starts_with("Sepal"))`  
Select columns whose name starts with a character string.

`select(iris, Sepal.Length:Petal.Width)`  
Select all columns between Sepal.Length and Petal.Width (inclusive).

`select(iris, -Species)`  
Select all columns except Species.

## Subset Variables (Columns)



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

`dplyr::select(iris, Sepal.Length, Petal.Length, Species)`  
Select columns by name or helper function.

`dplyr::filter(iris, contains("x"))`  
Select columns whose name contains a character string.

`dplyr::filter(iris, ends_with("Length"))`  
Select columns whose name ends with a character string.

`dplyr::filter(iris, everything())`  
Select every column.

`dplyr::filter(iris, matches("x:t"))`  
Select columns whose name matches a regular expression.

`dplyr::filter(iris, num_range("x", 1:5))`  
Select columns named x1,x2,x3,x4,x5.

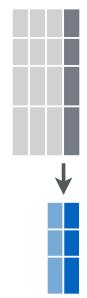
`dplyr::filter(iris, one_of(c("Species", "Genus")))`  
Select columns whose names are in a group of names.

`dplyr::filter(iris, starts_with("Sepal"))`  
Select columns whose name starts with a character string.

`dplyr::filter(iris, Sepal.Length:Petal.Width)`  
Select all columns between Sepal.Length and Petal.Width (inclusive).

`dplyr::filter(iris, -Species)`  
Select all columns except Species.

## Summarise Data



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

Summarise data into single row of values.

**dplyr::summarise\_each(iris, funs(mean))**

Apply summary function to each column.

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

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



summary  
function

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

**min**

Minimum value in a vector.

**max**

Maximum value in a vector.

**mean**

Mean value of a vector.

**median**

Median value of a vector.

**var**

Variance of a vector.

**sd**

Standard deviation of a vector.

**IQR**

IQR of a vector.

## Group Data

**dplyr::group\_by(iris, Species)**

Group data into rows with the same value of Species.

**dplyr::ungroup(iris)**

Remove grouping information from data frame.

**iris %>% group\_by(Species) %>% summarise(...)**

Compute separate summary row for each group.



window  
function

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

**dplyr::lead**

Copy with values shifted by 1.

**dplyr::lag**

Copy with values lagged by 1.

**dplyr::dense\_rank**

Ranks with no gaps.

**dplyr::min\_rank**

Ranks. Ties get min rank.

**dplyr::percent\_rank**

Ranks rescaled to [0, 1].

**dplyr::row\_number**

Ranks. Ties get to first value.

**dplyr::ntile**

Bin vector into n buckets.

**dplyr::between**

Are values between a and b?

**dplyr::cume\_dist**

Cumulative distribution.

**Binding**

**iris %>% group\_by(Species) %>% mutate(...)**

Compute new variables by group.



join data. Retain all values, all rows.  
**dplyr::inner\_join(a, b, by = "x1")**

Join data. Retain only rows in both sets.  
**dplyr::full\_join(a, b, by = "x1")**

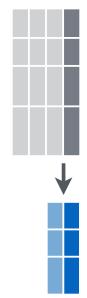
Join matching rows from a to b.  
**dplyr::left\_join(a, b, by = "x1")**

Join matching rows from b to a.  
**dplyr::right\_join(a, b, by = "x1")**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>a</b> |
| A         | 1         | T        |
| B         | 2         | F        |
| C         | 3         | NA       |
| D         | T         |          |

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x3</b> | <b>b</b> |
| A         | 1         | T        |
| B         | 2         | F        |
| C         | 3         | NA       |
| D         | T         |          |



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

Compute and append one or more new columns.

**dplyr::mutate\_each(iris, funs(min\_rank))**

Apply window function to each column.

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

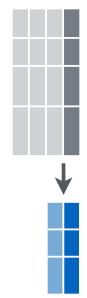
Compute one or more new columns. Drop original columns.

## Combine Data Sets

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>a</b> |
| A         | 1         | T        |
| B         | 2         | F        |
| C         | 3         | NA       |
| D         | T         |          |

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x3</b> | <b>b</b> |
| A         | 1         | T        |
| B         | 2         | F        |
| C         | 3         | NA       |
| D         | T         |          |



**dplyr::semi\_join(a, b, by = "x1")**

All rows in a that have a match in b.

**dplyr::anti\_join(a, b, by = "x1")**

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

Rows that appear in both y and z.  
**dplyr::intersect(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>z</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Rows that appear in either or both y and z.  
**dplyr::union(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>z</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Rows that appear in y but not z.  
**dplyr::setdiff(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Append z to y as new columns.  
**dplyr::bind\_cols(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>z</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         |          |

Caution: matches rows by position.  
**dplyr::bind\_rows(y, z)**

|           |           |          |
|-----------|-----------|----------|
| <b>x1</b> | <b>x2</b> | <b>y</b> |
| A         | 1         |          |
| B         | 2         |          |
| C         | 3         |          |
| D         | 4         | </       |

# Data Transformation With dplyr :: CHEAT SHEET



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

## Manipulate Cases



&



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

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

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



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

### Summarise Cases

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

### filter

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

Row functions return a subset of rows as a new table. Use a variant that ends in \_ for non-standard evaluation friendly code.

Column functions return a set of columns as a new table. Use a variant that ends in \_ for non-standard evaluation friendly code.

**distinct**(data, ..., keep\_all = FALSE) Remove rows with duplicate values. Also **distinct\_**().

**sample\_frac**(tbl, size = 1, replace = FALSE, weight = NULL, env = parent.frame()) Randomly select size rows. `sample_n(iris, 10, replace = TRUE)`

**slice**(data, ...) Select rows by position. Also `slice(1:5, iris[1:15])`

**summarise**(data, ...) Compute table of summaries. Also **summarise\_**().

`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 created 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))`

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

**Logical and boolean operators to use with filter()**

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

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

## ARRANGE CASES

**arrange**(data, ...) Order rows by values of a column (low to high), use with **desc**() to order from high to low.

`arrange(mtcars, mpg)`

`arrange(mtcars, desc(mpg))`

**group\_by**(data, ..., add = FALSE)

Returns copy of table grouped by ...  
`iris <- group_by(iris, Species)`

**ungroup**(x, ...)

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

## ADD CASES

**add\_row**(data, ..., before = NULL, after = NULL)

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

**mutate**(tbl, .predicte, .funs, ...)

Apply funs to all columns of one type. Use with **funs**().  
`mutate_if(iris, is.numeric, funs(log(.)))`

**add\_column**(data, ..., before = NULL, .after = NULL)

Add new column(s).  
`add_column(mtcars, new = 1:32)`

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

Column functions return a set of columns as a new table. Use a variant that ends in \_ for non-standard evaluation friendly code.

**select**(data, ...) Extract columns by name. Also **select\_if**()

`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)

**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**(). `mutate_all(faithful, funs(log(.), log2(.)))`

**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(.)))`

**mutate\_if**(tbl, .predicte, .funs, ...)

Apply funs to all columns of one type. Use with **funs**().  
`mutate_if(iris, is.numeric, funs(log(.)))`

# Vectorized Functions

## Summary 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()  
dplyr::cummax() - Cumulative max()  
dplyr::cummean() - Cumulative mean()  
cummin() - Cumulative min()  
cumprod() - Cumulative prod()  
cumsum() - Cumulative sum()

### RANKINGS

dplyr::cume\_dist() - Proportion of all values <= gaps()  
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"

### MISC

+,-,\*/,^,%/%,%% - arithmetic ops  
<=,>,>=,!=,== - logical comparisons

### SPREAD

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

dplyr::between() - x >= left & x <= right  
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

## ROW NAMES

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

A B C  
1 a t  
2 b u  
3 c v  
**rownames\_to\_column()**  
Move row names into col.  
d <- rownames\_to\_column(iris, var = "C")  
= "C")

A B C  
1 a t  
2 b u  
3 c v  
**column\_to\_rownames()**  
Move col in row names.  
column\_to\_rownames(a, var = "C")

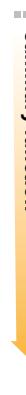
Also has **rownames()**, **remove\_rownames()**

## Combine Tables

### 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-NAs

### 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() - mean absolute deviation  
sd() - standard deviation  
var() - variance

### MISC

dplyr::between() - x >= left & x <= right  
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

### COMBINE CASES

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

### COMBINE VARIABLES

Use **bind\_cols(...)** to paste tables beside each other as they are.

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

### bind\_rows(..., id = NULL)

Use **bind\_rows()** to paste tables one on top of the other. BE SURE THAT ROWS ALIGN.

### bind\_cols(...)

Use **bind\_cols()** to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

### left\_join(x, y, by = NULL, copy = FALSE, suffix=c("x", "y"), ...)

Join matching values from x to y.

### inner\_join(x, y, by = NULL, copy = FALSE, suffix=c("x", "y"), ...)

Join data. Retain only rows with matches.

### full\_join(x, y, by = NULL, copy = FALSE, suffix=c("x", "y"), ...)

Join data. Retain all values, all rows.

### EXTRACT ROWS

### setdiff(x, y, ...)

Use **setdiff()** to specify the column(s) to match on.

Use **by = c("col1", "col2")** to specify the column(s) to match on.

Use **by = c("col1" = "col2")** to specify the column(s) to match on.

Use **by = c("col1" = "col2")** to specify the column(s) to match on.

### semi\_join(x, y, by = NULL, ...)

Use **semi\_join()** to filter one table against the rows of another.



# 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: **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 = <code><DATA></code>) +
 <code><GEOM_FUNCTION>(<POSITION> = aes(<MAPPINGS>,
 stat = <STAT>), position = <POSITION> + <COORDINATE_FUNCTION> + <SCALE_FUNCTION> + <THEME_FUNCTION>)</code>

```

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

`qplot`(`x` = cyl, `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 `ggplot`("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(mpg, aes(cty, hwy))

```

```

a + geom_blank()

```

(Used for expanding limits)

```

b + geom_curve(aes(vend = lat + 1,
 xend = long + 1, curvature = z))

```

(Used for expanding limits)

```

a + geom_path(lineend = "butt", linejoin = "round",
 lineirel = 1)

```

(Used for expanding limits)

```

a + geom_rect(aes(xmin = long, ymin = lat,
 xmax = long + 1, ymax = (lat + 1) - xmax, xmin, ymax,
 ymin, alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

b + geom_segment(aes(xend = lat + 1, xend = long + 1)

```

(Used for expanding limits)

```

b + geom_spoke(aes(angle = 1:115, radius = 1))

```

(Used for expanding limits)

### 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_smooth(method = lm, x, y, alpha,
 alpha, color, fill, group, linetype, size)

```

```

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

```

(Used for expanding limits)

```

c <- ggplot(mpg, aes(hwy))

```

(Used for expanding limits)

```

c + geom_area(stat = "bin")

```

(Used for expanding limits)

```

c + geom_boxplot(x, y, alpha, color, fill, group,
 ymin, ymax, alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

f + geom_dotplot(x, y, alpha, color, fill, group,
 x, y, alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

f + geom_freqpoly(x, y, alpha, color, group,
 linetype, size)

```

(Used for expanding limits)

```

c + geom_histogram(binwidth = 5) x, y, alpha,
 color, fill, linetype, size, weight

```

(Used for expanding limits)

**ONE VARIABLE** continuous

```

c <- ggplot(mpg, aes(hwy))

```

(Used for expanding limits)

```

c + geom_density(kernel = "gaussian")

```

(Used for expanding limits)

```

c + geom_hex(x, y, alpha, color, fill, group,
 alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

c + geom_point(x, y, alpha, color, fill, group,
 alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

c + geom_pointrange(x, y, alpha, color, fill, group,
 alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

j + geom_raster(aes(fill = z),
 state = tolower(trownames(USArrests)))

```

(Used for expanding limits)

**TWO VARIABLES**

continuous `x`, continuous `y`

```

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

```

(Used for expanding limits)

```

e + geom_label(aes(label = cyl), nudge_x = 1,
 alpha, angle, color, family, fontface, hjust,
 lineheight, size, vjust)

```

(Used for expanding limits)

```

e + geom_jitter(height = 2, width = 2)

```

(Used for expanding limits)

```

e + geom_hex(x, y, alpha, color, fill, shape,
 size, stroke)

```

(Used for expanding limits)

(Used for expanding limits)

(Used for expanding limits)

### CONTINUOUS BIVARIATE DISTRIBUTION

```

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

```

```

h + geom_bin2d(binwidth = c(0.25, 500))

```

(Used for expanding limits)

```

h + geom_hex2d(binwidth = c(0.25, 500))

```

(Used for expanding limits)

```

h + geom_hex(x, y, alpha, colour, group, linetype, size)

```

(Used for expanding limits)

```

h + geom_rug(sides = "bl")

```

(Used for expanding limits)

```

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

```

(Used for expanding limits)

```

h + geom_text(label = cyl, nudge_x = 1,
 nudge_y = 1, check_overlap = TRUE, x, y, label,
 alpha, angle, color, family, fontface, hjust,
 lineheight, size, vjust)

```

(Used for expanding limits)

### CONTINUOUS FUNCTION

```

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

```

```

i + geom_area()

```

(Used for expanding limits)

```

i + geom_line()

```

(Used for expanding limits)

```

i + geom_step(direction = "hv")

```

(Used for expanding limits)

### VISUALIZING ERROR

```

df <- data.frame(gp = c("A", "B"), fit = 4.5, se = 1.2)

```

```

j <- ggplot(df, aes(gp, fit, ymin = fit - se, ymax = fit + se))

```

(Used for expanding limits)

```

j + geom_crossbar(fatten = 2)

```

(Used for expanding limits)

```

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

```

(Used for expanding limits)

```

j + geom_pointrange(x, y, alpha, color, fill, group, linetype, size)

```

(Used for expanding limits)

```

j + geom_errorbar(x, y, alpha, color, fill, group, linetype, size, width)

```

(Used for expanding limits)



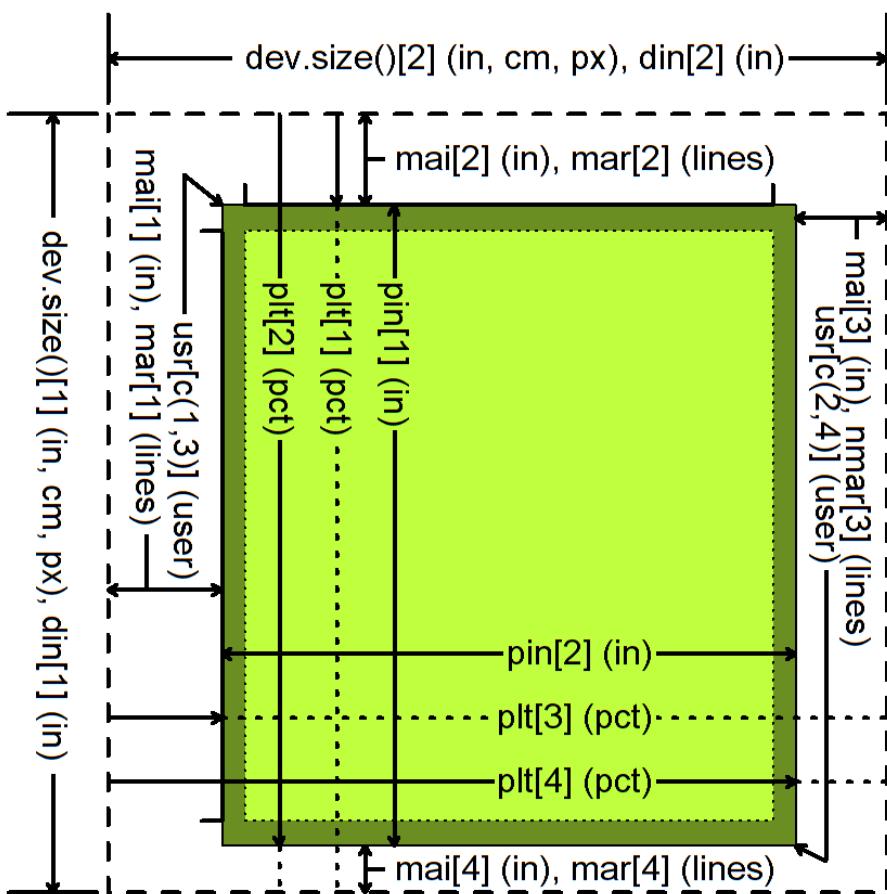


# How Big is Your Graph?

## 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 plotting region

**par("pin")** (width, height) in inches  
**par("plt")** (left, right, bottom, top) in pct

The **pin** argument gives you the size of the plotting region (the size of the device minus the size of the margins) in inches.

The **plt** argument gives you the percentage of the device from the left/bottom edge up to the left edge of the plotting region, the right edge, the bottom edge, and the top edge. The first and third values are equivalent to the percentage of space devoted to the left and bottom margins. Subtract the second and fourth values from 1 to get the percentage of space devoted to the right and top margins.

### Your x-y coordinates

**par("usr")** (xmin, ymin, xmax, ymax)

Your x-y coordinates are the values you use when plotting your data. This normally is not the same as the values you specified with the **xlim** and **ylim** arguments in **plot**. By default, R adds an extra 4% to the plotting range (see the dark green region on the figure) so that points right up on the edges of your plot do not get partially clipped. You can override this by setting **xaxs="in"** and/or the **yaxs="in"** parameter.

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")**  
**par("mar")** (bottom, left, top, right) in inches

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))**

### 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 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 <-
px.per.xy <- region.px / user.range
px.dev.size(units="px") * region.pct
```

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

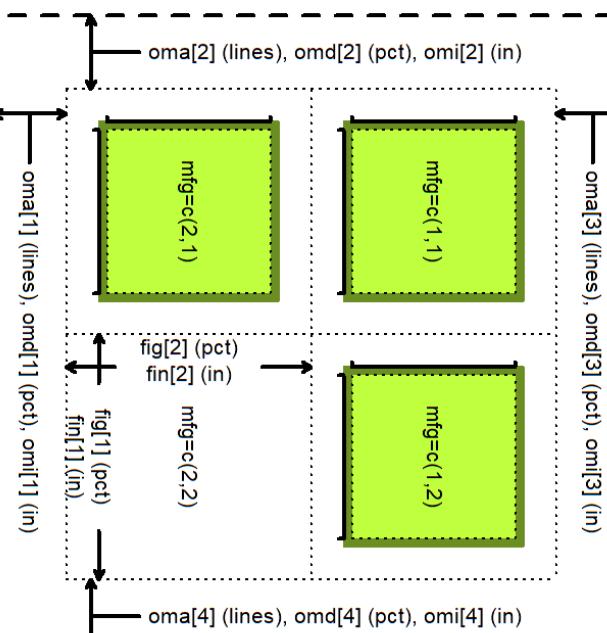
If you display multiple plots within a single graphics window (e.g., with the `mflow` or `mfcol` arguments of `par` or with the `layout` function), then the `fig` and `fin` arguments will tell you the size of the current subplot window in percent or inches, respectively.

`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.



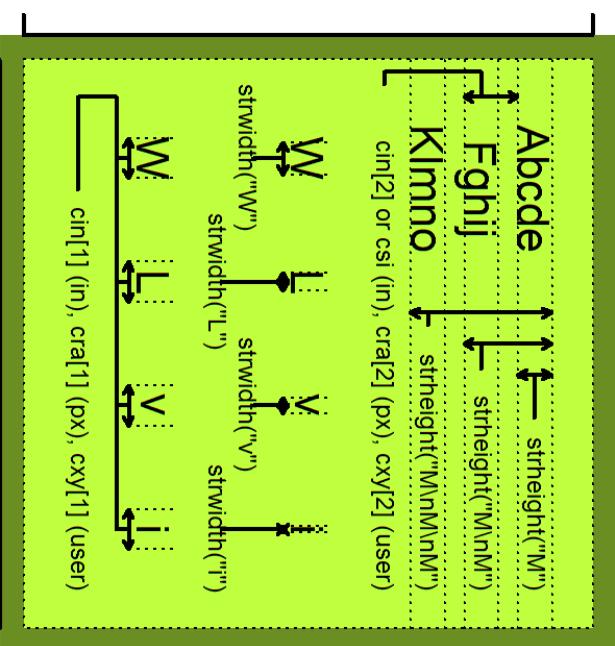
## stheight()

The `stheight` 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, `stheight` will give you the height of the letter "M". If you have a string with one or more linebreaks ("\\n"), the `stheight` 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 `height` argument of `par`. The default line height (`height=1`), corresponding to single spaced lines, produces a line height roughly 1.5 times the height of "M".

## strwidth()

The `strwidth` function will produce different widths to individual characters, representing the proportional spacing used by most fonts (a "W" using much more space than an "i"). For the width of a string, the `strwidth` function will sum up the lengths of the individual characters in the string.



## Character and string sizes

`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 `cra`, `cxy` gives you the height of a line, in inches, pixels, or xy (user) coordinates.

The first of the two values returned by the `cra`, `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. Trying opening your graphics device with different values for `height` and `width`. Fonts that look too big might be better
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 `mai` 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.

`las`

`las = 0`      `las = 1`

`las`

`las = 0`      `las = 1`

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.

`cex.main`

`cex.lab`

`cex.axis`

`cex.lab`

`cex.axis`

`cex.lab`

`cex.axis`

`cex.lab`

`cex.axis`

# Work With Strings 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. | TRUE → TRUE<br>FALSE → FALSE<br>TRUE → TRUE |
| <code>str_detect(fruit, "a")</code>      |                                                     |                                             |
| 0                                        | 0 → 1 → 2 → 4                                       |                                             |
| 1                                        |                                                     |                                             |
| 2                                        |                                                     |                                             |
| 3                                        |                                                     |                                             |
| 4                                        |                                                     |                                             |
| 5 → state end                            |                                                     |                                             |
| 6 → NA                                   |                                                     |                                             |
| 7 → NA NA                                |                                                     |                                             |
| 8 → NA NA                                |                                                     |                                             |
| 9 → NA NA                                |                                                     |                                             |
| 10 → NA NA                               |                                                     |                                             |
| 11 → NA NA                               |                                                     |                                             |
| 12 → NA NA                               |                                                     |                                             |
| 13 → NA NA                               |                                                     |                                             |
| 14 → NA NA                               |                                                     |                                             |
| 15 → NA NA                               |                                                     |                                             |
| 16 → NA NA                               |                                                     |                                             |
| 17 → NA NA                               |                                                     |                                             |

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

`str_count(fruit, "a")`

0

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

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

## Subset Strings

|                                                       |                                             |                |
|-------------------------------------------------------|---------------------------------------------|----------------|
| <code>str_sub(string, start = 1L, end = -1L)</code>   | Extract substrings from a character vector. | 4 → 6<br>2 → 3 |
| <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. | 4 → 6<br>2 → 3 |
| <code>str_subset(fruit, "a")</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> | 4 → 6<br>2 → 3 |
| <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> .                                | 4 → 6<br>2 → 3 |

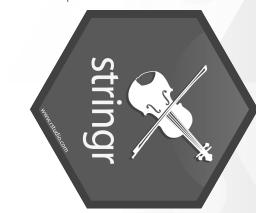
|                                                      |  |  |
|------------------------------------------------------|--|--|
| <code>str_match(sentences, "(a the) (/n  +)")</code> |  |  |
|                                                      |  |  |

## 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> | 4 → 6<br>2 → 3 |
| <code>str_sub(fruit, 1, 3), str_sub(fruit, -2)</code> |                                                                                                                                    |                |

|                                                                                              |                                                                                                  |                |
|----------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------|----------------|
| <code>str_pad(string, width, side = c("left", "right", "both"))</code>                       | Pad strings to constant width. <code>str_pad(fruit, 17)</code>                                   | 4 → 6<br>2 → 3 |
| <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> | 4 → 6<br>2 → 3 |

|                                                                  |                                                                                     |                |
|------------------------------------------------------------------|-------------------------------------------------------------------------------------|----------------|
| <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> | 4 → 6<br>2 → 3 |
| <code>str_trim(fruit)</code>                                     |                                                                                     |                |



## Mutate Strings

|                                                              |                                                                                             |                |
|--------------------------------------------------------------|---------------------------------------------------------------------------------------------|----------------|
| <code>str_sub(string, start = 1L, end = -1L)</code>          | Identifying the substrings with <code>str_sub()</code> and assigning into the results.      | 4 → 6<br>2 → 3 |
| <code>str_sub(fruit, 1, 3) &lt;- "Str"</code>                |                                                                                             |                |
| <code>str_replace(string, pattern, replacement)</code>       | Replace the first matched pattern in each string. <code>str_replace(fruit, "a", "-")</code> | 4 → 6<br>2 → 3 |
| <code>str_replace_all(string, pattern, replacement)</code>   | Replace all matched patterns in each string. <code>str_replace_all(fruit, "a", "-")</code>  | 4 → 6<br>2 → 3 |
| <code>str_to_lower(string, locale = "en")<sup>1</sup></code> | Convert strings to lower case. <code>str_to_lower(sentences)</code>                         | 4 → 6<br>2 → 3 |
| <code>str_to_upper(string, locale = "en")<sup>1</sup></code> | Convert strings to upper case. <code>str_to_upper(sentences)</code>                         | 4 → 6<br>2 → 3 |
| <code>str_to_title(string, locale = "en")<sup>1</sup></code> | Convert strings to title case. <code>str_to_title(sentences)</code>                         | 4 → 6<br>2 → 3 |

|                                                    |                                                    |                |
|----------------------------------------------------|----------------------------------------------------|----------------|
| <code>str_c(..., sep = "", collapse = NULL)</code> | Join multiple strings into a single string.        | 4 → 6<br>2 → 3 |
| <code>str_c(letters, LETTERS)</code>               |                                                    |                |
| <code>str_c(..., sep = "", collapse = NULL)</code> | Collapse a vector of strings into a single string. | 4 → 6<br>2 → 3 |
| <code>str_c(letters, collapse = "")</code>         |                                                    |                |
| <code>str_dup(string, times)</code>                | Repeat strings times                               | 4 → 6<br>2 → 3 |
| <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). Also <code>str_split</code> to return a list of substrings.                 | 4 → 6<br>2 → 3 |
| <code>str_split_fixed(fruit, "", n=2)</code>                                                                  |                                                                                                                                                                            |                |
| <code>glue:glue(..., sep = "", .envir = parent.frame(), .open = "'", .close = "'")<sup>1</sup></code>         | Create a string from strings and {expressions} to evaluate. <code>glue:glue("{'pi is {pi}")</code>                                                                         | 4 → 6<br>2 → 3 |
| <code>glue:glue_data(x, ..., sep = "", .envir = parent.frame(), .open = "'", .close = "'")<sup>1</sup></code> | Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate. <code>glue:glue_data(mtcars, "rownames(mtcars)) has {hp} hp")</code> | 4 → 6<br>2 → 3 |
| <code>str_view_all(string, pattern, match = NA)</code>                                                        | View HTML rendering of all regex matches.                                                                                                                                  | 4 → 6<br>2 → 3 |
| <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>                                                                                        | 4 → 6<br>2 → 3 |

## Join and Split

### Order Strings

### Helpers

`str_order(x, decreasing = FALSE, na.last = TRUE, locale = "en", numeric = FALSE, ...)` Return the vector of indexes that sorts a character vector. `x$str_order(x)`

`str_sort(x, decreasing = FALSE, na.last = TRUE, locale = "en", numeric = FALSE, ...)` Sort a character vector. `x$str_order(x)`

`str_sort(x, decreasing = FALSE, na.last = TRUE, locale = "en", numeric = FALSE, ...)` Sort a character vector. `x$str_order(x)`

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

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

`str_view_all(string, pattern, match = NA)` View HTML rendering of all regex matches.

`str_view_all(fruit, "[aeiou]")`

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

# Need to Know

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

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

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

| Special Character | Represents                                |
|-------------------|-------------------------------------------|
| \                 | a backslash                               |
| \n                | a new line                                |
| \t                | a tab                                     |
| \r                | a carriage return                         |
| \f                | a form feed                               |
| \v                | a vertical tab                            |
| \b                | a word boundary                           |
| \w                | any word character (W for non-word chars) |
| \s                | any whitespace (\S for non-whitespaces)   |
| \d                | any digit (\D for non-digits)             |
| \w                | any word character (W for non-word chars) |
| \b                | word boundaries                           |
| \d                | digits                                    |
| \w                | letters                                   |
| \b                | lowercase letters                         |
| \B                | uppercase letters                         |
| \w                | letters and numbers                       |
| \punct            | punctuation                               |
| \graph            | letters, numbers, and punctuation         |
| \space            | space characters (i.e. \s)                |
| \blank            | space and tab (but not new line)          |
| .                 | every character except a new line         |

## INTERPRETATION

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

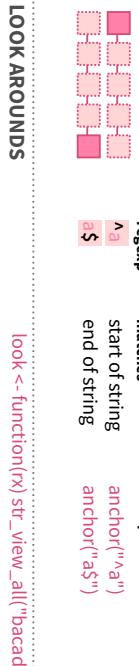
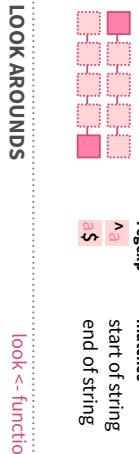
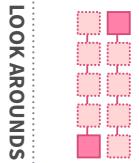
```
rexpx(pattern, ignore_case = FALSE, multiline = FALSE, comments = FALSE, dotall = FALSE,...)
Modifies a regex to ignore cases, match end of lines as well or end of strings, allow R comments within regex's, and/or have . match everything including \n.
str_detect("/", regex("/", TRUE))
```

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

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

**boundary()** Matches boundaries between characters, line breaks, sentences, or words.

str\_split(sentences, boundary("word"))



# Regular Expressions -

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

## MATCH CHARACTERS

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

| string (type <code>regexp</code> this) | matches                             | example                         |
|----------------------------------------|-------------------------------------|---------------------------------|
| <code>a (etc.)</code>                  | <code>matches</code> (to mean this) | <code>a (etc.)</code>           |
| <code>\.</code>                        | <code>(which matches this)</code>   | <code>see("a")</code>           |
| <code>!</code>                         |                                     | <code>see("!\")</code>          |
| <code>?</code>                         |                                     | <code>see("!\?")</code>         |
| <code>\?</code>                        |                                     | <code>see("!\?")</code>         |
| <code>\{</code>                        |                                     | <code>see("!\{\")</code>        |
| <code>\}</code>                        |                                     | <code>see("!\}")</code>         |
| <code>\(</code>                        |                                     | <code>see("!\()")</code>        |
| <code>\)</code>                        |                                     | <code>see("!\)")</code>         |
| <code>\{</code>                        |                                     | <code>abc ABC 123 .\r\n0</code> |
| <code>\}</code>                        |                                     | <code>abc ABC 123 .\r\n0</code> |
| <code>\(</code>                        |                                     | <code>abc ABC 123 .\r\n0</code> |
| <code>\)</code>                        |                                     | <code>abc ABC 123 .\r\n0</code> |
| <code>\[</code>                        |                                     | <code>abc ABC 123 .\r\n0</code> |
| <code>\]</code>                        |                                     | <code>abc ABC 123 .\r\n0</code> |
| <code>\n</code>                        | <code>new line</code>               | <code>\n</code>                 |
| <code>\t</code>                        | <code>tab</code>                    | <code>\t</code>                 |
| <code>\r</code>                        | <code>space</code>                  | <code>\r</code>                 |

see `<- function(rx) str_view_all("0 1 2 3 4 5 6 7 8 9", rx)`

| string (type <code>regexp</code> this) | matches                                  | example                          |
|----------------------------------------|------------------------------------------|----------------------------------|
| <code>[0-9]</code>                     | <code>matches</code> digits              | <code>0 1 2 3 4 5 6 7 8 9</code> |
| <code>[a-f]</code>                     | <code>matches</code> letters             | <code>[a-f]</code>               |
| <code>[A-F]</code>                     | <code>matches</code> uppercase letters   | <code>A B C D E F</code>         |
| <code>[0-9A-F]</code>                  | <code>matches</code> letters and numbers | <code>g h i j k l</code>         |
| <code>[[:alpha:]]</code>               | <code>matches</code> punctuation         | <code>G H I J K L</code>         |
| <code>[[:lower:]]</code>               | <code>matches</code> punctuation         | <code>m n o p q r</code>         |
| <code>[[:upper:]]</code>               | <code>matches</code> punctuation         | <code>M N O P Q R</code>         |
| <code>[[:alnum:]]</code>               | <code>matches</code> punctuation         | <code>s t u v w x</code>         |
| <code>[[:graph:]]</code>               | <code>matches</code> punctuation         | <code>S T U V W X</code>         |
| <code>[[:space:]]</code>               | <code>matches</code> punctuation         | <code>z</code>                   |

see `<- function(rx) str_view_all("a[bc]", rx)`

| string (type <code>regexp</code> this) | matches                                  | example                         |
|----------------------------------------|------------------------------------------|---------------------------------|
| <code>a[bc]</code>                     | <code>matches</code> a or b or c         | <code>a b c d e f</code>        |
| <code>[ab]</code>                      | <code>matches</code> either a or b       | <code>abc ABC 123 .\r\n0</code> |
| <code>[^ab]</code>                     | <code>matches</code> anything but a or b | <code>g h i j k l</code>        |
| <code>[a-c]</code>                     | <code>matches</code> range               | <code>M N O P Q R</code>        |

see `<- function(rx) str_view_all("aa[a]", rx)`

| string (type <code>regexp</code> this) | matches                                         | example                  |
|----------------------------------------|-------------------------------------------------|--------------------------|
| <code>aa[a]</code>                     | <code>matches</code> a or aa                    | <code>A B C D E F</code> |
| <code>aa[a-]</code>                    | <code>matches</code> a or aa or aaa             | <code>g h i j k l</code> |
| <code>aa[a-z]</code>                   | <code>matches</code> a or aa or aaa or ... or z | <code>M N O P Q R</code> |

see `<- function(rx) str_view_all("a{2,4}", rx)`

| string (type <code>regexp</code> this) | matches                                | example                      |
|----------------------------------------|----------------------------------------|------------------------------|
| <code>a{2,4}</code>                    | <code>matches</code> zero or one a     | <code>a a a a a a a a</code> |
| <code>a{2,4}</code>                    | <code>matches</code> zero or more a    | <code>a a a a a a a a</code> |
| <code>a{2,4}</code>                    | <code>matches</code> one or more a     | <code>a a a a a a a a</code> |
| <code>a{2,4}</code>                    | <code>matches</code> exactly n a       | <code>a a a a a a a a</code> |
| <code>a{2,4}</code>                    | <code>matches</code> n or more a       | <code>a a a a a a a a</code> |
| <code>a{2,4}</code>                    | <code>matches</code> between n and m a | <code>a a a a a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`

| string (type <code>regexp</code> this) | matches                                                  | example                        |
|----------------------------------------|----------------------------------------------------------|--------------------------------|
| <code>a{2,4}.</code>                   | <code>matches</code> zero or one a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> zero or more a followed by a dot    | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> one or more a followed by a dot     | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> exactly n a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> n or more a followed by a dot       | <code>a a a a . a a a a</code> |
| <code>a{2,4}.</code>                   | <code>matches</code> between n and m a followed by a dot | <code>a a a a . a a a a</code> |

see `<- function(rx) str_view_all("a{2,4}.aaa", rx)`



# 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 render 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

Populate the file with a template

### 2 Write document by editing template

### 3 Knit document to create report, use knit button or render() to knit

### 4 Preview Output in IDE window

### 5 Publish (optional) to web server

### 6 Examine build log in R Markdown console

- 7 Use output file that is saved along side .Rmd

## Embed code with knitr syntax

### INLINE CODE

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

Built with `r getRVersion()`

Built with 3.2.3

## render

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

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

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

```{r echo=TRUE}``

```{r}``

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- File Menu:** Open in Browser, Publish, Publish to shinylabps.io, RStudio Connect, Reload document.
- Code Editor:** Shows R Markdown code for "report.Rmd".
- Preview Area:** Shows the rendered HTML output of the R Markdown code.
- Console:** Shows the R Markdown build log.
- Plots:** Shows a histogram of cars data.
- Help:** Shows the R Markdown documentation page.

## .rmd Structure

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

At start of file

Between lines of - - -

**Text**  
Narration formatted with markdown, mixed with:

**Code Chunks**  
Chunks of embedded code. Each chunk:

Begins with ```{r}```

ends with ```{}```

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

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

2. Call parameters • Call parameter values in code as params\$`name`

3. Set parameters • Set values with knit with parameters or the params argument of render():

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

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE

3. Call Shiny render functions to embed reactive output.

4. Render with rmarkdown:::run or click Run Document in RStudio IDE



## Interactive Documents

Embed a complete app into your document with

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.

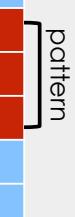


# Basic Regular Expressions in R

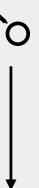
## Cheat Sheet

| Character Classes |                                                               |
|-------------------|---------------------------------------------------------------|
| [:digit:] or \d   | Digits; [0-9]                                                 |
| \D                | Non-digits; [^0-9]                                            |
| [:lower:]         | Lower-case letters; [a-z]                                     |
| [:upper:]         | Upper-case letters; [A-Z]                                     |
| [:alpha:]         | Alphabetic characters; [a-zA-Z]                               |
| [:alnum:]         | Alphanumeric characters; [a-zA-Z0-9]                          |
| \W                | Non-word characters                                           |
| [[blank]]         | Word characters; [A-Za-f]                                     |
| [[space]] or \s   | Space, tab, vertical tab, newline, form feed, carriage return |
| \S                | Not space; [^:space:]                                         |
| [[punct]]         | Punctuation characters; !#\$& ()*,.-/:;<>?@]{_,`{ ~           |
| [[graph]]         | Graphical char.; [:alnum:]                                    |
| [[print]]         | Printable characters; [:alnum:]                               |
| [[ctrl]] or \c    | Control characters; \n, \r etc.                               |

## Functions for Pattern Matching



### Detect pattern



| pattern                                                                                                                                                                                                                                          | Extract Patterns                                                                                                                                                                 |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <code>regmatches(string, grepexpr(pattern, string))</code>                                                                                                                                                                                       | extract first match                                                                                                                                                              |
| <code>regmatches(string, gregexpr(pattern, string))</code>                                                                                                                                                                                       | extracts all matches, outputs a list                                                                                                                                             |
| <code>[[1]] "tam" [[2]] character(0) [[3]] "tim" "tom"</code>                                                                                                                                                                                    | <code>stringr::str_extract(string, pattern)</code>                                                                                                                               |
| <code>[[1]] "tam" NA "tim"</code>                                                                                                                                                                                                                | <code>stringr::str_extract_all(string, pattern)</code>                                                                                                                           |
| <code>stringr::str_extract_all(string, pattern)</code>                                                                                                                                                                                           | extract all matches, outputs a list                                                                                                                                              |
| <code>stringr::str_match(string, pattern)</code>                                                                                                                                                                                                 | extract first match + individual character groups                                                                                                                                |
| <code>stringr::str_match_all(string, pattern)</code>                                                                                                                                                                                             | extract all matches + individual character groups                                                                                                                                |
| Locate Patterns                                                                                                                                                                                                                                  | Locate Patterns                                                                                                                                                                  |
| <code>gregexpr(pattern, string)</code>                                                                                                                                                                                                           | find starting position and length of first match                                                                                                                                 |
| <code>gregexpr(pattern, string)</code>                                                                                                                                                                                                           | find starting position and length of all matches                                                                                                                                 |
| <code>gregexpr(pattern, string)</code>                                                                                                                                                                                                           | find starting and end position of first match                                                                                                                                    |
| <code>gregexpr(pattern, string)</code>                                                                                                                                                                                                           | find starting and end position of all matches                                                                                                                                    |
| Replace Patterns                                                                                                                                                                                                                                 | Replace Patterns                                                                                                                                                                 |
| <code>gsub(pattern, replacement, string)</code>                                                                                                                                                                                                  | replace first match                                                                                                                                                              |
| <code>gsub(pattern, replacement, string)</code>                                                                                                                                                                                                  | replace all matches                                                                                                                                                              |
| <code>stringr::str_replace(string, pattern, replacement)</code>                                                                                                                                                                                  | replace first match                                                                                                                                                              |
| <code>stringr::str_replace_all(string, pattern, replacement)</code>                                                                                                                                                                              | replace all matches                                                                                                                                                              |
| Anchors                                                                                                                                                                                                                                          | Quantifiers                                                                                                                                                                      |
| <code>^</code>                                                                                                                                                                                                                                   | <code>*</code>                                                                                                                                                                   |
| Start of the string                                                                                                                                                                                                                              | Matches at least 0 times                                                                                                                                                         |
| <code>\$</code>                                                                                                                                                                                                                                  | <code>+</code>                                                                                                                                                                   |
| End of the string                                                                                                                                                                                                                                | Matches at least 1 time                                                                                                                                                          |
| <code>\b</code>                                                                                                                                                                                                                                  | <code>?</code>                                                                                                                                                                   |
| Empty string at either edge of a word                                                                                                                                                                                                            | Matches at most 1 time; optional string                                                                                                                                          |
| <code>\B</code>                                                                                                                                                                                                                                  | <code>{n}</code>                                                                                                                                                                 |
| NOT the edge of a word                                                                                                                                                                                                                           | Matches exactly n times                                                                                                                                                          |
| <code>\K</code>                                                                                                                                                                                                                                  | <code>{n,}</code>                                                                                                                                                                |
| Beginning of a word                                                                                                                                                                                                                              | Matches at least n times                                                                                                                                                         |
| <code>\D</code>                                                                                                                                                                                                                                  | <code>{,n}</code>                                                                                                                                                                |
| End of a word                                                                                                                                                                                                                                    | Matches at most n times                                                                                                                                                          |
| <code>\N</code> where N is an integer                                                                                                                                                                                                            | Matches between n and m times                                                                                                                                                    |
| CaseConversions                                                                                                                                                                                                                                  | Greedy Matching                                                                                                                                                                  |
| Regular expressions can be made case insensitive using <code>(?i)</code> . In backreferences, the strings can be converted to lower or upper case using <code>\L</code> or <code>\U</code> (e.g. <code>\L\1</code> ). This requires PERL = TRUE. | By default the asterisk <code>*</code> is greedy, i.e. it always matches the longest possible string. It can be used in lazy mode by adding <code>? i.e. <code>*?</code>.</code> |
| Greedy mode can be turned off using <code>(?U)</code> . This switches the syntax, so that <code>(?U)a*</code> is lazy and <code>(?U)a*?</code> is greedy.                                                                                        | Regular expressions can conveniently be created using <code>rex::rex()</code> .                                                                                                  |

| Special Metacharacters |                 |
|------------------------|-----------------|
| <code>\n</code>        | New line        |
| <code>\r</code>        | Carriage return |
| <code>\t</code>        | Tab             |
| <code>\v</code>        | Vertical tab    |
| <code>\f</code>        | Form feed       |

| Character Classes and Groups |                                                                   |
|------------------------------|-------------------------------------------------------------------|
| <code>-</code>               | Any character except \n                                           |
| <code> </code>               | Or, e.g. (a b)                                                    |
| <code>[...]</code>           | List permitted characters, e.g. [abc]                             |
| <code>[a-zA-Z]</code>        | Specify character ranges                                          |
| <code>[^...]</code>          | List excluded characters                                          |
| <code>(...)</code>           | Grouping, enables back referencing using \N where N is an integer |

## Lookarounds and Conditionals\*

| General Modes              |                                                                                            |
|----------------------------|--------------------------------------------------------------------------------------------|
| <code>(?=</code>           | Lookahead (requires PERL = TRUE), e.g. <code>(?=&gt;xy)</code> ; position followed by 'xy' |
| <code>(?!=</code>          | Negative lookahead (PERL = TRUE); position NOT followed by pattern                         |
| <code>(?&lt;=</code>       | Lookbehind (PERL = TRUE), e.g. <code>(?&lt;=yx)</code> ; position following 'xy'           |
| <code>(?&lt;!</code>       | Negative lookbehind (PERL = TRUE); position NOT following pattern                          |
| <code>?{if}then</code>     | If-then-condition (PERL = TRUE); use lookahead, optional char. etc in if-clause            |
| <code>?{if}thenelse</code> | If-then-else-condition (PERL = TRUE)                                                       |

\*see, e.g. <http://www.regular-expressions.info/lookaround.html>

# R Markdown Cheat Sheet

learn more at [rmarkdown.rstudio.com](http://rmarkdown.rstudio.com)

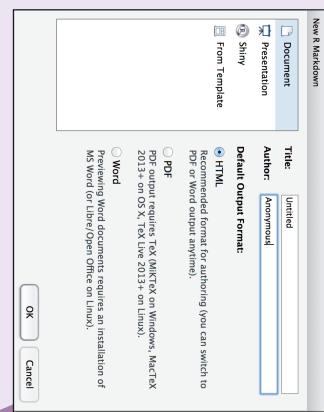
rmarkdown 0.2.50 Updated: 8/14



## 2. Open File

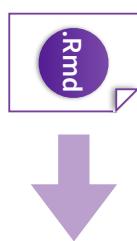
Start by saving a text file with the extension .Rmd, or open an RStudio Rmd template

- In the menu bar, click **File ▶ New File ▶ R Markdown...**
- A window will open. Select the class of output you would like to make with your .Rmd file
- Select the specific type of output to make with the radio buttons (you can change this later)
- Click OK

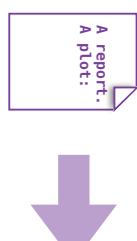


**1. Workflow** R Markdown is a format for writing reproducible, dynamic reports with R. Use it to embed R code and results into slideshows, pdfs, html documents, Word files and more. To make a report:

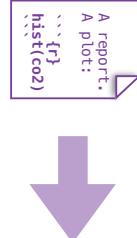
i. **Open** - Open a file that uses the .Rmd extension.



ii. **Write** - Write content with the easy to use R Markdown syntax



iii. **Embed** - Embed R code that creates output to include in the report



iv. **Render** - Replace R code with its output and transform the report into a slideshow, pdf, html or ms Word file.



## 3. Markdown

Next, write your report in plain text. Use markdown syntax to describe how to format text in the final report.

### Syntax

### becomes

Plain text  
End a line with two spaces to start a new paragraph.  
\*italics\* and \_italics\_  
\*\*bold\*\* and \_\_bold\_\_  
superscript<sup>A<sup>2</sup></sup>  
~~strikethrough~~  
[link](www.rstudio.com)

Plain text  
End a line with two spaces to start a new paragraph.  
**bold** and **bold**  
**superscript**  
~~strikethrough~~  
**link**

```
Header 1
Header 2
Header 3
Header 4
Header 5
Header 6
```

```
Header 1
Header 2
Header 3
Header 4
Header 5
Header 6
```

Header 1  
Header 2  
Header 3  
Header 4  
Header 5  
Header 6

## 4. Choose Output

Write a YAML header that explains what type of document to build from your R Markdown file.

### YAML

A YAML header is a set of key:value pairs at the start of your file. Begin and end the header with a line of three dashes (---)

The RStudio template writes the YAML header for you

```

title: "Untitled"
author: "Anonymous"
output: html_document

```

This is the start of my report. The above is metadata saved in a YAML header.

The output value determines which type of file R will build from your .Rmd file (in Step 6)

5

```
output: html_document html file (web page)
output: pdf_document pdf document
output: word_document Microsoft Word docx
output: beamer_presentation beamer slideshow (pdf)
output: ioslides_presentation ioslides slideshow (html)
```

**5. Embed Code** Use knitr syntax to embed R code into your report. R will run the code and include the results when you render your report.

### inline code

Surround code with back ticks and r. R replaces inline code with its results.

```
Two plus two
equals `r 2 + 2`.
```

Two plus two  
equals 4.

### code chunks

Start a chunk with ``{r}. End a chunk with ``}``.

```
Here's some code
```{r eval=FALSE}
dim(iris)
```

Here's some code

```
## [1] 150 5
```

display options

Use knitr options to style the output of a chunk.

Place options in brackets above the chunk.

display options

Use knitr options to style the output of a chunk.

Place options in brackets above the chunk.

6. Render Use your .Rmd file as a blueprint to build a finished report.

Render your report in one of two ways

1. Run `rmarkdown::render("<file path>")`
2. Click the **knit HTML** button at the top of the RStudio scripts pane

When you render R will

- execute each embedded code chunk and insert the results into your report
- build a new version of your report in the output file type
- open a preview of the output file in the viewer pane
- save the output file in your working directory

7. Interactive Docs

Turn your report into an interactive Shiny document in 3 steps

1 Add runtime:shiny to the YAML header

```
title: "line graph"
output: html_document
runtime: shiny
```

2 In the code chunks, add Shiny input functions to embed widgets. Add Shiny render functions to embed reactive output

```
title: "Line graph"
output: html_document
runtime: shiny

Choose a time series:
```{r echo = FALSE}
selectInput("data", "",
c("CO2", "LH"))
```
See a plot:
```{r echo = FALSE}
renderPlot({
d <- get(input$data)
plot(d)
})
```

3 Render with `rmarkdown::run` or click Run Document in RStudio

### Line graph

Choose a time series:

CO2

See a plot:

LH

d

Time

## 8. Publish

Share your report where users can visit it online

### Rpubs.com

Share non-interactive documents on RStudio's free R Markdown publishing site

[www.rpubs.com](http://www.rpubs.com)

### ShinyApps.io

Host an interactive document on RStudio's server. Free and paid options

[www.shinyapps.io](http://www.shinyapps.io)

## 9. Learn More

Documentation and examples - [rmarkdown.rstudio.com](http://rmarkdown.rstudio.com)

Further Articles - [shiny.rstudio.com/articles](http://shiny.rstudio.com/articles)

[@rstudio](http://blog.rstudio.com)

R Studio®

RStudio™ and Shiny™ are trademarks of RStudio, Inc.

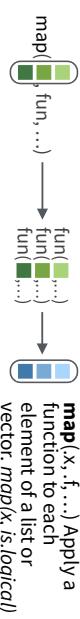
CC BY RStudio info@rstudio.com

844-448-1212 rstudio.com

# Apply functions with purrr :: CHEAT SHEET

## Apply Functions

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



**map(x, f, ...)** Apply a function to each element of a list or vector. `map(x, is.logical)`



**map2(x, y, f, ...)** Apply a function to pairs of elements from two lists. `vectors, map2(x, y, sum)`



**pmap(l, f, ...)** Apply a function to groups of elements from list of lists, vectors. `pmap(list(x, y, z), sum, na.rm = TRUE)`



**invoke\_map(f, x = list(NULL), ..., eneNULLL)** Run each function in a list. `Also invoke_map(l <- list(var, sd); invoke_map(l, x = 1:9))`

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

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

**OUTPUT**

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

**map\_dbl** double (numeric) vector  
**map\_dfc** data frame (column bind)  
**map\_dfr** data frame (row bind)

**map\_int** integer vector  
**map\_lgl** logical vector

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

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

~.1..2 etc becomes **function(...1, ...2, etc)**. `1..2 etc` e.g. `map(l, ~2+)` becomes `map(l, function(x) 2+x)`

~.becomes **function(x)x**, e.g. `map(l, "a")` becomes `map(l, function(x) a)`

## Work with Lists

### FILTER LISTS

**pluck(x, ...)** 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)`

**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)`

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

**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")`

### WORK WITH LISTS

**array** Turn array into list. Also **array**, **branch**. `array_tree(array, margin = 3)`

**cross2(x, y, filter = NULL)** All combinations of x and y. Also **cross**, **cross3**, **cross\_dfc**. `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"))`

### 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, "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)`



## 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)`

**compose()** Compose multiple functions. `compose(a, b, c)(x)`

**lift()** Change the type of input a function takes. Also **lift\_dfl**, **lift\_dfl**, **lift\_lv**, **lift\_lv\_dfl**, **lift\_lv\_id**, **lift\_lv\_id**.

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

**safely()** Modify func to return list of results and errors.

## Modify function behavior

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

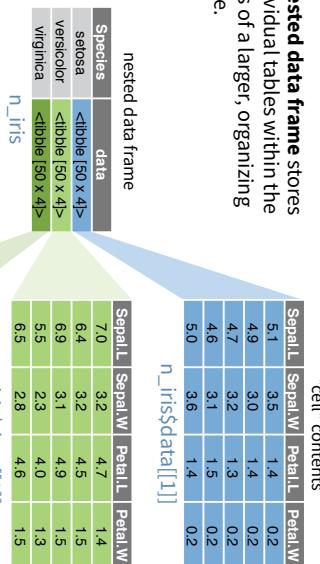
**quietly()** Modify function to return list of results, output, messages, warnings.

**partial()** Partially apply a function, filling in some args.

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

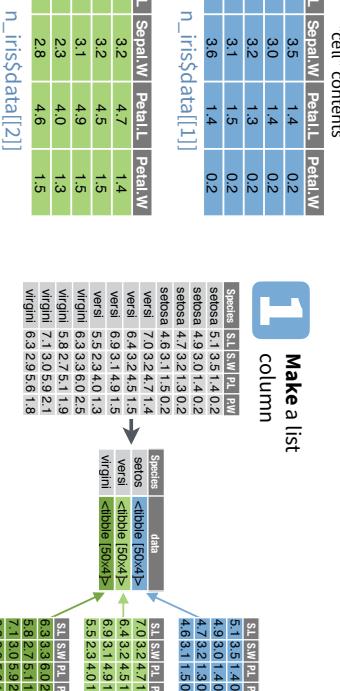
# Nested Data

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



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

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



- 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

**1. MAKE A LIST COLUMN** - You can create list columns with functions in the **tibble** and **dplyr** packages, as well as **tidyverse**'s **nest()**

**tibble::tribble(...)**  
Makes list column when needed  
**tibble(~max, ~seq, ...)**  
max: seq: n: max = c(3, 4), seq = list(1:3, 1:4, 1:5))  
**tibble::nest()**  
Creates multi-level list to tibble with list cols  
enframe(list(3=1:3, '4'=1:4, 5=1:5), 'max', '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.  
**dplyr::mutate(data, ...)** Also **transmute()**  
Returns list col when result is wrapped with **list()**  
**mtcars %>% mutate(seq = map(cyl, seq))**  
**tibble::enframe(x, name = "name", value = "value")**  
Converts multi-level list to tibble with list cols  
**summarise(q = list(quantile(mpg)))**  
**dplyr::summarise(data, ...)**  
Returns list col when result is wrapped with **list()**  
**mtcars %>% group\_by(cyl) %>% summarise(q = list(quantile(mpg)))**

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

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

# 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**

**2 Work with list columns**

**3 Simplify the list column**



# Caret Package

## Cheat Sheet

### Specifying the Model

Possible syntaxes for specifying the variables in the model:

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

```
train(y ~ x1 + x2, data = dat, ...)
```

```
train(x = predictor, df, y = outcome, vector, ...)
```

```
train(recipe_object, data = dat, ...)
```

```
rfe, sbf, gafs, and safs only have the x/y interface.
```

```
The train formula method will always create dummy
```

```
variables.
```

```
The x/y interface to train will not create dummy variables
```

```
(but the underlying model function might).
```

**Remember** to:

- Have column names in your data.

- Use factors for a classification outcome (not 0/1 or integers).

- Have valid R names for class levels (not "0"/"1")

- Set the random number seed prior to calling train repeatedly to get the same resamples across calls.

- Use the train option na.action = na.pass if you will be imputing missing data. Also, use this option when predicting new data containing missing values.

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

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

### Parallel Processing

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

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

The function parallel::detectCores can help too.

## Preprocessing

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

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

Methods include:

- "center", "scale", and "range" to normalize predictors.

- "BoxCox", "YeoJohnson", or "expoTrans" to transform predictors.

- "knnImpute", "bagImpute", or "medianImpute" to impute.

- "corr", "nzv", "zv", and "conditionalX" to filter.

- "pca", "ica", or "spatialSign" to transform groups.

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

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

### Adding Options

Many train options can be specified using the trainControl function:

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

### Random Search

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

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

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

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

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

## Performance Metrics

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

```
trainControl(summaryFunction = <R function>,
```

```
classProbs = <logical>)
```

Custom R functions can be used but caret includes several:

- "defaultSummary" (for accuracy, RMSE, etc), "twoClassSummary" (for ROC curves), and "prSummary" (for information retrieval). For the last two functions, the option classProbs must be set to TRUE.

### Grid Search

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

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

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

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

### Subsampling

# 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:

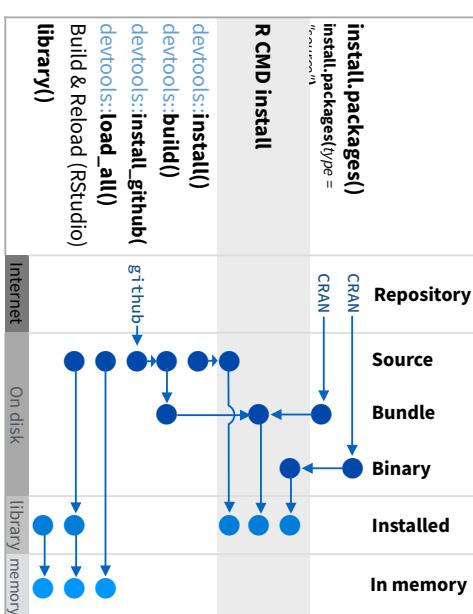
Package	DESCRIPTION	SETUP
R/	WRITE CODE	TEST
tests/	DOCUMENT	TEACH
man/	ADD DATA	
vignettes/		
data/		
NAMESPACE	ORGANIZE	

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)

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::add\_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.

<input checked="" type="checkbox"/> You must have a DESCRIPTION file
Add the packages that yours relies on with <code>devtools::use_package()</code>
Adds a package to the imports or Suggests field
No strings attached. MIT license applies to your code if re-shared.
GPL-2 license applies to your code, and all code anyone bundles with it, if re-shared.

## 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.

## WORKFLOW

1. Modify your code or tests.
2. Test your code with one of  
`devtools::test()`  
Runs all tests in tests/  
`Ctrl/Cmd + Shift + T`  
(keyboard shortcut)
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)
})

```

## Test (tests/)

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

- Add a tests/ directory
- Import tests/ 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/

```

Package: mypackage
Title: Title of Package
Version: 0.1.0
Authors@R: person("Hadley", "Wickham", email = "hadley@ml.com", role = c("aut", "cre"))
Description: What the package does (one paragraph)
Depends: R (>= 3.1.0)
License: GPL-2
LazyData: true
Imports: ggplot2 (>= 0.4.0),
knitr (>= 0.1.0)
Suggests: ggvis (>= 0.2.2)
Imports: knitr (>= 0.1.0)

```



## Test (tests/)

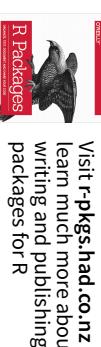
Use tests/ to store tests that will alert you if your code breaks.

- Add a tests/ directory
- Import tests/ 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/

```

Package: mypackage
Title: Title of Package
Version: 0.1.0
Authors@R: person("Hadley", "Wickham", email = "hadley@ml.com", role = c("aut", "cre"))
Description: What the package does (one paragraph)
Depends: R (>= 3.1.0)
License: GPL-2
LazyData: true
Imports: ggplot2 (>= 0.4.0),
knitr (>= 0.1.0)
Suggests: ggvis (>= 0.2.2)
Imports: knitr (>= 0.1.0)

```



## Document (□ man/)

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

- ✓ Use roxygen comments to document each function
- ✓ 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)}
\strong{bold(text)}
\code{function(args)}
\pkgs{package}
```

### COMMON ROXYGEN TAGS

```
\dontrun{code}
\dontshow{code}
\donttest{code}
\deqn{a + b}{block}
\eqn{a + b}{inline}
}
add <- function(x, y) {
 x + y
}
```

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



## 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 (<http://markdown.rstudio.com/>)

## Add Data (□ data/)

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

- ✓ Save data as .Rdata files (suggested)  
(R/Sysdata.rda if **internal = TRUE**)
- ✓ 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

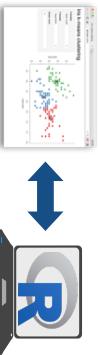
## SUBMIT YOUR PACKAGE



# 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()
server <- function(input, output){}
shinyApp(ui = ui, server = server)
```

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

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

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

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

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

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

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

works with

**shinyapps.io** 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/
```

## Inputs

collect values from the user

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

`actionLink`(inputId, label, icon, choices, selected, inline)

`checkboxGroupInput`(inputId, label, value)

`checkboxInput`(inputId, label, checked)

`dateRangeInput`(inputId, label, value, min, max, format, startView, weekStart, language, separator)

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

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

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

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

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

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

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

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

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

`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**.



<b>Trigger</b>	<code>run(this)</code>
<b>arbitrary code</b>	<code>invalidatedLater()</code>
<b>reactiveFileReader()</b>	<code>observeEvent()</code>
<b>observed</b>	<code>reactiveValues()</code>
<b>input\$</b>	<code>update</code>
<b>expression()</b>	<code>expression()</code>
<b>Modularize reactions</b>	<code>modularize()</code>
<b>reactive()</b>	<code>reactive()</code>
<b>prevent reactions</b>	<code>prevent()</code>
<b>isolate()</b>	<code>isolate()</code>
<b>CREATE YOUR OWN REACTIVE VALUES</b>	<code>reactiveValues(...)</code>
<b>* Input() functions</b> (see front page)	<code>ui &lt;- fluidPage(ui, ..., reactiveValues(...))</code>
<b>reactiveValues(...)</b>	<code>server &lt;- function(input, output){ reactiveValues(...); ... } reactiveValues() creates a list of reactive values whose values you can set.</code>
<b>RENDER REACTIVE OUTPUT</b>	<code>library(shiny) ui &lt;- fluidPage(ui, ..., renderText({ output\$ &lt;- reactiveValues(...); ... })) shinyApp(ui, server)</code>
<b>render*() functions</b> (see front page)	<code>library(shiny) ui &lt;- fluidPage(ui, ..., renderText({ output\$ &lt;- reactiveValues(...); ... })) shinyApp(ui, server)</code>
<b>renderText()</b>	<code>Builds an object to display. Will rerun code in body to rebuild the object whenever a reactive value in the code changes.</code>
<b>Save the results to output\$&lt;outputid&gt;</b>	<code>Save the results to output\$&lt;outputid&gt;</code>
<b>PREPARE REACTIONS</b>	<code>library(shiny) ui &lt;- fluidPage(ui, ..., reactiveValues(...), server &lt;- function(input, output){ reactiveValues(...); ... }) shinyApp(ui, server)</code>
<b>isolate(expr)</b>	<code>Runs a code block. Returns a <b>non-reactive</b> copy of the results.</code>
<b>RENDER ARBITRARY CODE</b>	<code>library(shiny) ui &lt;- fluidPage(ui, ..., observeEvent(eventExpr, handlerExpr, eventQuoted, handlerEnv, ...), server &lt;- function(input, output){ observeEvent(input\$go, { print(input\$go); print(input\$go); }) observeEvent(input\$go, { print(input\$go); }) ... }) shinyApp(ui, server)</code>
<b>MODULARIZE REACTIONS</b>	<code>library(shiny) ui &lt;- fluidPage(ui, ..., reactiveFileReader("a", "a", "a", reactiveValues(...)), server &lt;- function(input, output){ reactiveValues(...); ... }) shinyApp(ui, server)</code>
<b>DELAY REACTIONS</b>	<code>library(shiny) ui &lt;- fluidPage(ui, ..., eventReactive(eventExpr, eventEnv, eventQuoted, valueEnv, domain, ignoreNULL), server &lt;- function(input, output){ eventReactive(input, output, re) ... }) shinyApp(ui, server)</code>

# 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>

```

**HTML** Add static HTML elements with `tags`, a list of `tags$xa()`. Unnamed arguments will be passed into the tag; named arguments will become tag attributes.

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

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

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

**fluidRow(column)** `ui <- fluidPage( fluidRow(column(width=2, offset=3), fluidRow(column(width=4), # object1, fluidRow(column(width=2, offset=3), # object2, fluidRow(column(width=12), # object3, ...))) )`

**fluidRow(column(width=2, offset=3))** `ui <- fluidPage( fluidRow(column(width=12), sidebarLayout(sidePanel("object 1", object2, object3), mainPanel("main panel", sidebarLayout("side panel", "panel"))), ...))) )`

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

```

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

```

**Header 1**

**bold**

**italic**

**code**

**link**

**Raw html**

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

**Header 1**

**verticalLayout()** 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")))
)

```

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

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

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

**Header 1**

**verticalLayout()**



# Layouts

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

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

<b>wellPanel()</b>	<code>wellPanel(dateInput("a", "")), wellPanel(submitButton(), mainPanel())</code>
<b>navListPanel()</b>	<code>navListPanel()</code>
<b>fixedPanel()</b>	<code>fixedPanel()</code>
<b>headerPanel()</b>	<code>headerPanel()</code>
<b>tabPanel()</b>	<code>tabPanel()</code>
<b>titlePanel()</b>	<code>titlePanel()</code>
<b>inputPanel()</b>	<code>inputPanel()</code>
<b>mainPanel()</b>	<code>mainPanel()</code>

# 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**. With sparklyr, you can orchestrate distributed machine learning using either Spark's MLlib or H2O Sparkling Water.

Starting with **version 1.0.44**, **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.



## Data Science Toolchain with Spark + sparklyr

### Using 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 YARN MANAGED CLUSTER

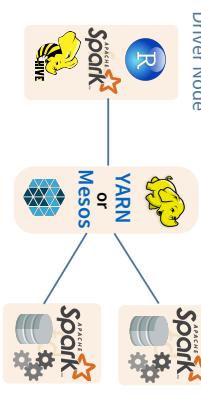
1. Install RStudio Server or RStudio Pro on one of the existing nodes, preferably an edge node
2. Locate path to the cluster's Spark Home Directory, it normally is `/usr/lib/spark`
3. Open a connection  
`spark_connect(master="yarn-client", version = "1.6.2", spark_home = [Cluster's Spark path])`

```
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(partition_iris,
 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
```

## Cluster Deployment

### MANAGED CLUSTER

Driver Node → Cluster Manager → Worker Nodes



### ON A SPARK STANDALONE CLUSTER

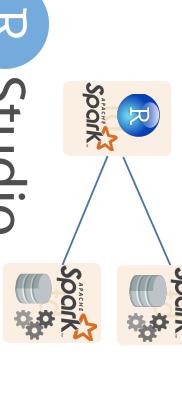
1. Install RStudio Server or RStudio Pro on one of the existing nodes or a server in the same LAN

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",
 host:port, spark_home=spark_home_dir())
```

## STAND ALONE CLUSTER

Driver Node → Worker Nodes



## 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.executor.instances
- spark.executor.extra.lavaOptions
- spark.network.timeout **120s**
- spark.executor.heartbeatInterval **10s**
- spark.executor.memory **1g**
- spark.yarn.shell.executor-memory
- spark.executor.cores **1**
- spark.yarn.driver.memory

```
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()
```

Bring data back into R memory for plotting

```
spark_disconnect(sc) Disconnect
```



# Reactivity

## COPY A DATA FRAME INTO SPARK

```
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
```

## FROM A TABLE IN HIVE

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

**CSV**

```
spark_read_csv(header = TRUE,
 columns = NULL, inferSchema = TRUE,
 delimiter = "", quote = "", escape = "\\",
 charset = "UTF-8", nullValue = NULL)
```

**JSON**

```
spark_read_json()
```

**PARQUET**

```
spark_read_parquet()
```

## 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")
```

### 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")
```

```
sf_predict(object, newdata)
```

Spark DataFrame with predicted values

# Visualize & Communicate

## SPARK SQL COMMANDS

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

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

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

Download a Spark DataFrame to an R DataFrame

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

value)

Returns contents of a single column to R

## SAVE FROM SPARK TO FILE SYSTEM

Arguments that apply to all functions: x, path

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

**JSON**

```
spark_read_json(mode = NULL)
```

**PARQUET**

```
spark_read_parquet(mode = NULL)
```

Creates a reference to the table without loading it into memory

## Reading & Writing from Apache Spark

### ML TRANSFORMERS

```
ft_binarizer(my_table$input.col == "Petal_Length",
 output.col == "betaLarge",
 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)

Column of labels into a column of label indices.

ft\_vectorAssembler()

Combine vectors into single row-vector

# Model (MLlib)

## ml\_decision\_tree

```
(my_table, response = "Species", features =
 c("PetalLength", "PetalWidth"))
```

rating.column = "rating", itemcolumn = "item", rank = 10L, regularization.parameter = 0.1, iter.max = 10L, ml.options = ml\_options())

ml\_factorization(x, usercolumn = "user", itemcolumn = "rating", itemcolumn = "item", rank = 10L, regularization.parameter = 0.1, iter.max = 10L, ml.options = ml\_options())

ml\_generalized\_linear\_regression(x, response, features, intercept = TRUE, family = gaussian(link = "identity"), iter.max = 100L, ml.options = ml\_options())

ml\_kmeans(x, centers, iter.max = 100, features = dplyr::tbl\_vars(x), 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 = 100, 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, classifier, response, features, lambda = 0, ml.options = 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 = 32, 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\_ll, metric = "f1")

ml\_tree\_feature\_importance(sc, model)

is an R interface for

Apache Spark



# The eurostat package

## R tools to access open data from Eurostat database

### Search and download

Data in the Eurostat database is stored in tables. Each table has an identifier, a short `table_code`, and a description (e.g. `tsdr420 - People killed in road accidents`). Key eurostat functions allow to find the `table_code`, download the eurostat table and polish labels in the table.

### Find the table code

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

```
library("eurostat")
query <- search_eurostat("road", type = "table")
Goods transport by road title code
1 People killed in road accidents tsdr420
2 Enterprises with broadband access tim00090
```

### Download the table

The `get_eurostat(id, time_format = "date", filters = "none", type = "code", cache = TRUE, ...)` function downloads the requested table from the Eurostat bulk download facility or from The Eurostat Web Service JSON API (if `filters` are defined). Downloaded data is cached (if column `time_format` and if table dimensions shall be kept as codes or converted to labels `type`).

```
dat <- get_eurostat(id = "tsdr420", time_format = "num")
head(dat)
#> #> unit sex geo time values
#> 1 NR T BE 1999 1079
#> 2 NR T CZ 1999 1455
#> 3 NR T DK 1999 514
#> 4 NR T EL 1999 2116
#> 5 NR T ES 1999 5738
#> 6
```

### Add labels

The `label_eurostat(x, lang = "en", ...)` gets definitions for Eurostat codes and replace them with labels in given language ("en", "fr" or "de").

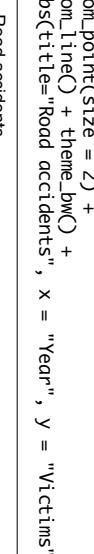
```
dat <- label_eurostat(ddat)
head(ddat)
#> #> unit sex geo time values
#> 1 Number Total Austria 1999 1079
#> 2 Number Total Belgium 1999 1397
#> 3 Number Total Czech Republic 1999 1455
#> 4 Number Total Denmark 1999 514
#> 5 Number Total Greece 1999 2116
#> 6 Number Total Spain 1999 5738
```

## eurostat and plots

The `get_eurostat()` function returns tibbles in the long format. Packages `dplyr` and `tidyR` are well suited to transform these objects. The `ggplot2` package is well suited to plot these objects.

```
t1 <- get_eurostat("tsdr420", filters = "ES", "PT"))
list(geo = c("UK", "FR", "PL", "ES", "PT"))
```

```
library("ggplot2")
ggplot(t1, aes(x = time, y = values, color = geo,
geom_point(size = 2) +
geom_line() + theme_bw() +
labs(title = "Road accidents", x = "Year", y = "Victims"))
```



## eurostat and maps

There are three functions to work with geospatial data from GISCO. The `get_eurostat_geospatial()` returns preprocessed spatial data as `sp-objects` or as data frames. The `merge_eurostat_geospatial()` both download and merges the geospatial data with a preloaded tabular data. The `cut_to_classes()` is a wrapper for `cut()` - function and is used for categorizing data for maps with tidy labels.

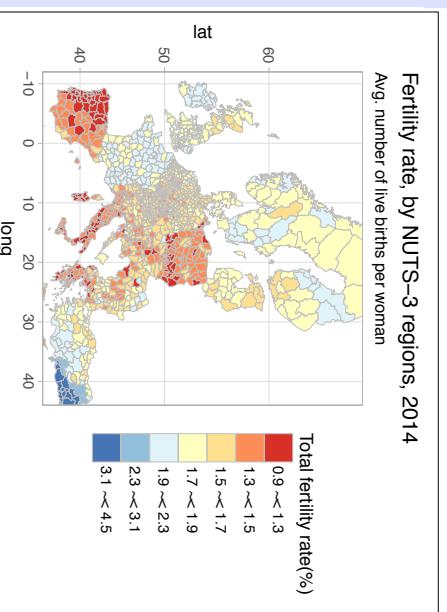
```
library("eurostat")
fertility <- get_eurostat("demo_r_frate3") %>%
filter(time == "2014-01-01") %>%
mutate(cat = cut_to_classes(values, n=7, decimal=1))
```

```
mapdata <- merge_eurostat_geodata(fertility,
resolution = "20")
head(select(mapdata, geo, values, cat, long, lat, order, id))
#> #> geo values cat long lat order id
#> 1 AT124 1.39 1.3 ~ 1.5 15.54245 48.90770 214 10
#> 2 AT124 1.39 1.3 ~ 1.5 15.75363 48.85218 215 10
#> 3 AT124 1.39 1.3 ~ 1.5 15.88763 48.78511 216 10
#> 4 AT124 1.39 1.3 ~ 1.5 15.8135 48.69270 217 10
#> 5 AT124 1.39 1.3 ~ 1.5 15.94094 48.67173 218 10
#> 6 AT124 1.39 1.3 ~ 1.5 15.90833 48.58815 219 10
```

### Draw a cartogram

The object returned by `merge_eurostat_geospatial()` are ready to be plotted with ggplot2 package. The `coord_map()` function is useful to set the projection while `labs()` adds annotations to the plot.

```
library("ggplot2")
ggplot(mapdata, aes(x = long, y = lat, group = group))+
geom_polygon(aes(fill=cat), color="grey", size = .1)+
scale_fill_brewer(palette = "RdYBu")+
labs(title="Fertility rate, by NUTS-3 regions, 2014",
subtitle="Avg. number of live births per woman",
fill="Total fertility rate(%)")+
coord_map(xlim=c(-12,44), ylim=c(35,67))
```



## Fetch and process data

There are three functions to work with geospatial data from GISCO. The `get_eurostat_geospatial()` returns preprocessed spatial data as `sp-objects` or as data frames. The `merge_eurostat_geospatial()` both download and merges the geospatial data with a preloaded tabular data. The `cut_to_classes()` is a wrapper for `cut()` - function and is used for categorizing data for maps with tidy labels.

```
library("eurostat")
fertility <- get_eurostat("demo_r_frate3") %>%
filter(time == "2014-01-01") %>%
mutate(cat = cut_to_classes(values, n=7, decimal=1))
```

```
mapdata <- merge_eurostat_geodata(fertility,
resolution = "20")
head(select(mapdata, geo, values, cat, long, lat, order, id))
#> #> geo values cat long lat order id
#> 1 AT124 1.39 1.3 ~ 1.5 15.54245 48.90770 214 10
#> 2 AT124 1.39 1.3 ~ 1.5 15.75363 48.85218 215 10
#> 3 AT124 1.39 1.3 ~ 1.5 15.88763 48.78511 216 10
#> 4 AT124 1.39 1.3 ~ 1.5 15.8135 48.69270 217 10
#> 5 AT124 1.39 1.3 ~ 1.5 15.94094 48.67173 218 10
#> 6 AT124 1.39 1.3 ~ 1.5 15.90833 48.58815 219 10
```

```
mapdata <- merge_eurostat_geodata(fertility,
resolution = "20")
head(select(mapdata, geo, values, cat, long, lat, order, id))
#> #> geo values cat long lat order id
#> 1 AT124 1.39 1.3 ~ 1.5 15.54245 48.90770 214 10
#> 2 AT124 1.39 1.3 ~ 1.5 15.75363 48.85218 215 10
#> 3 AT124 1.39 1.3 ~ 1.5 15.88763 48.78511 216 10
#> 4 AT124 1.39 1.3 ~ 1.5 15.8135 48.69270 217 10
#> 5 AT124 1.39 1.3 ~ 1.5 15.94094 48.67173 218 10
#> 6 AT124 1.39 1.3 ~ 1.5 15.90833 48.58815 219 10
```

```
mapdata <- merge_eurostat_geodata(fertility,
resolution = "20")
head(select(mapdata, geo, values, cat, long, lat, order, id))
#> #> geo values cat long lat order id
#> 1 AT124 1.39 1.3 ~ 1.5 15.54245 48.90770 214 10
#> 2 AT124 1.39 1.3 ~ 1.5 15.75363 48.85218 215 10
#> 3 AT124 1.39 1.3 ~ 1.5 15.88763 48.78511 216 10
#> 4 AT124 1.39 1.3 ~ 1.5 15.8135 48.69270 217 10
#> 5 AT124 1.39 1.3 ~ 1.5 15.94094 48.67173 218 10
#> 6 AT124 1.39 1.3 ~ 1.5 15.90833 48.58815 219 10
```

# Data & Variable Transformation with sjmisc Cheat Sheet



## Descriptives and Summaries

## Recode and Transform Variables

## Summarise Variables and Cases

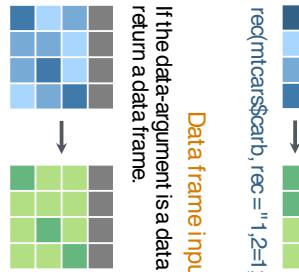
sjmisc complements dplyr, and helps with data transformation tasks and recoding variables.

sjmisc works together seamlessly with dplyr and pipes. All functions are designed to support labelled data.

### Design Philosophy

The design of sjmisc functions follows the tidyverse-approach: first argument is always the data (either a data frame or vector), followed by variable names to be processed by the functions. The returned object for each function equals the type of the data-argument.

If the data-argument is a vector, functions return a vector.



### Data Frame Input

- If the data-argument is a data frame, functions return a data frame.



Apply functions to a single variable, selected variables or to a complete data frame.

Variable selection is powered by dplyr's `select()`: Separate variables with comma, or use dplyr's select-helpers to select variables, e.g. `?select`:

```
rec(mtcars, one_of(c("gear", "carb"))),
 rec = "min:3=1; 4:max=2")
rec(mtcars, gear, carb, rec = "min:3=1; 4:max=2")
```

Most of the sjmisc functions (including recode-functions) also work on grouped data frames:

```
library(dplyr)
library(sjmisc)
group_by(e16sex, c172code) %>%
 frq(e42dep)
```

### Frequency Tables

Print frequency tables of (labelled) vectors. Uses variable labels as table header.

```
data(efc); frq(efc, e42dep, c16sex)
```

Use this data set in examples!

```
flat_table(data, ..., margin = c("counts",
 "cell", "row", "col"), digits = 2,
 show_values = FALSE)
flat_table(efc, e42dep, c172code, e16sex)
```

Print contingency tables of (labelled) vectors.

Uses value labels.

```
count_n(a, ...)
```

Print frequency/table of tagged NAs values.

```
library(lavaan); x <- labelled(c(1:3,
 tagged_na("a", "a", "z")), labels =
 c("Refused" = tagged_na("a"), "NA" =
 tagged_na("z")))
count_na(x)
```

Count NAs.

### Descriptive Summary

Descriptive summary of data frames, including variable labels in output.

```
descri(x, ..., max.length = NULL)
```

Descriptive summary of data frames, including

variable labels in output.

```
descri(efc, contains("cop"), max.length = 20)
```

### Finding Variables in a Data Frame

Use `find_var()` to search for variables by names, value or variable labels. Returns vector/data frame.

```
variables with "cop" in names and variable labels
find_var(efc, pattern = "cop", out = "gf")
variables with "level" in names and value labels
find_var(efc, "level", search = "name_value")
```

Recode functions add a suffix to new variables, so original variables are preserved. By default, only the new created variables are returned. Use `append = TRUE` to return the original input data frame as well.

```
rec(x, ..., rec, as.numeric = TRUE, var.label =
 NULL, val.labels = NULL, append =
 FALSE, suffix = "_r")
recode values, return result as numeric, character or categorical (factor).
```

```
rec(mtcars, carb, rec = "1,2=r; 3,4=2; else=3")
```

```
dichot(x, ..., dich, by = "median", as.numeric =
 FALSE, var.label = NULL, val.labels = NULL,
 append = FALSE, suffix = "d")
dichotomise variable by median, mean or specific value.
```

```
dichot(mtcars, disp)
```

```
split_var(x, ..., n, as.numeric = FALSE,
 val.labels = NULL, var.label = NULL,
 inclusive = FALSE, append = FALSE,
 suffix = "g")
split variable into equal sized groups. Unlike dplyr::unite(), does not split original categories into different values (see examples in split_var).
```

```
split_var(mtcars, mpg, disp, n = 3)
```

Split variable into groups with equal value range, or into a max # of groups (value range per group is adjusted to match # of groups).

```
group_var(mtcars, mpg, disp, size = 5)
group_var(mtcars, mpg, size = "auto", n = 4)
group_var(mtcars, mpg, size = "auto", n = 4)
```

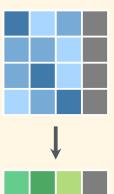
### Other Useful Functions

- `add_columns()` and `replace_columns()` to combine data frames, but either `replace` or `preserve` existing columns.
- `set_na()` and `replace_na()` to convert regular into missing values, or vice versa. `replace_na()` also replaces specific `tagged NA` values only.
- `remove_var()` and `var_rename()` to remove variables from data frames, or rename variables.
- `group_str()` to group similar string values. Useful for variables with similar but not identically spelled string values that should be "merged".
- `merge_df()` to full join data frames and preserve value and variable labels.
- `to_long()` to gather multiple columns in data frames from wide into long format.

### Use with %>% and dplyr

```
use sjmisc-functions in pipes
mtcars %>% select(gear, carb) %>%
 rec(rec = "min:3=1; 4:max=2")
use sjmisc-function inside mutate
mtcars %>% select(gear, carb) %>% mutate(
 carb2 = rec(carb, rec = "1,2=0;3:8=1"),
 gear2 = rec(gear, rec = "3=1;4:max=2"))
```

The summary functions mostly mimic base R equivalents, but are designed to work together with pipes and dplyr.



# Leaflet Cheat Sheet



## Markers

an open-source JavaScript library for mobile-friendly interactive maps

### Quick Start

Use `install.packages("leaflet")` to install the package or directly from Github `devtools::install_github("rstudio/leaflet")`

```
m <- leaflet() %>%
 addTiles() %>%
 # leaflet works with the pipe operator
 # setup the default OpenStreetMap map tiles
 addMarkers(lat = 174.768, lon = -36.852, popup = "The birthplace of R")
 # add a single point layer
```

#### Installation

Use `install.packages("leaflet")` to install the package or directly from Github `devtools::install_github("rstudio/leaflet")`

### Icon Markers

*Regular icons: default and simple*  
**addMarkers()** Add basic icon markers  
`markerOptions(...)` Create a list of icons  
`markerClusterOptions()` Marker Clusters; option of `addMarkers()`  
`clusterOptions(...)` Add awesome markers with colors and icons  
`addAwesomeMarkers(...)` Awesome icons; customizable with colors and icons  
`makeIcon(...)` Create a list of icons  
`iconList()` ShadowURL, shadowWidth, shadowHeight, ... ) customize marker icons

### Shiny Integration

To integrate a Leaflet map into an app:  
 \* In the UI, call `leafletProxy("name")`  
 \* On the server side, assign a `renderLeaflet(...)` call to the output  
 Inside the `renderLeaflet` expression, return a Leaflet map object  
 Other packages including `RJSONIO` and `jsonlite` can help fast parse or generate the data needed.

### Map Widget

#### Map Widget

`m <- leaflet(options = leafletOptions(...))`  
 Initialization  
 center Initial geographic center of the map  
 zoom Minimum map zoom level  
 minZoom Maximum zoom level of the map  
 maxZoom Maximum zoom level of the map

### Popups and Labels

`addPopups(lat, lng, ...)` Add standalone popups  
`options = popupOptions(closeDialog = FALSE)`  
`addMarkers(..., popup = ...)` Show popups with markers or shapes  
`addLabels(..., label, labelOptions = ...)` Show labels with markers or shapes  
`labelOptions = labelOptions(noHide, textOnly, textSize, direction, style)`  
`addLabelOnlyMarkers()` Add labels without markers

### Modification

To modify an existing map or add incremental changes to the map, you can use `leafletProxy()`. This should be performed in an observer on the server side.

Other useful functions to edit your map:

`fitBounds(lng1, lat1, lng2, lat2)` Similar to `setView`  
 Fit the view into the rectangle [lng1, lat1] - [lng2, lat2]  
`m %>% clearBounds()` Clear the bound, automatically determine from the map elements

### Lines and Shapes

`Polygons and PolyLines`  
`addPolygons(...)` Add standalone polygons  
`fillColor = ~colorQuantile(~order, ALAND(ALAND), highlightOptions, ...)`  
`highlightOptions(..., weight = 2, bringToFront = TRUE)` highlight shapes  
 Use `rmapshaper::ms_simplify` to simplify complex shapes  
`Circles` Add circles  
`addCircles(lng, lat, weight = 1, radius, ...)`  
`Rectangles` Add rectangles  
`addRectangles(lng1, lat1, lng2, lat2, fillColor = "transparent", ...)`

### Inputs/Events

*Object Events*  
 Object event names generally use this pattern:

`inputs$MAPID_OBJECT_EVENTNAME`  
 Trigger an event changes the value of the Shiny input at this variable.

Valid values for `OBJECTCATEGORY` are `marker`, `shape`, `geojson` and `topojson`.

Valid values for `EVENTNAME` are `click`, `mouseover` and `mouseout`.

All of these events are set to either `NULL` if the event has never happened, or a `list` that includes:

- \* `lat` The latitude of the object, if available; otherwise, the mouse cursor
- \* `lng` The longitude of the object, if available; otherwise, the mouse cursor
- \* `id` The layerId, if any

GeoJSON events also include additional properties:

`* featureId` The feature ID, if any

`* properties` The feature properties

`Map Events`

`inputs$MAPID_CLICK` When the map background or basemap is clicked

`inputs$MAPID_BOUNDS` Provide the lat/lng bounds of the visible map area

`value -- a list with north, east, south and west`

`inputs$MAPID_ZOOM` An integer indicates the zoom level

## GeoJSON and TopoJSON

There are two options to use the GeoJSON/TopoJSON data:

\* To read into `sp` objects with the `geojsonio` or `rgdal` package:  
`geojsonio::geojson_read(..., what = "sp")` `rgdal::readOGR(..., "OGRGeoJSON")`  
 \* Or to use the `addGeoJSON()` and `addTopoJSON()` functions:  
`addTopoJSON/addGeoJSON(..., weight, color, fill, opacity, fillOpacity, ...)`

### Basemaps

`Base R` `addTiles()`  
`sp package` providers\$Stamen.Toner, CartoDB.Positron, Esri.NatGeoWorldMap  
`maps package` Default tiles  
 Third-party tiles `addProviderTiles()` to add a custom map tile URL template, use `addWMSTiles()` to add WMPS (Web Map Service) tiles

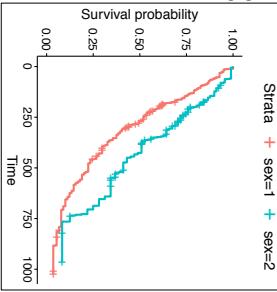
# Creating Survival Plots

## Informative and Elegant with survminer

### Survival Curves

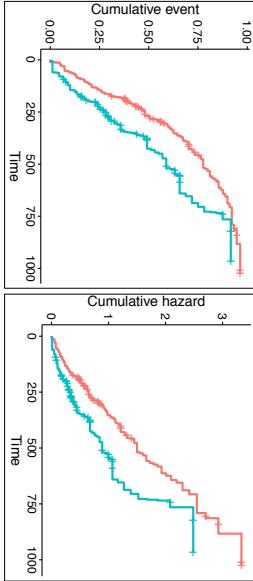
The `ggsurvplot()` function creates `ggplot2` plots from `survfit` objects.

```
library("survival")
fit <- survfit(Surv(time, status) ~ sex, data = lung)
class(fit)
[1] "survfit"
library("survminer")
ggsurvplot(fit, data = lung)
```



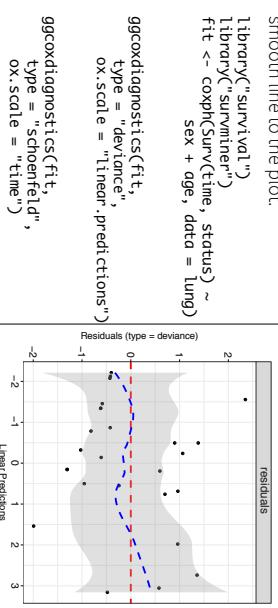
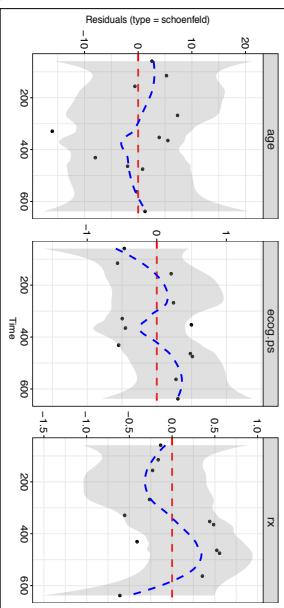
Use the `fun` argument to set the transformation of the survival curve. E.g. `"event"` for cumulative events, `"cumhaz"` for the cumulative hazard function or `"pct"` for survival probability in percentage.

```
ggsurvplot(fit, data = lung, fun = "event")
ggsurvplot(fit, data = lung, fun = "cumhaz")
```



With lots of graphical parameters you have full control over look and feel of the survival plots; position and content of the legend; additional annotations like p-value, title, subtitle.

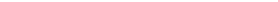
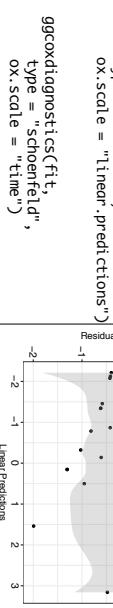
```
ggsurvplot(fit, data = lung,
conf.int = TRUE,
pval = "pct",
risk.table = TRUE,
size = 1,
line_type = "strata",
pallete = c("#728800",
 "#2E9DDE"),
legend.title = "Sex",
legend.labs = c("Male",
 "Female"))
```



The function `ggcoxdiagnostics()` plots different types of residuals as a function of time, linear predictor or observation id. The type of residual is selected with `type` argument. Possible values are "martingale", "deviance", "score", "schoenfeld", "dfbeta", "dfbeta", "dfbeta", "scaled.sch". The `ox.scale` argument defines what shall be plotted on the OX axis. Possible values are "linear.predictions", "observation.id", "time". Logical arguments `hline` and `sline` may be used to add horizontal line or smooth line to the plot.

```
library("survival")
library("survminer")
fit <- coxph(Surv(time, status) ~ sex + age, data = lung)
```

```
ggcoxdiagnostics(fit,
type = "deviance",
ox.scale = "linear.predictions")
```



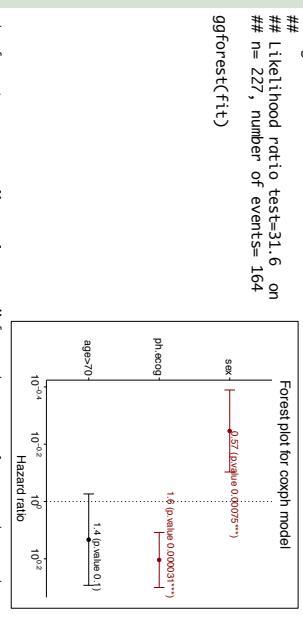
The function `ggcoxadjustedcurves()` from the `survminer` package plots Adjusted Survival Curves for Cox Proportional Hazards Model. Adjusted Survival Curves show how a selected factor influences survival estimated from a Cox model.

Note that these curves differ from Kaplan Meier estimates since they present expected survival based on given Cox model.

```
library("survival")
lung$sex <- ifelse(lung$sex == 1, "Male", "Female")
```

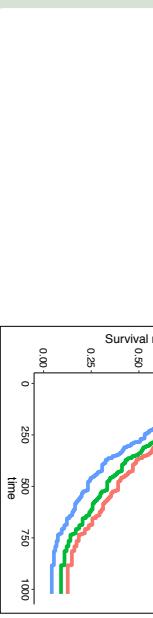
```
fit <- coxph(Surv(time, status) ~ sex + ph.ecog + age,
data = lung)
```

```
ggcoxadjustedcurves(fit, data=lung,
variable=lung$sex)
```



Note that it is not necessary to include the grouping factor in the Cox model. Survival curves are estimated from Cox model for each group defined by the factor independently.

```
lung$age3 <- cut(lung$age,
c(35, 55, 65, 85))
ggcoxadjustedcurves(fit, data=lung,
variable=lung$age3)
```

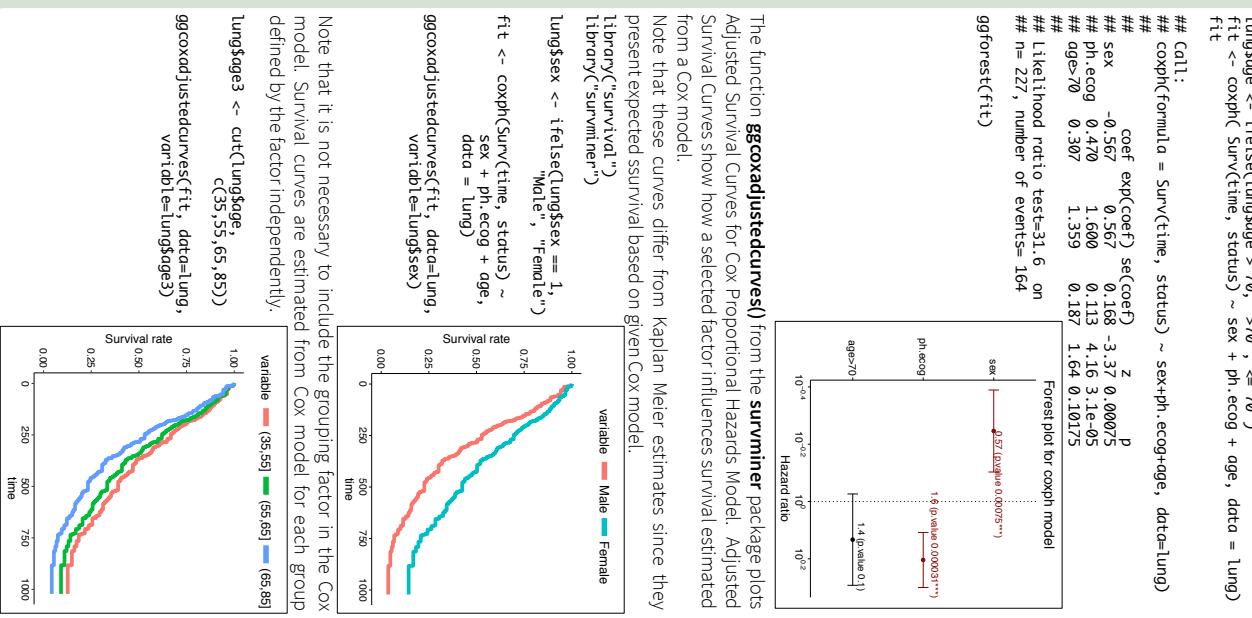


## Diagnostics of Cox Model

The function `cox.zph()` from `survival` package may be used to test the proportional hazards assumption for a Cox regression model fit. The graphical verification of this assumption may be performed with the function `ggcoxph()` from the `survminer` package. For each covariate it produces plots with scaled Schoenfeld residuals against the time.

```
library("survival")
fit <- cox.zph(Surv(time, status) ~ sex + age, data = lung)
ftest <- cox.zph(fit)

sex
rho chisq p
0.1235 2.432 0.117
age -0.0275 1.129 0.719
GLOBAL NA 2.651 0.266
library("survminer")
ggcoxph(ftest)
```



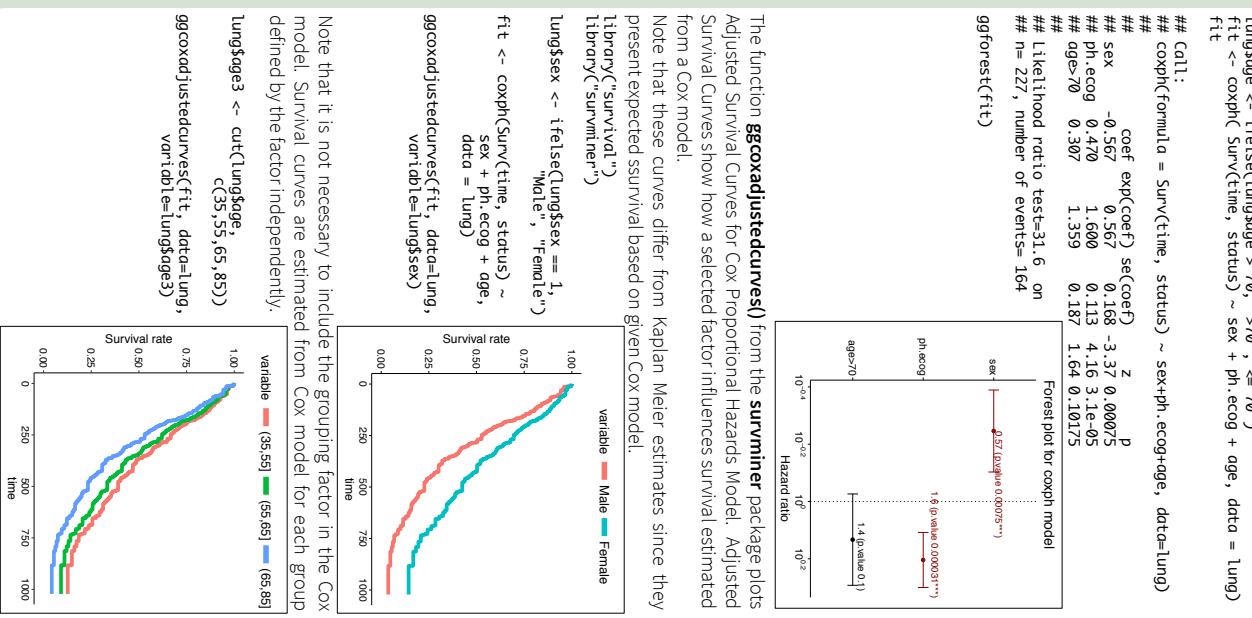
## Summary of Cox Model

The function `ggforest()` from the `survminer` package creates a forest plot for a Cox regression model fit. Hazard ratio estimates along with confidence intervals and p-values are plotted for each variable.

```
library("survival")
lung$age <- ifelse(lung$age > 70, ">70", "<= 70")
fit <- coxph(Surv(time, status) ~ sex + ph.ecog + age, data = lung)
fit
```

```
call:
coxph(formula = Surv(time, status) ~ sex+ph.ecog+age, data=lung)
coef exp(coef) se(coef) z p
sex -0.567 0.567 0.168 -3.37 0.00075
ph.ecog 0.470 1.600 0.113 4.16 3.1e-05
age>70 0.367 1.359 0.187 1.64 0.10175
n = 227, number of events = 164
```

```
ggforest(fit)
```



# Quantateda 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

- quantateda** works well with these companion packages:
- readtext**: An easy way to read text data
  - spacyr**: NLP using the spaCy library data
  - quantatedaData**: additional textual data
  - LWCalike**: R implementation of the Linguistic Inquiry and Word Count approach

## Create a corpus from texts (corpus\_\*)

```
Read texts (txt, pdf, csv, doc, docx, json, xml)
```

**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.

# 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", use_docvars = TRUE))
```

### 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, nfeture = 4)
Document-feature matrix of: 2 documents, 4 features (41.7% sparse).
```

```
docs fellow-citizens senate house representatives
1789-Washington 1 1 2 2
1793-Washington 0 0 0 0
```

### Create a dictionary

```
dictionary(x, 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)
```

### Compress a dfm by combining identical elements

```
dfm_compress(x, margin = c("both", "documents", "features"))
```

### Randomly sample documents or features

```
dfm_sample(x, what = c("documents", "features"))
```

### Weight or smooth the feature frequencies

```
dfm_weight(x, type = "relfreq") | 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

```
kwi(x, data.corpus_inaugural, "america")
```

### Utility functions

texts(corpus)	Show texts of a corpus
ndoc(corpus/dfm/tokens)	Count documents/features
nfeature(corpus/dfm/tokens)	Count features
summary(corpus/dfm)	Print summary
head(corpus/dfm)	Return first part
tail(corpus/dfm)	Return last part

## Tokenize a set of texts (tokens\_\*)

### Tokenize texts from a character vector or corpus

```
x <- tokens("Powerful tool for text analysis.",
 remove_punct = TRUE, stem = TRUE)
```

### Convert sequences into compound tokens

```
mysegs <- phrase(c("powerful", "tool", "text analysis"))
tokens_compound(x, mysegs)
```

### Select tokens

```
tokens_select(x, c("powerful", "text"), selection = "keep")
```

### Create ngrams and skipgrams from tokens

```
tokens_ngrams(x, n = 1:3)
tokens_skipgrams(toks, n = 2, skip = 0:1)
```

### Convert case of tokens

```
tokens_tolower(x) | tokens_toupper(x)
```

### Stem the terms in an object

```
tokens_wordstem(x)
```

## Calculate text statistics (textstat\_\*)

### Tabulate feature frequencies from a dfm

```
textstat_frequency(x) | topfeatures(x)
```

### Identify and score collocations from a tokenized text

```
toks <- tokens(cc("quanteda is a pkg for quant text analysis",
 "quant text analysis is a growing field"))
textstat_collocations(toks, size = 3, min_count = 2)
```

### Calculate readability of a corpus

```
textstat_reddibility(data_corpus_inaugural, measure = "Flesch")
```

### Calculate lexical diversity of a dfm

```
textstat_lexdiv(x, measure = "TRR")
```

### Measure distance or similarity from a dfm

```
textstat_simil(x, "2017-Trump", method = "cosine")
textstat_dist(x, "2017-Trump", margin = "features")
```

### Calculate keyness statistics

```
textstat_keyness(x, target = "2017-Trump")
```

## Fit text models based on a dfm (textmodel\_\*)

### Correspondence Analysis (CA)

```
textmodel_ca(x, threads = 2, sparse = TRUE, residual_floor = 0.1)
```

### Naïve Bayes classifier for texts

```
textmodel_NB(x, y = training_labels, distribution = "multinomial")
```

### Wordscores text model

```
refscores <- c(seq(-1.5, 1.5, .75), NA)
textmodel_wordscores(data_dfm_lbgexample, refscores)
```

### Wordfish Poisson scaling model

```
textmodel_wordfish(dfm(data_corpus_irishbudget2010), dir = c(6,5))
```

### Textmodel methods: predict(), coef(), summary(), print()

## Plot features or models (textplot\_\*)

### Plot features as a wordcloud

```
data_corpus_inaugural %>%
corpus_subset(Year == "Obama") %>%
dfm(remove = stopwords("english")) %>%
textplot_wordcloud()
```

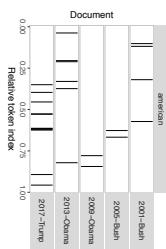
### Plot the dispersion of key word(s)

```
data_corpus_inaugural %>%
corpus_subset(Year > 1945) %>%
kwic("american") %>%
textplot_xray()
```



### Plot word keyness

```
data_corpus_inaugural %>%
corpus_subset(President %in%
c("Obama", "Trump")) %>%
dfm(groups = "President",
remove = stopwords("english")) %>%
textstat_keyness(target = "Trump") %>%
textplot_keyness()
```



### Plot Wordfish, Wordscores or CA models

```
textplot_scale1d(scaling_model,
groups = party,
margin = c("documents"))
```

## Convert dfm to a non-quanteda format

```
convert(x, to = c("lda", "tm", "stm", "austin", "topicmodels",
"lsa", "matrix", "data.frame"))
```

# R For Data Science Cheat Sheet

## data.table

Learn R for data science interactively at [www.DataCamp.com](http://www.DataCamp.com)



## data.table

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

### Load the package:

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

## Creating A data.table

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

## Subsetting Rows Using i

```
> DT[3:5]
```

Select 3rd to 5th row  
Select all rows that have value A in column V2  
Select all rows that have value A or C in column V2

## Manipulating on Columns in j

```
> DT[, V2:= "A"]
> DT[, (V2,V3)]
> DT[, sum(V1)]
> DT[, .N]
> DT[, .N, by=V1]
```

Create a data.table and call it DT  
Return v2 as a vector  
Return the sum of all elements of V1 in a vector  
Return the sum of all elements of V1, and the std. dev. of V3 in a data.table

```
> DT[, .N, by=V1]
> DT[, .N, by=V1]
```

Return v2 and v3 as a data.table  
Return the sum of all elements of V1, vector  
Aggregate sd.V3

```
> DT[, .N, by=V1]
> DT[, .N, by=V1]
```

Which returns a single value and gets recycled

## Doing j by Group

```
> DT[, .(V4 .Sum=sum(V4)), by=V1]
```

Calculate sum of V4 for every group in V1  
Calculate sum of v4 for every group in v1  
Calculate sum of v4 for every group in v1  
Calculate sum of v4 for every group in sign(v1-1)

```
> DT[, .(V4 .Sum=sum(V4)),
 by=(V1, V2)]
> DT[, .(V4 .Sum=sum(V4)),
 by=sign(V1-1)]
```

The same as the above, with new name  
for the variable you're grouping by  
Calculate sum of v4 for every group in  
sign(v1-1)

```
> DT[, .(V4 .Sum=sum(V4)),
 by=(V1-1)]
> DT[, .(V4 .Sum=sum(V4)),
 by=V1]
```

The same as the above, with new name  
for the variable you're grouping by  
Calculate sum of v4 for every group in v1  
after subsetting on the first 5 rows  
Count number of rows for every group in

```
> DT[, N, by=V1]
```

V1

## General form: DT[i, j, by] → ↴

"Take DT, subset rows using i, then calculate j grouped by by"

## Adding/Updating Columns By Reference in j Using ::

```
> DT[, V1:=round(exp(V1), 2)]
```

Columns V1 and V2 are updated by what is after ::

```
> DT[, V2:=round(exp(V1-2), LETTERS[4:6])]
```

Alternative to the above one. With [], you print the result to the screen

```
> DT[, V1:=NULL]
```

```
> DT[, c("V1", "V2") :=NULL]
```

```
> DT[, cols.choosen:=NULL]
```

```
> DT[, (cols.choosen) :=NULL]
```

Delete the column with column name

cols.choosen  
Delete the columns specified in the variable cols.choosen

## Indexing And Keys

```
> setkey(DT, V2)
```

Select key is set on V2; output is returned invisibly

```
> DT["A"]
```

Return all rows where the key column (set to V2) has the value A

```
> DT[V2=="A"]
```

Return first row of all rows that match value A in key

```
> DT["A", mult="first"]
```

Return last row of all rows that match value A in key

```
> DT["A", mult="last"]
```

Return all rows where the key column (V2) has value A or C

```
> DT[c("A", "C")]
```

Return all rows where the key column (V2) has value A or D

```
> DT[, .(Agg=aggregate=sd(V3))]
```

The same as the above, with new names

```
> DT[, .(sum(V1), sd(V3))]
```

Return the sum of all elements of V1, and the std. dev. of V3,

```
> DT[, .(V1, sd.V3=sd(V3))]
```

Select column V2 and compute std. dev. of V3,

```
> DT[, .(sum(V2), sd.V3)]
> DT[, .(sum(V2), sd.V3)]
> DT[, .(sum(V2), sd.V3)]
```

Print column V2 and plot V3

```
> DT[, .(sum(V2), sd.V3)]
```

Print column V2 and plot V3

```
> DT[, .(sum(V2), sd.V3)]
```

Print column V2 and plot V3

```
> DT[, .(sum(V2), sd.V3)]
```

Print column V2 and plot V3

## Advanced Data Table Operations

Return the penultimate row of the DT

Return the number of rows

Return v2 and v3 as a data.table

Return the result of j-grouped by all possible combinations of groups specified in by

Look at what .SD contains

Select the first and last row grouped by v2

Calculate sum of columns in .SD grouped by v2

Calculate sum of v3 and v4 in .SD grouped by v2

Look at what .SDcols contains

Select the first group of which the sum is >40 (chaining)

Select that group of which the sum is >40 (chaining)

Calculate sum of v4, grouped by v1, ordered on v1

Calculate sum of v4, grouped by v1, ordered on v1

Calculate sum of v4, grouped by v1, ordered on v1

Calculate sum of v4, grouped by v1, ordered on v1

Syntax: for (i in from:to) set(DT, row, column, new value)

rows <- list(3:4, 5:6)

cols <- 1:12

for(i in seq\_along(rows))

{set(DT,

i=rows[[i]],

j=cols[[i]],

value=NA)}

Syntax: setnames(DT, "old", "new") []

Set names of v2 to Rating (invisibly)

Sequence along the values of rows, and for the values of cols, set the values of those elements equal to NA (invisibly)

Syntax: setnames(DT, "old", "new") []

Set names of v2 to Rating (invisibly)

Sequence along the values of rows, and for the values of cols, set the values of those elements equal to NA (invisibly)

Syntax: setnames(DT, "old", "new") []

Change 2 column names (invisibly)

Sequence along the values of rows, and for the values of cols, set the values of those elements equal to NA (invisibly)

Syntax: setcolororder(DT, "neworder")

Change column ordering to contents

Sequence along the values of rows, and for the values of cols, set the values of those elements equal to NA (invisibly)

Syntax: setcolororder(DT, "neworder")

Change column ordering to contents

Sequence along the values of rows, and for the values of cols, set the values of those elements equal to NA (invisibly)

Syntax: setcolororder(DT, "neworder")

Change column ordering to contents

Sequence along the values of rows, and for the values of cols, set the values of those elements equal to NA (invisibly)

Learn Python for Data Science interactively



# R For Data Science Cheat Sheet

Learn R for data science [Interactively](#) at [www.DataCamp.com](http://www.DataCamp.com)

## xts

**eXtensible Time Series (xts)** is a powerful package that provides an extensible time series class, enabling uniform handling of many R time series classes by extending zoo.

Load the package as follows:

```
> library(xts)
```

### xts Objects

xts objects have three main components:

- coredata: always a matrix for xts objects, while it could also be a vector for zoo objects
- index: vector of any Date, POSIXct, chron, yearmon, yearqtr, or DateTime classes
- xtsAttributes: arbitrary attributes

## Creating xts Objects

```
> xts2[dates] <- 0
> xts2["1961"] <- NA
> xts2["2016-05-02"] <- NA
```

## Convert To And From xts

```
> data(AirPassengers)
> xts5 <- as.xts(AirPassengers)
```

## Import From Files

```
> dat <- read.csv(tmp_file)
> dat$dat_order <- Sys.Date()
> dat$dat_order <- as.POSIXct(dat$dat_order, format = "%Y-%m-%d")
> dat$dat_order <- as.POSIXct(dat$dat_order, format = "%Y-%m-%d %H:%M:%S")
> dat_zoo <- read.zoo(tmp_file,
 index.column=0,
 sep=",",
 format="%m/%d/%Y")
> dat_zoo <- read.zoo(tmp, sep=",", FUN=as.yearmon)
> dat_zts <- as.xts(dat_zoo)
```

## Inspect Your Data

```
> core_data <- coredata(xts5)
> index(xts5) Extract core data of objects
 Extract index of objects
```

### Class Attributes

> indexClass(xts2)	Get index class
> indexIndex(xts5)	Replacing index class
> indexFormat(xts5) <- "%Y-%m-%d"	Change format of time display

### Time Zones

> zone(xts1) <- "Asia/Hong_Kong"	Change the time zone
> tzzone(xts1)	Extract the current time zone



Xts

## Export xts Objects

```
> date_xts <- as.xts(matrix())
> tmp <- tempfile()
> write.zoo(date_xts, sep=",", file=tmp)
```

## Replace & Update

> xts2[dates] <- 0	Replace values in xts2 on dates with 0
> xts2["1961"] <- NA	Replace dates from 1961 with NA
> xts2["2016-05-02"] <- NA	Replace the value at 1 specific index with NA

## Applying Functions

> ep1 <- endpoints(xts4, on="weeks", k=2)	Take index values by time
[1] 0 5 10	
> ep2 <- endpoints(xts5, on="years")	
[1] 0 12 24 36 48 60 72 84 96 108 120 132 144	
> period.apply(xts5, INDEX=ep2, FUN=mean)	Calculate the yearly mean
> xts5\$yearly <- split(xts5, f="years")	Split xts by year
> lapply(xts5\$yearly, FUN=mean)	Create a list of yearly means
> do.call(rbind,	Find the last observation in each year in xts5
lapply(split(xts5, "years"),	
function(w) last(w, n="1 month"))	Calculate cumulative annual passengers
> rollapply(xts5, 3, sd)	Apply sd to rolling margins of xts5

## Select, Subsetting & Indexing

> mar55 <- xts5["1955-03"]	Get value for March 1955
> xts5[1954]	Get all data from 1954
> xts5["1954/1954-03"]	Extract data from Jan to March '54
> xts5[ep1]	Subset xts4 using ep2

### first() and last()

> first(xts4, "1 week")	Extract first week
> first(last(xts4), "1 week")	Get first 3 days of the last week of data

### Indexing

> xts2[index(xts3)]	Extract rows with the index of xts3
> xts3[days]	Extract rows using the vector days
> xts2(as.POSIXct(days, tz="UTC"))	Extract rows using days as POSIXct
> index <- which(indexday(xts1)==0, indexday(xts1)==6)	Index of weekend days
> xts1[index]	Extract weekend days of xts1

## Periods, Periodicity & Timestamps

> periodicity(xts5)	Estimate frequency of observations
> to_yearly(xts5)	Convert xts5 to yearly OHLC
> to_monthly(xts3)	Convert xts3 to monthly OHLC
> to_quarterly(xts5)	Convert xts5 to quarterly OHLC
> to_period(xts5, period="quarters")	Convert to quarterly OHLC
> to.period(xts5, period="years")	Convert to yearly OHLC
> nmonths(xts5)	Count the months in xts5
> nquarters(xts5)	Count the quarters in xts5
> nyears(xts5)	Count the years in xts5
> makeIndex.unique(xts3, eps=1e-4)	Remove duplicate times
> alignTime(xts3, n=3600)	Round index time to the next n seconds

## Other Useful Functions

> .index(xts4)	Extract raw numeric index of xts4
> .indexDay(xts3)	Value of week/day starting on Sunday, in index of xts3
> .indexHour(xts3)	Value of hour in index of xts3
> start(xts3)	Extract first observation of xts4
> end(xts4)	Display structure of xts3
> time(xts3)	Extract raw numeric index of xts1
> head(xts2)	First part of xts2
> tail(xts2)	Last part of xts2

## Missing Values

> na.omit(xts5)	Omit NA values in xts5
> xts_last <- na.locf(xts2)	Fill missing values in xts2 using last observation
> xts_last <- na.locf(xts2, fromLast=TRUE)	Fill missing values in xts2 using next observation
> na.approx(xts2)	Interpolate NAs using linear approximation

## Arithmetic Operations

> coredata() or as.numeric()	Addition
> coredata(xts4) * xts3	Multiplication
> coredata(xts4) - xts3	Subtraction
> coredata(xts4) / xts3	Division

## Shifting Index Values

> xts5 + lag(xts5) <sup>ep1</sup>	Period-over-period differences
> diff(xts5, lag=1, differences=1)	Lagged differences

## Reindexing

> xts5 + merge(xts2, index(xts1), fill=0)	Addition
> xts1 - merge(xts2, index(xts1), fill=na.locf)	Subtraction
> xts5 - lag(xts5) <sup>ep1</sup>	Subtraction
> xts5 - 1:5	Subtraction
> xts5[1:5]	Subtraction
> merge(xts2, xts1, join="left", fill=0)	Left join of xts2 and xts1, fill empty spots with 0
> merge(xts2, xts1, join="left", fill=0)	Inner join of xts2 and xts1
> rbind(xts1, xts4)	Combine xts1 and xts4 by rows

