



READY RECKONER

Compiled by

MACHINE LEARNING SHAstra

$y = f(x)$
MLShastra.com

RStudio IDE :: CHEAT SHEET

Documents and Apps

   Open Shiny, R Markdown, knitr, Sweave, LaTeX, .Rd files and more in Source Pane

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

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

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

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

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

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

Debug Mode

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

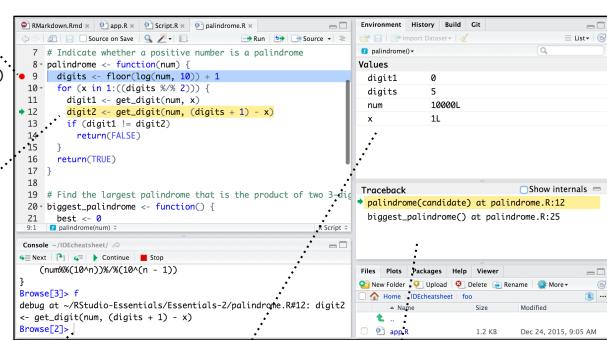
Click next to line number to add/remove a breakpoint.

Highlighted line shows where execution has paused

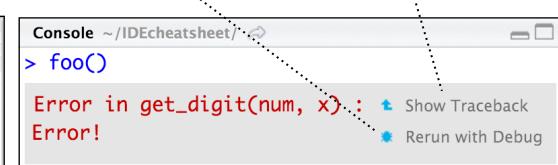
Run commands in environment where execution has paused

Examine variables in executing environment

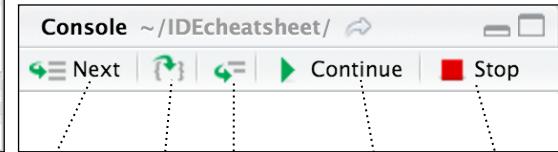
Select function in traceback to debug



Launch debugger mode from origin of error



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



Step through code one line at a time

Step into and out of functions to run

Resume execution mode

R Support

 Import data with wizard

History of past commands to run/copy

Display .RPres slideshows
File > New File > R Presentation

 Load workspace

Save workspace

Delete all saved objects

Search inside environment

Choose environment to display from list of parent environments

Display objects as list or grid

 Data

Values

Functions

View in data viewer

View function source code

Displays saved objects by type with short description

Change file type

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

Jump to function in file

Working Directory

Maximize, minimize panes

Drag pane boundaries

Path to displayed directory

Press ↑ to see command history

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

Pro Features

 Share Project

Active shared with Collaborators



Start new R Session in current project

Close R Session in project

Select R Version

PROJECT SYSTEM
File > New Project

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

RStudio opens plots in a dedicated Plots pane

Files Plots Packages Help Viewer
Run Plot Open in recent plots Export Delete plot Delete all plots

GUI Package manager lists every installed package

Files Plots Packages Help Viewer
Install Packages Update Packages Create reproducible package library for your project

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

Package version installed Delete from library

RStudio opens documentation in a dedicated Help pane

Files Plots Packages Help Viewer
R: Arithmetic Operators Find in Topic Home page of helpful links Search within help file Search for help file

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

Files Plots Packages Help Viewer
Stop Shiny app Publish to shinyapps.io, rpubs, RSConnect, ... Refresh

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

Filter	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
All	5.1	3.5	1.4	0.2	setosa
1					
2					
3					
4					
Filter rows by value or value range	Sort by values	Search for value			

1 LAYOUT

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

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

2 RUN CODE

Search command history

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

Interrupt current command

Clear console

Quit Session (desktop only)

Restart R Session

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

Source the current file

Source with echo

Windows/Linux Mac

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

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

Ctrl+Shift+F10 Cmd+Shift+F10

Ctrl+Enter Cmd+Enter
Alt+Enter Option+Enter
Ctrl+Alt+E Cmd+Option+E
Ctrl+Alt+F Cmd+Option+F
Ctrl+Alt+G Cmd+Option+G

Ctrl+Shift+S Cmd+Shift+S

Ctrl+Shift+Enter Cmd+Shift+Enter

3 NAVIGATE CODE

Goto File/Function

Fold Selected Alt+L
Unfold Selected Shift+Alt+L
Fold All Alt+O
Unfold All Shift+Alt+O
Go to line Shift+Alt+G
Jump to Shift+Alt+J
Switch to tab Ctrl+Shift+.
Previous tab Ctrl+F11
Next tab Ctrl+F12
First tab Ctrl+Shift+F11
Last tab Ctrl+Shift+F12
Navigate back Ctrl+F9
Navigate forward Ctrl+F10
Jump to Brace Ctrl+P
Select within Braces Ctrl+Shift+Alt+E
Use Selection for Find Ctrl+F3
Find in Files Ctrl+Shift+F
Find Next Win: F3, Linux: Ctrl+G
Find Previous W: Shift+F3, L:
Jump to Word Ctrl+←/→
Jump to Start/End Ctrl+↑/↓
Toggle Outline Ctrl+Shift+O

Windows /Linux

Mac Ctrl+.
Ctrl+.
Cmd+Option+L
Cmd+Shift+Option+L
Cmd+Option+O
Cmd+Shift+Option+O
Cmd+Shift+Option+G
Cmd+Shift+Option+J
Ctrl+Shift+.
Ctrl+F11
Ctrl+F12
Ctrl+Shift+F11
Ctrl+Shift+F12
Ctrl+F9
Ctrl+F10
Ctrl+P
Ctrl+Shift+Alt+E
Cmd+E
Cmd+Shift+F
Cmd+G
Cmd+Shift+G
Option+←/→
Cmd+↑/↓
Ctrl+Shift+O

4 WRITE CODE

Attempt completion

Navigate candidates
Accept candidate
Dismiss candidates
Undo Ctrl+Z
Redo Ctrl+Shift+Z
Cut Ctrl+X
Copy Ctrl+C
Paste Ctrl+V
Select All Ctrl+A
Delete Line Ctrl+D

Select Shift+[Arrow]
Select Word Ctrl+Shift+←/→

Select to Line Start Alt+Shift+←

Select to Line End Alt+Shift+→

Select Page Up/Down Shift+PageUp/Down

Select to Start/End Shift+Alt+↑/↓

Delete Word Left Ctrl+Backspace

Delete Word Right

Delete to Line End

Delete to Line Start

Indent Tab (at start of line)

Outdent Shift+Tab

Yank line up to cursor Ctrl+U

Yank line after cursor Ctrl+K

Insert yanked text Ctrl+Y

Insert <-

Insert %>%

Show help for function F1

Show source code F2

New document Ctrl+Shift+N

New document (Chrome) Ctrl+Alt+Shift+N

Open document Ctrl+O

Save document Ctrl+S

Close document Ctrl+W

Close document (Chrome) Ctrl+Alt+W

Close all documents Ctrl+Shift+W

Extract function Ctrl+Alt+X

Extract variable Ctrl+Alt+V

Reindent lines Ctrl+I

(Un)Comment lines

Reflow Comment Ctrl+Shift+/

Reformat Selection Ctrl+Shift+A

Select within braces Ctrl+Shift+E

Show Diagnostics Ctrl+Shift+Alt+P

Transpose Letters Ctrl+T

Move Lines Up/Down Alt+↑/↓

Copy Lines Up/Down Shift+Alt+↑/↓

Add New Cursor Above Ctrl+Alt+Up

Add New Cursor Below Ctrl+Alt+Down

Move Active Cursor Up Ctrl+Alt+Shift+Up

Move Active Cursor Down Ctrl+Alt+Shift+Down

Find and Replace Ctrl+F

Use Selection for Find Ctrl+F3

Replace and Find Ctrl+Shift+J

Windows /Linux

Tab or Ctrl+Space

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

Option+Shift+ ←/→

Alt+Shift+←

Alt+Shift+→

Shift+PageUp/Down

Shift+Alt+↑/↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Option+-

Ctrl+Shift+M

F1

F2

Cmd+Shift+N

Cmd+Shift+Opt+N

Cmd+O

Cmd+S

Cmd+W

Cmd+Option+W

Cmd+Shift+W

Cmd+Alt+X

Cmd+Option+X

Cmd+Option+V

Cmd+I

Ctrl+Shift+C

Ctrl+Shift+/

Ctrl+Shift+A

Ctrl+Shift+E

Ctrl+Shift+Opt+P

Ctrl+T

Option+↑/↓

Shift+Alt+↑/↓

Cmd+Option+↑/↓

Ctrl+Alt+Up

Ctrl+Alt+Down

Ctrl+Alt+Shift+Up

Ctrl+Alt+Shift+Down

Ctrl+F

Ctrl+F3

Ctrl+Shift+J

Mac

Tab or Cmd+Space

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

Option+Shift+ ←/→

Alt+Shift+←

Alt+Shift+→

Shift+PageUp/Down

Shift+Alt+↑/↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Option+-

Ctrl+Shift+M

F1

F2

Cmd+Shift+N

Cmd+Shift+Opt+N

Cmd+O

Cmd+S

Cmd+W

Cmd+Option+W

Cmd+Shift+W

Cmd+Alt+X

Cmd+Option+X

Cmd+I

Cmd+Shift+C

Cmd+Shift+/

Cmd+Shift+A

Cmd+Shift+E

Cmd+Shift+Opt+P

Cmd+T

Option+↑/↓

Shift+Alt+↑/↓

Cmd+Option+↑/↓

Ctrl+Alt+Up

Ctrl+Alt+Down

Ctrl+Alt+Shift+Up

Ctrl+Alt+Shift+Down

Ctrl+F

Ctrl+F3

Ctrl+Shift+J

WHY RSTUDIO SERVER PRO?

RSP extends the open source server with a commercial license, support, and more:

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



5 DEBUG CODE

Toggle Breakpoint
Execute Next Line
Step Into Function

Base R Cheat Sheet

Getting Help

Accessing the help files

?mean

Get help of a particular function.

help.search('weighted mean')

Search the help files for a word or phrase.

help(package = 'dplyr')

Find help for a package.

More about an object

str(iris)

Get a summary of an object's structure.

class(iris)

Find the class an object belongs to.

Using Packages

install.packages('dplyr')

Download and install a package from CRAN.

library(dplyr)

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

dplyr::select

Use a particular function from a package.

data(iris)

Load a built-in dataset into the environment.

Working Directory

getwd()

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

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

Change the current working directory.

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

Vectors

Creating Vectors

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

Vector Functions

sort(x)

Return x sorted.

rev(x)

Return x reversed.

table(x)

See counts of values.

unique(x)

See unique values.

Selecting Vector Elements

By Position

x[4]

The fourth element.

x[-4]

All but the fourth.

x[2:4]

Elements two to four.

x[!(2:4)]

All elements except two to four.

x[c(1, 5)]

Elements one and five.

By Value

x[x == 10]

Elements which are equal to 10.

x[x < 0]

All elements less than zero.

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

Elements in the set 1, 2, 5.

Named Vectors

x['apple']

Element with name 'apple'.

Programming

For Loop

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

Example

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

While Loop

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

Example

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

Functions

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

Example

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

Reading and Writing Data

Also see the **readr** package.

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

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

Types

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

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

Maths Functions

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

Variable Assignment

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

The Environment

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

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

Matrices

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

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

Lists

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

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

Also see the `dplyr` package.

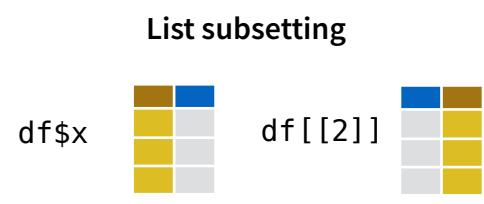
Data Frames

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

x	y
1	a
2	b
3	c

Matrix subsetting

<code>df[, 2]</code>	
<code>df[2,]</code>	
<code>df[2, 2]</code>	



Understanding a data frame
`View(df)` See the full data frame.
`head(df)` See the first 6 rows.

`nrow(df)` Number of rows.
`ncol(df)` Number of columns.
`dim(df)` Number of columns and rows.

`cbind` - Bind columns.

`rbind` - Bind rows.

Values of x in order.

Strings

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

Factors

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

Statistics

<code>lm(y ~ x, data=df)</code>	Linear model.
<code>glm(y ~ x, data=df)</code>	Generalised linear model.
<code>summary</code>	Get more detailed information out a model.
<code>pairwise.t.test</code>	Perform a t-test for paired data.

Distributions

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

Plotting

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

Dates

See the `lubridate` package.

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()`

1. **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')`

2. **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

1. **Empty environment** – ultimate ancestor of all environments
 - Parent: none
 - Access with: `emptyenv()`

2. **Base environment** - environment of the base package
 - Parent: empty environment
 - Access with: `baseenv()`

3. **Global environment** – the interactive workspace that you normally work in
 - Parent: environment of last attached package
 - Access with: `globalenv()`

4. **Current environment** – environment that R is currently working in (may be any of the above and others)
 - Parent: empty environment
 - Access with: `environment()`

Environments

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')`

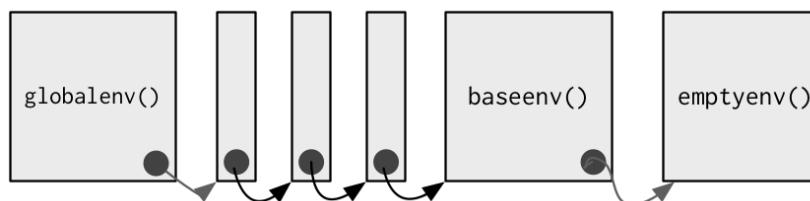


Figure 1 – The Search Path

- 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'  
library(reshape2); search()  
'.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

Binding Names to Values

Assignment – act of binding (or rebinding) a name to a value in an environment.

1. `<-` (Regular assignment arrow) – always creates a variable in the current environment
2. `<--` (Deep assignment arrow) - modifies an existing variable found by walking up the parent environments

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

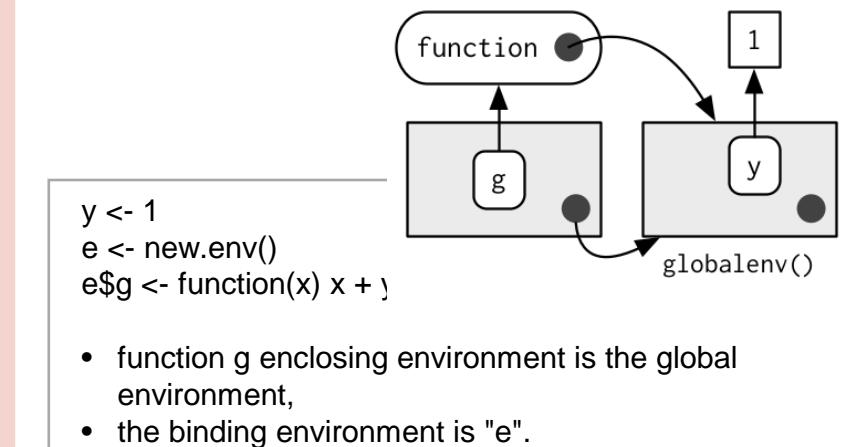
Function Environments

1. **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.
 - Access with: `environment('func1')`

2. **Binding environment** - all environments that the function has a binding to. It determines how we find the function.

- Access with: `pryr::where('func1')`

Example (for enclosing and binding environment):



- function g enclosing environment is the global environment,
- the binding environment is "e".

3. **Execution environment** - new created environments to host a function call execution.

- Two parents :
 - I. Enclosing environment of the function
 - II. Calling environment of the function
- Execution environment is thrown away once the function has completed.

4. **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

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
- Access with: **typeof()**

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

- 'type' of the object from R's object-oriented programming point of view
- Access with: **class()**

	typeof()	class()
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:	canvas.drawRect('blue')
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')`.

S3

1. 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.

2. Notation :

- **generic.class()**

mean.Date()	Date method for the generic - mean()
-------------	--------------------------------------

3. Useful 'Generic' Operations

- Get all methods that belong to the 'mean' generic:
 - **Methods('mean')**
- 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) or pryr::oGetType()</code>
Check if object inherits from a specific class	<code>inherits(x, 'classname')</code>
Determine class of any object	<code>class(x)</code>

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

- **Internal representation** : C structure (or struct) that includes :

- Contents of the object
- Memory Management Information
- Type
 - Access with: **typeof()**

Functions

Function Basics

Functions – objects in their own right

All R functions have three parts:

body()	code inside the function
formals()	list of arguments which controls how you can call the function
environment()	“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.

Function Arguments

Arguments – passed by reference and copied on modify

1. Arguments are matched first by exact name (perfect matching), then by prefix matching, and finally by position.
2. Check if an argument was supplied : **missing()**

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

3. Lazy evaluation – since x is not used **stop("This is an error!")** never get evaluated.

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

4. Force evaluation

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

5. Default arguments evaluation

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

f() -> 'a' 'x'	ls() evaluated inside f
f(ls())	ls() evaluated in global environment

Return Values

- **Last expression evaluated or explicit return()**. Only use explicit return() when returning early.
- **Return ONLY single object**. Workaround is to return a list containing any number of objects.
- **Invisible return object value** - not printed out by default when you call the function.

```
f1 <- function() invisible(1)
```

Primitive Functions

What are Primitive Functions?

1. Call C code directly with **.Primitive()** and contain no R code

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

2. **formals()**, **body()**, and **environment()** are all NULL

3. Only found in base package

4. More efficient since they operate at a low level

Influx Functions

What are Influx Functions?

1. Function name comes in between its arguments, like + or -
2. All user-created infix functions must start and end with %.

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

3. 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
```

Replacement Functions

What are Replacement Functions?

1. Act like they modify their arguments in place, and have the special name xxx <-
2. 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
```

Subsetting

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

Simplifying vs. Preserving Subsetting

1. Simplifying subsetting

- Returns the **simplest** 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
Vector	x[[1]]	x[1]
List	x[[1]]	x[1]
Factor	x[1:4, drop = T]	x[1:4]
Array	x[1,] or x[, 1]	x[1, , drop = F] or x[, 1, drop = F]
Data frame	x[, 1] or x[[1]]	x[, 1, drop = F] or x[1]

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(df1[, 'col1']) -> int [1:3]
```

- Result: the result is a vector

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

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.

```
df1$countCol is c(3, 5, 1)
rep(1:nrow(df1), df1$countCol)
> 1 1 1 2 2 2 2 2 3
```

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$col1 == 5 & df1$col2 == 4, ]
```

Subsetting continued

Boolean Algebra vs. Sets (Logical and Integer Subsetting)

1. **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.

2. **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 :
 - I. Use **x[which(y)]** instead of **x[y]**
 - II. **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

Subsetting with Assignment

1. All subsetting operators can be combined with assignment to modify selected values of the input vector.

```
df1$col1[df1$col1 < 8] <- 0
```

2. Subsetting with nothing in conjunction with assignment :

- Why : Preserve original object class and structure

```
df1[] <- lapply(df1, as.integer)
```

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
- Creates **last.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)
```

In a later interactive session :

```
load("last.dump.rda")  
debugger()
```

Condition Handling of Expected Errors

1. Communicating potential problems to users:

I. stop()

- Action : raise fatal error and force all execution to terminate
- Example usage : when there is no way for a function to continue

II. warning()

- Action : generate warnings to display potential problems
- Example usage : when some of elements of a vectorized input are invalid

III. message()

- Action : generate messages to give informative output
- Example usage : when you would like to print the steps of a program execution

2. Handling conditions programmatically:

I. try()

- Action : gives you the ability to continue execution even when an error occurs

II. 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

1. **stopifnot()** or use 'assertthat' package - check inputs are correct

2. **Avoid subset(), transform() and with()** - these are non-standard evaluation, when they fail, often fail with uninformative error messages.

3. **Avoid [and sapply()** - functions that can return different types of output.

- Recommendation : Whenever subsetting a data frame in a function, you should always use **drop = FALSE**

caret Package

Cheat Sheet

Specifying the Model

Possible syntaxes for specifying the variables in the model:

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

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

Remember to:

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

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

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

Parallel Processing

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

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

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

Preprocessing

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

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

Methods include:

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

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

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

Adding Options

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

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

Resampling Options

`trainControl` is used to choose a resampling method:

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

Methods and options are:

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

Performance Metrics

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

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

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

Grid Search

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

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

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

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

Random Search

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

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

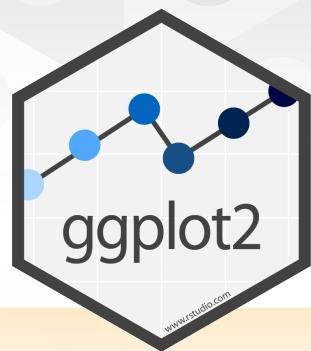
Subsampling

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

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

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

Data Visualization with ggplot2 :: CHEAT SHEET



Basics

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



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



Complete the template below to build a graph.

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

required

Not required, sensible defaults supplied

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

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

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

last_plot() Returns the last plot

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

Geoms

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

GRAPHICAL PRIMITIVES

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

LINE SEGMENTS

- common aesthetics: x, y, alpha, color, linetype, size
- b + geom_abline(aes(intercept=0, slope=1))**
 - b + geom_hline(aes(yintercept = lat))**
 - b + geom_vline(aes(xintercept = long))**
 - b + geom_segment(aes(yend=lat+1, xend=long+1))**
 - b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

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

discrete

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

TWO VARIABLES

continuous x , continuous y

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

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

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

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

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

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

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

discrete x , continuous y

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

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

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

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

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

discrete x , discrete y

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

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

THREE VARIABLES

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

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

continuous bivariate distribution

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

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

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

continuous function

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

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

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

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

visualizing error

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

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

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

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

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

maps

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

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

Data Import :: CHEAT SHEET

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

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

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

OTHER TYPES OF DATA

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

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

Save Data

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

Comma delimited file

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

File with arbitrary delimiter

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

CSV for excel

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

String to file

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

String vector to file, one element per line

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

Object to RDS file

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

Tab delimited files

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



Read Tabular Data

- These functions share the common arguments:

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

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

Comma Delimited Files

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

To make file.csv run:

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

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

Semi-colon Delimited Files

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

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

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

Files with Any Delimiter

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

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

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

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")
```

a\tb\tc 1\t2\t3 4\t5\tNA	A B C 1 2 3 4 5 NA
a\tb\tc 1\t2\t3 4\t5\tNA	A B C 1 2 3 4 5 NA

Tab Delimited Files

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

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

USEFUL ARGUMENTS

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

Example file

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

1 2 3
4 5 NA

Skip lines

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

A B C
1 2 3
4 5 NA

No header

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

A B C
1 2 3

Read in a subset

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

A B C
NA 2 3
4 5 NA

Provide header

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

A B C
NA 2 3
4 5 NA

Missing Values

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

Read Non-Tabular Data

Read a file into a single string

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

Read each line into its own string

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

Read Apache style log files

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

Read a file into a raw vector

```
read_file_raw(file)
```

Read each line into a raw vector

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

Data types

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

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

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

age is an integer

sex is a character

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

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

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

- **col_guess()** - the default

- **col_character()**

- **col_double()**, **col_euro_double()**

- **col_datetime(format = "")** Also **col_date(format = "")**, **col_time(format = "")**

- **col_factor(levels, ordered = FALSE)**

- **col_integer()**

- **col_logical()**

- **col_number()**, **col_numeric()**

- **col_skip()**

```
x <- read_csv("file.csv", col_types = cols(  
  A = col_double(),  
  B = col_logical(),  
  C = col_factor()))
```

Tibbles - an enhanced data frame

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

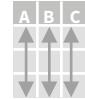
- **Subsetting** - [always returns a new tibble, [[and \$ always return a vector.
- **No partial matching** - You must use full column names when subsetting
- **Display** - When you print a tibble, R provides a concise view of the data that fits on one screen

# A tibble: 234 x 6	manufacturer	model	displ	year	cyl
1 audi	a4	1.8	1999	4	
2 audi	a4	1.8	1999	4	
3 audi	a4	2.0	1999	4	
4 audi	a4	2.0	1999	4	
5 audi	a4	2.0	1999	4	
6 audi	a4	2.0	1999	4	
7 audi	a4	3.1	1999	6	
8 audi	a4 quattro	1.8	1999	4	
9 audi	a4 quattro	1.8	1999	4	
10 audi	a4 quattro	1.8	1999	4	
11 audi	a4 quattro	1.8	1999	4	
12 audi	a4 quattro	1.8	1999	4	
13 audi	a4 quattro	1.8	1999	4	
14 audi	a4 quattro	1.8	1999	4	
15 audi	a4 quattro	1.8	1999	4	
16 audi	a4 quattro	1.8	1999	4	
17 audi	a4 quattro	1.8	1999	4	
18 audi	a4 quattro	1.8	1999	4	
19 audi	a4 quattro	1.8	1999	4	
20 audi	a4 quattro	1.8	1999	4	
21 audi	a4 quattro	1.8	1999	4	
22 audi	a4 quattro	1.8	1999	4	
23 audi	a4 quattro	1.8	1999	4	
24 audi	a4 quattro	1.8	1999	4	
25 audi	a4 quattro	1.8	1999	4	
26 audi	a4 quattro	1.8	1999	4	
27 audi	a4 quattro	1.8	1999	4	
28 audi	a4 quattro	1.8	1999	4	
29 audi	a4 quattro	1.8	1999	4	
30 audi	a4 quattro	1.8	1999	4	
31 audi	a4 quattro	1.8	1999	4	
32 audi	a4 quattro	1.8	1999	4	
33 audi	a4 quattro	1.8	1999	4	
34 audi	a4 quattro	1.8	1999	4	
35 audi	a4 quattro	1.8	1999	4	
36 audi	a4 quattro	1.8	1999	4	
37 audi	a4 quattro	1.8	1999	4	
38 audi	a4 quattro	1.8	1999	4	
39 audi	a4 quattro	1.8	1999	4	
40 audi	a4 quattro	1.8	1999	4	
41 audi	a4 quattro	1.8	1999	4	
42 audi	a4 quattro	1.8	1999	4	
43 audi	a4 quattro	1.8	1999	4	
44 audi	a4 quattro	1.8	1999	4	
45 audi	a4 quattro	1.8	1999	4	
46 audi	a4 quattro	1.8	1999	4	
47 audi	a4 quattro	1.8	1999	4	
48 audi	a4 quattro	1.8	1999	4	
49 audi	a4 quattro	1.8	1999	4	
50 audi	a4 quattro	1.8	1999	4	
51 audi	a4 quattro	1.8	1999	4	
52 audi	a4 quattro	1.8	1999	4	
53 audi	a4 quattro	1.8	1999	4	
54 audi	a4 quattro	1.8	1999	4	
55 audi	a4 quattro	1.8	1999	4	
56 audi	a4 quattro	1.8	1999	4	
57 audi	a4 quattro	1.8	1999	4	
58 audi	a4 quattro	1.8	1999	4	
59 audi	a4 quattro	1.8	1999	4	
60 audi	a4 quattro	1.8	1999	4	
61 audi	a4 quattro	1.8	1999	4	
62 audi	a4 quattro	1.8	1999	4	
63 audi	a4 quattro	1.8	1999	4	
64 audi	a4 quattro	1.8	1999	4	
65 audi	a4 quattro	1.8	1999	4	
66 audi	a4 quattro	1.8	1999	4	
67 audi	a4 quattro	1.8	1999	4	
68 audi	a4 quattro	1.8	1999	4	
69 audi	a4 quattro	1.8	1999	4	
70 audi	a4 quattro	1.8	1999	4	
71 audi	a4 quattro	1.8	1999	4	
72 audi	a4 quattro	1.8	1999	4	
73 audi	a4 quattro	1.8	1999	4	
74 audi	a4 quattro	1.8	1999	4	
75 audi	a4 quattro	1.8	1999	4	
76 audi	a4 quattro	1.8	1999	4	
77 audi	a4 quattro	1.8	1999	4	
78 audi	a4 quattro	1.8	1999	4	
79 audi	a4 quattro	1.8	1999	4	
80 audi	a4 quattro	1.8	1999	4	
81 audi	a4 quattro	1.8	1999	4	
82 audi	a4 quattro	1.8	1999	4	
83 audi	a4 quattro	1.8	1999	4	
84 audi	a4 quattro	1.8	1999	4	
85 audi	a4 quattro	1.8	1999	4	
86 audi	a4 quattro	1.8	1999	4	
87 audi	a4 quattro	1.8	1999	4	
88 audi	a4 quattro	1.8	1999	4	
89 audi	a4 quattro	1.8	1999	4	
90 audi	a4 quattro	1.8	1999	4	
91 audi	a4 quattro	1.8	1999	4	
92 audi	a4 quattro	1.8	1999	4	
93 audi	a4 quattro	1.8	1999	4	
94 audi	a4 quattro	1.8	1999	4	
95 audi	a4 quattro	1.8	1999	4	
96 audi	a4 quattro	1.8	1999	4	
97 audi	a4 quattro	1.8	1999	4	
98 audi	a4 quattro	1.8	1999	4	
99 audi	a4 quattro	1.8	1999	4	
100 audi	a4 quattro	1.8	1999	4	
101 audi	a4 quattro	1.8	1999	4	
102 audi	a4 quattro	1.8	1999	4	
103 audi	a4 quattro	1.8	1999	4	
104 audi	a4 quattro	1.8	1999	4	
105 audi	a4 quattro	1.8	1999	4	
106 audi	a4 quattro	1.8	1999	4	
107 audi	a4 quattro	1.8	1999	4	
108 audi	a4 quattro	1.8	1999	4	
109 audi	a4 quattro	1.8	1999	4	
110 audi	a4 quattro	1.8	1999	4	
111 audi	a4 quattro	1.8	1999	4	
112 audi	a4 quattro	1.8	1999	4	
113 audi	a4 quattro	1.8	1999	4	
114 audi	a4 quattro	1.8	1999	4	
115 audi	a4 quattro	1.8	1999	4	
116 audi	a4 quattro	1.8	1999	4	
117 audi	a4 quattro	1.8	1999	4	
118 audi	a4 quattro	1.8	1999	4	
119 audi	a4 quattro	1.8	1999	4	
120 audi	a4 quattro	1.8	1999	4	
121 audi	a4 quattro	1.8	1999	4	
122 audi	a4 quattro	1.8	1999	4	
123 audi	a4 quattro	1.8	1999	4	
124 audi	a4 quattro	1.8	1999	4	
125 audi	a4 quattro	1.8	1999	4	
126 audi	a4 quattro	1.8	1999	4	
127 audi	a4 quattro	1.8	1999	4	
128 audi	a4 quattro	1.8	1999	4	
129 audi	a4 quattro	1.8	1999	4	
130 audi	a4 quattro	1.8	1999	4	
131 audi	a4 quattro	1.8	1999	4	
132 audi	a4 quattro	1.8	1999	4	
133 audi	a4 quattro	1.8	1999	4	
134 audi	a4 quattro	1.8	1999	4	
135 audi	a4 quattro	1.8	1999	4	
136 audi	a4 quattro	1.8	1999	4	
137 audi	a4 quattro	1.8	1999	4	
138 audi	a4 quattro	1.8	1999	4	
139 audi	a4 quattro	1.8	1999	4	
140 audi	a4 quattro	1.8	1999	4	
141 audi	a4 quattro	1.8	1999	4	
142 audi	a4 quattro	1.8	1999	4	
143 audi	a4 quattro	1.8	1999	4	
144 audi	a4 quattro	1.8	1999	4	
145 audi	a4 quattro	1.8	1999	4	
146 audi	a4 quattro	1.8	1999	4	
147 audi	a4 quattro	1.8	1999	4	
148 audi	a4 quattro	1.8	1999	4	
149 audi	a4 quattro	1.8	1999	4	
150 audi	a4 quattro	1.8	1999	4	
151 audi	a4 quattro	1.8	1999	4	
152 audi	a4 quattro	1.8	1999	4	
153 audi	a4 quattro	1.8	1999	4	
154 audi	a4 quattro	1.8	1999	4	
155 audi	a4 quattro	1.8	1999	4	
156 audi	a4 quattro	1.8	1999	4	
157 audi	a4 quattro	1.8	1999	4	
158 audi	a4 quattro	1.8	1999	4	
159 audi	a4 quattro	1.8	1999	4	
160 audi	a4 quattro	1.8	1999	4	
161 audi	a4 quattro	1.8	1999	4	
162 audi	a4 quattro	1.8	1999	4	
163 audi	a4 quattro	1.8	1999	4	
164 audi	a4 quattro	1.8	1999	4	
165 audi	a4 quattro	1.8	1999	4	
166 audi	a4 quattro	1.8	1999	4	
167 audi	a4 quattro	1.8	1999	4	
168 audi	a4 quattro	1.8	1999	4	
169 audi	a4 quattro	1.8	1999	4	
170 audi	a4 quattro	1.8	1999	4	
171 audi	a4 quattro	1.8	1999	4	
172 audi	a4 quattro	1.8	1999		

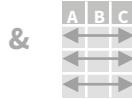
Data Transformation with dplyr :: CHEAT SHEET



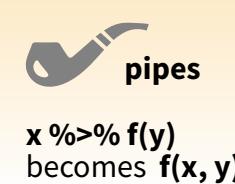
dplyr functions work with pipes and expect **tidy data**. In tidy data:



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



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



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).

	summary function
	summarise(.data, ...) Compute table of summaries. Also summarise_() . <code>summarise(mtcars, avg = mean(mpg))</code>
	count(x, ..., wt = NULL, sort = FALSE) Count number of rows in each group defined by the variables in ... Also tally() . <code>count(iris, Species)</code>

VARIATIONS

- summarise_all()** - Apply funs to every column.
- summarise_at()** - Apply funs to specific columns.
- summarise_if()** - Apply funs to all cols of one type.

Group Cases

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

	<code>mtcars %>% group_by(cyl) %>% summarise(avg = mean(mpg))</code>
--	--

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

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

Manipulate Cases

EXTRACT CASES

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

	filter(.data, ...) Extract rows that meet logical criteria. Also filter_() . <code>filter(iris, Sepal.Length > 7)</code>
	distinct(.data, ..., .keep_all = FALSE) Remove rows with duplicate values. Also distinct_() . <code>distinct(iris, Species)</code>
	sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select fraction of rows. <code>sample_frac(iris, 0.5, replace = TRUE)</code>
	sample_n(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select size rows. <code>sample_n(iris, 10, replace = TRUE)</code>
	slice(.data, ...) Select rows by position. Also slice_() . <code>slice(iris, 10:15)</code>
	top_n(x, n, wt) Select and order top n entries (by group if grouped data). <code>top_n(iris, 5, Sepal.Width)</code>

Logical and boolean operators to use with filter()

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

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

ARRANGE CASES

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

ADD CASES

	add_row(.data, ..., .before = NULL, .after = NULL) Add one or more rows to a table. <code>add_row(faithful, eruptions = 1, waiting = 1)</code>
--	---

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() . <code>select(iris, Sepal.Length, Species)</code>
--	---

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

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

MAKE NEW VARIABLES

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

vectorized function

	mutate(.data, ...) Compute new column(s). <code>mutate(mtcars, gpm = 1/mpg)</code>
--	---

	transmute(.data, ...) Compute new column(s), drop others. <code>transmute(mtcars, gpm = 1/mpg)</code>
--	--

	mutate_all(.tbl, .funs, ...) Apply funs to every column. Use with funs() . <code>mutate_all(faithful, funs(log(.), log2(.)))</code>
--	--

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

	mutate_if(.tbl, .predicate, .funs, ...) Apply funs to all columns of one type. Use with funs() . <code>mutate_if(iris, is.numeric, funs(log(.)))</code>
--	--

	add_column(.data, ..., .before = NULL, .after = NULL) Add new column(s). <code>add_column(mtcars, new = 1:32)</code>
--	--

	rename(.data, ...) Rename columns. <code>rename(iris, Length = Sepal.Length)</code>
--	---



Vector Functions

TO USE WITH MUTATE ()

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

vectorized function →

OFFSETS

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

CUMULATIVE AGGREGATES

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

RANKINGS

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

MATH

+, -, *, /, ^, %/%, %% - arithmetic ops
log(), **log2()**, **log10()** - logs
<, <=, >, >=, !=, == - logical comparisons

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

Summary Functions

TO USE WITH SUMMARISE ()

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

summary function →

COUNTS

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

LOCATION

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

LOGICALS

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

POSITION/ORDER

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

RANK

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

SPREAD

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

Row Names

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

A	B
1	a t
2	b u
3	c v

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

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

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

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

Combine Tables

COMBINE VARIABLES

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

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

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

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

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

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

A	B	C	D
a	t	1	3
b	u	2	2
d	w	NA	1

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

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

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

A	B	C	D
a	t	1	3
b	u	2	2
c	v	3	NA
d	w	NA	1

full_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ...)
Join data. Retain all values, all rows.

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

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

Use **suffix** to specify suffix to give to duplicate column names.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

COMBINE CASES

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

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

bind_rows(..., .id = NULL)
Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured)

intersect(x, y, ...)
Rows that appear in both x and y.

setdiff(x, y, ...)
Rows that appear in x but not y.

union(x, y, ...)
Rows that appear in x or y.
(Duplicates removed). **union_all()** retains duplicates.

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

EXTRACT ROWS

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

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

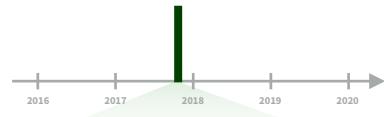
semi_join(x, y, by = NULL, ...)
Return rows of x that have a match in y.
USEFUL TO SEE WHAT WILL BE JOINED.

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

Dates and times with lubridate :: CHEAT SHEET



Date-times



2017-11-28 12:00:00

2017-11-28 12:00:00

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

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

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

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

2017-11-28T14:02:00

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

2017-22-12 10:00:00

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

11/28/2017 1:02:03

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

1 Jan 2017 23:59:59

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

20170131

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

July 4th, 2000

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

4th of July '99

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

2001: Q3

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

2:01

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

2017.5

date_decimal(decimal, tz = "UTC") Q for quarter. date_decimal(2017.5)



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

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

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

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

2017-11-28

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

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

12:00:00

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

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

GET AND SET COMPONENTS

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

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

2018-01-31 11:59:59

date(x) Date component. date(dt)

2018-01-31 11:59:59

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

2018-01-31 11:59:59

month(x, label, abbr) Month.

month(dt)

2018-01-31 11:59:59

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

2018-01-31 11:59:59

hour(x) Hour. hour(dt)

2018-01-31 11:59:59

minute(x) Minutes. minute(dt)

2018-01-31 11:59:59

second(x) Seconds. second(dt)

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

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

2018-01-31 11:59:59

leap_year(x) Is it a leap year?

leap_year(dt)

2018-01-31 11:59:59

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

Round Date-times



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

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

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

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

Tip: use a date with day > 12

Stamp Date-times

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

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

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

Time Zones

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

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

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

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

PT MT CT ET
7:00 Pacific 7:00 Mountain 7:00 Central

7:00 Pacific 7:00 Mountain 7:00 Central

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

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



Math with Date-times

— Lubridate provides three classes of timespans to facilitate math with dates and date-times

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

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

The start of daylight savings (spring forward)
`gap <- ymd_hms("2018-03-11 01:30:00", tz="US/Eastern")`

The end of daylight savings (fall back)
`lap <- ymd_hms("2018-11-04 00:30:00", tz="US/Eastern")`

Leap years and leap seconds
`leap <- ymd("2019-03-01")`

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

`normal + minutes(90)`

`gap + minutes(90)`

`lap + minutes(90)`

`leap + years(1)`

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

`normal + dminutes(90)`

`gap + dminutes(90)`

`lap + dminutes(90)`

`leap + dyears(1)`

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

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

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

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

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

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

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

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

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

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

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

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

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

PERIODS

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

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

```
p <- months(3) + days(12)
p
## "3m 12d 0H 0M 0S"
```

Number of months Number of days etc.

`years(x = 1) x years.`
`months(x) x months.`
`weeks(x = 1) x weeks.`
`days(x = 1) x days.`
`hours(x = 1) x hours.`
`minutes(x = 1) x minutes.`
`seconds(x = 1) x seconds.`
`milliseconds(x = 1) x milliseconds.`
`microseconds(x = 1) x microseconds`
`nanoseconds(x = 1) x milliseconds.`
`picoseconds(x = 1) x picoseconds.`

`period(num = NULL, units = "second", ...)`
An automation friendly period constructor.
`period(5, unit = "years")`

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

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

DURATIONS

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

Diftimes are a class of durations found in base R.

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

```
dd <- ddays(14)
dd
## "1209600s (~2 weeks)"
```

Exact length in seconds Equivalent in common units

`dyears(x = 1) 31536000x seconds.`
`dweeks(x = 1) 604800x seconds.`
`ddays(x = 1) 86400x seconds.`
`dhours(x = 1) 3600x seconds.`
`dminutes(x = 1) 60x seconds.`
`dseconds(x = 1) x seconds.`
`dmilliseconds(x = 1) x × 10⁻³ seconds.`
`dmicroseconds(x = 1) x × 10⁻⁶ seconds.`
`dnanoseconds(x = 1) x × 10⁻⁹ seconds.`
`dpicoseconds(x = 1) x × 10⁻¹² seconds.`

`duration(num = NULL, units = "second", ...)`
An automation friendly duration constructor. `duration(5, unit = "years")`

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

`make_difftime(x)` Make difftime with the specified number of units.
`make_difftime(99999)`

INTERVALS

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

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

```
i <- interval(ymd("2017-01-01"), d)
j <- d %--% ymd("2017-12-31")
## 2017-01-01 UTC--2017-11-28 UTC
## 2017-11-28 UTC--2017-12-31 UTC
```

Start Date End Date

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

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

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

`int_diff(times)` Make the intervals that occur between the date-times in a vector.
`v <- c(dt, dt + 100, dt + 1000)); int_diff(v)`

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

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

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

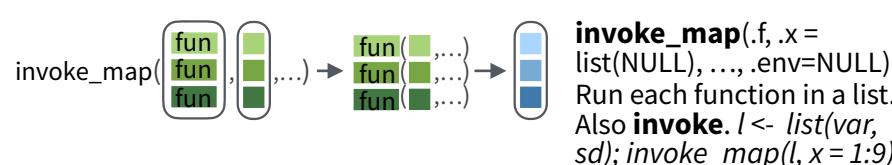
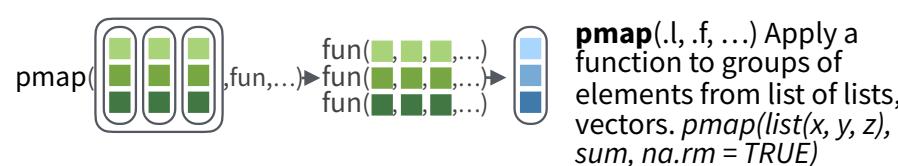
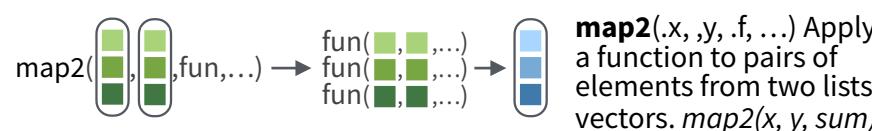
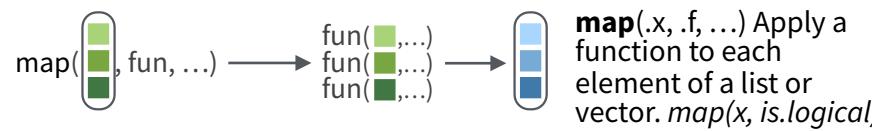
`as.interval(x, start, ...)` Coerce a timespans to an interval with the start date-time. Also `is.interval()`. `as.interval(days(1), start = now())`

Apply functions with purrr :: CHEAT SHEET



Apply Functions

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



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

OUTPUT

map(), **map2()**, **pmap()**, **imap** and **invoke_map** each return a list. Use a suffixed version to return the results as a specific type of flat vector, e.g. **map2_chr**, **pmap_lgl**, etc.

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

SHORTCUTS - within a purrr function:

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

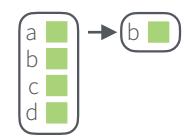
~.x becomes **function(x)x**, e.g. **map(l, ~.x + y)** becomes **map2(l, p, function(l, p) l + p)**

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

function	returns
map	list
map_chr	character vector
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
walk	triggers side effects, returns the input invisibly

Work with Lists

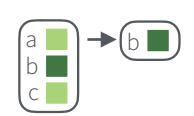
FILTER LISTS



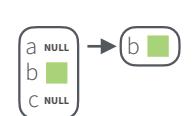
pluck(x, ..., .default=NULL) Select an element by name or index, **pluck(x, "b")**, or its attribute with **attr_getter**. **pluck(x, "b", attr_getter("n"))**



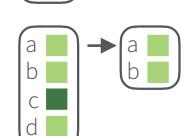
keep(x, .p, ...) Select elements that pass a logical test. **keep(x, is.character)**



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

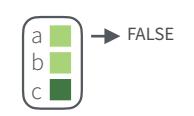


compact(x, .p = identity) Drop empty elements. **compact(x)**

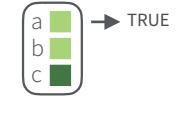


head_while(x, .p, ...) Return head elements until one does not pass. Also **tail_while**. **head_while(x, is.character)**

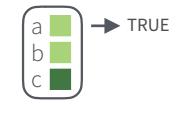
SUMMARISE LISTS



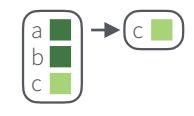
every(x, .p, ...) Do all element pass a test? **every(x, is.character)**



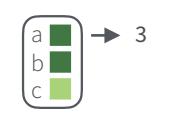
some(x, .p, ...) Do some elements pass a test? **some(x, is.character)**



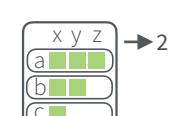
has_element(x, .y) Does a list contain an element? **has_element(x, "foo")**



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

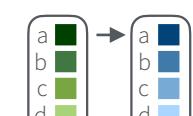


detect_index(x, .f, ..., .right=FALSE, .p) Find index of first element to pass. **detect_index(x, is.character)**

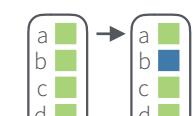


vec_depth(x) Return depth (number of levels of indexes). **vec_depth(x)**

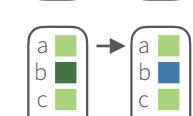
TRANSFORM LISTS



modify(x, .f, ...) Apply function to each element. Also **map**, **map_chr**, **map_dbl**, **map_dfc**, **map_dfr**, **map_int**, **map_lgl**. **modify(x, ~.+2)**



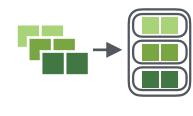
modify_at(x, .at, .f, ...) Apply function to elements by name or index. Also **map_at**. **modify_at(x, "b", ~.+2)**



modify_if(x, .p, .f, ...) Apply function to elements that pass a test. Also **map_if**. **modify_if(x, is.numeric, ~.+2)**

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

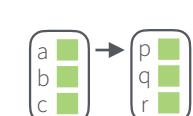
WORK WITH LISTS



array_tree(array, margin = NULL) Turn array into list. Also **array_branch**. **array_tree(x, margin = 3)**



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

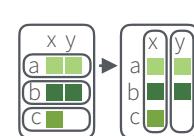


set_names(x, nm = x) Set the names of a vector/list directly or with a function. **set_names(x, c("p", "q", "r"))** **set_names(x, tolower)**

RESHAPE LISTS

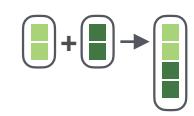


flatten(x) Remove a level of indexes from a list. Also **flatten_chr**, **flatten_dbl**, **flatten_dfc**, **flatten_dfr**, **flatten_int**, **flatten_lgl**. **flatten(x)**

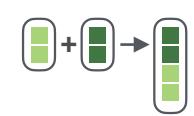


transpose(x, .names = NULL) Transposes the index order in a multi-level list. **transpose(x)**

JOIN (TO) LISTS



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

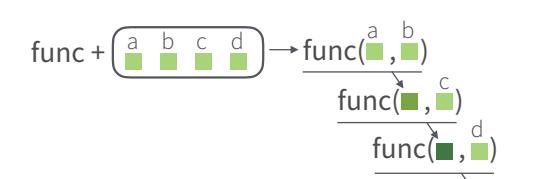


prepend(x, values, before = 1) Add to start of list. **prepend(x, list(d = 1))**

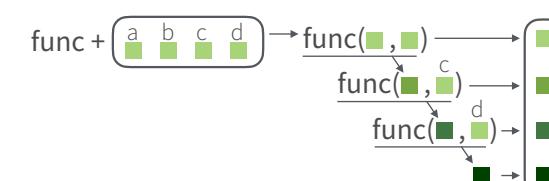


splice(...) Combine objects into a list, storing S3 objects as sub-lists. **splice(x, y, "foo")**

Reduce Lists



reduce(x, .f, ..., .init) Apply function recursively to each element of a list or vector. Also **reduce_right**, **reduce2**, **reduce2_right**. **reduce(x, sum)**



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

compose() Compose multiple functions.

lift() Change the type of input a function takes. Also **lift_dl**, **lift_lv**, **lift_vl**.

rerun() Rerun expression n times.

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

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

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

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

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



Nested Data

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

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

`n_iris$data[[1]]`

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

`n_iris`

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

`n_iris$data[[2]]`

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

`n_iris$data[[3]]`

Use a nested data frame to:

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

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

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

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

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

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

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

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

Species	data	Species	S.L	S.W	P.L	P.W
setos	<tibble [50x4]>	setosa	5.1	3.5	1.4	0.2
versi	<tibble [50x4]>	setosa	4.9	3.0	1.4	0.2
virgini	<tibble [50x4]>	setosa	4.7	3.2	1.3	0.2

`n_iris %>% unnest()`

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

Unnests a nested data frame.

List Column Workflow

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

1 Make a list column

Species	S.L	S.W	P.L	P.W
setosa	5.1	3.5	1.4	0.2
setosa	4.9	3.0	1.4	0.2
setosa	4.7	3.2	1.3	0.2
setosa	4.6	3.1	1.5	0.2

Species	data
setosa	<tibble [50x4]>
versi	<tibble [50x4]>
virgini	<tibble [50x4]>

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

2 Work with list columns

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

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

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

3 Simplify the list column

Species	beta
setos	2.35
versi	1.89
virgini	0.69

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

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

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

`tibble::tribble(...)`

Makes list column when needed

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

`tibble::tibble(...)`

Saves list input as list columns

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

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

Converts multi-level list to tibble with list cols
`enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')`

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

Returns list col when result returns list.

`mtcars %>% group_by(cyl) %>%`

`summarise(q = list(quantile(mpg)))`

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.



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

Subset Strings

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

Manage Lengths

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

Mutate Strings

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

Join and Split

	<code>str_c(..., sep = "", collapse = NULL)</code> Join multiple strings into a single string. <code>str_c(letters, LETTERS)</code>
	<code>str_c(..., sep = "", collapse = NULL)</code> Collapse a vector of strings into a single string. <code>str_c(letters, collapse = "")</code>
	<code>str_dup(string, times)</code> Repeat strings times times. <code>str_dup(fruit, times = 2)</code>
	<code>str_split_fixed(string, pattern, n)</code> Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also <code>str_split</code> to return a list of substrings. <code>str_split_fixed(fruit, " ", n=2)</code>
	<code>glue::glue(..., .sep = "", .envir = parent.frame(), .open = "{", .close = "}")</code> Create a string from strings and {expressions} to evaluate. <code>glue::glue("Pi is {pi}")</code>
	<code>glue::glue_data(.x, ..., .sep = "", .envir = parent.frame(), .open = "{", .close = "}")</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>

Order Strings

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

Helpers

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

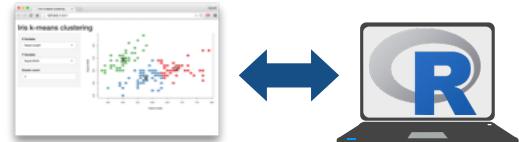
¹ See bit.ly/ISO639-1 for a complete list of locales.



Shiny :: CHEAT SHEET

Basics

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



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

APP TEMPLATE

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

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

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

SHARE YOUR APP

 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/



Building an App

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

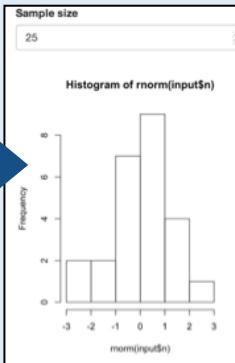
Add inputs to the UI with `*Input()` functions

Add outputs with `*Output()` functions

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

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

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



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

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

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

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

ui.R contains everything you would save to ui.

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

No need to call **shinyApp()**.

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

● ● ● **app-name**
 .r
 app.R
 global.R
 DESCRIPTION
 README
 <other files>
 www

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

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

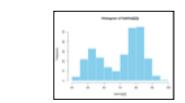
Outputs - `render*`() and `*Output()` functions work together to add R output to the UI



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



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



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



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

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

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

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

works with

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

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

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

`verbatimTextOutput(outputId)`

`tableOutput(outputId)`

`textOutput(outputId, container, inline)`

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

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

Inputs

collect values from the user

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

Action

Link

- Choice 1
 Choice 2
 Choice 3
 Check me

`checkboxInput(inputId, label, value)`

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

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

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

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

`passwordInput(inputId, label, value)`

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

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

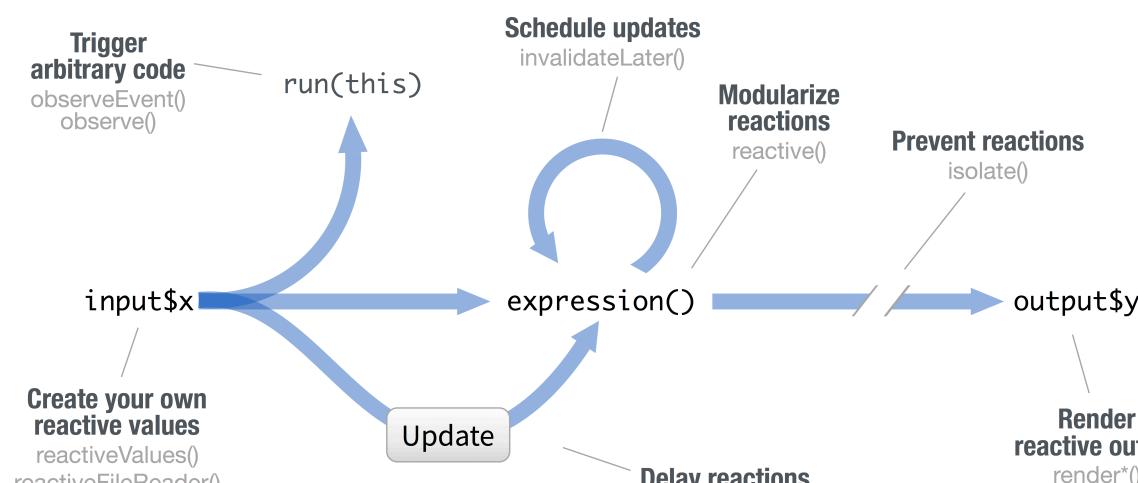
`sliderInput(inputId, label, min, max, value, step, round, format, locale, ticks, animate, width, sep, pre, post)`

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

`textInput(inputId, label, value)`

Reactivity

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



CREATE YOUR OWN REACTIVE VALUES

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

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

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

PREVENT REACTIONS

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

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

shinyApp(ui, server)
```

isolate(expr)
Runs a code block. Returns a **non-reactive** copy of the results.

RENDERS REACTIVE OUTPUT

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

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

shinyApp(ui, server)
```

render*() functions
(see front page)

Builds an object to display. Will rerun code in body to rebuild the object whenever a reactive value in the code changes.
Save the results to **output\$<outputId>**

TRIGGER ARBITRARY CODE

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

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

shinyApp(ui, server)
```

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

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

MODULARIZE REACTIONS

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

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

shinyApp(ui, server)
```

reactive(x, env, quoted, label, domain)
Creates a **reactive expression** that

- **caches** its value to reduce computation
- can be called by other code
- notifies its dependencies when it has been invalidated

Call the expression with function syntax, e.g. **re()**

DELAY REACTIONS

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

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

shinyApp(ui, server)
```

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

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

UI

An app's UI is an HTML document.

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

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



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

tags\$a	tags\$data	tags\$h6	tags\$nav	tags\$span
tags\$abbr	tags\$datalist	tags\$head	tags\$noscript	tags\$strong
tags\$address	tags\$dd	tags\$header	tags\$object	tags\$style
tags\$area	tags\$del	tags\$hgroup	tags\$ol	tags\$sub
tags\$article	tags\$details	tags\$hrg	tags\$optgroup	tags\$summary
tags\$aside	tags\$dfn	tags\$HTML	tags\$option	tags\$sup
tags\$audio	tags\$div	tags\$si	tags\$output	tags\$table
tags\$b	tags\$dl	tags\$iframe	tags\$p	tags\$tbody
tags\$base	tags\$dt	tags\$img	tags\$param	tags\$td
tags\$bdi	tags\$em	tags\$input	tags\$pre	tags\$textarea
tags\$bdo	tags\$embed	tags\$ins	tags\$progress	tags\$tfoot
tags\$blockquote	tags\$eventsouce	tags\$kbd	tags\$q	tags\$th
tags\$body	tags\$fieldset	tags\$keygen	tags\$ruby	tags\$thead
tags\$br	tags\$figcaption	tags\$label	tags\$rp	tags\$time
tags\$button	tags\$figure	tags\$legend	tags\$rt	tags\$title
tags\$canvas	tags\$footer	tags\$li	tags\$ss	tags\$track
tags\$caption	tags\$form	tags\$link	tags\$small	tags\$u
tags\$cite	tags\$h1	tags\$mark	tags\$meta	tags\$ul
tags\$code	tags\$h2	tags\$map	tags\$select	tags\$var
tags\$col	tags\$h3	tags\$menu	tags\$small	tags\$video
tags\$colgroup	tags\$h4	tags\$meta	tags\$source	tags\$wbr
tags\$command	tags\$h5	tags\$script		

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

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

Header 1

bold
italic
code
link
Raw html

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

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

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

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

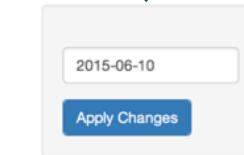
IMAGES To include an image 1. Place the file in the **www** subdirectory 2. Link to it with **img(src=<file name>")**

Layouts



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

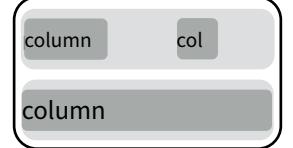
```
wellPanel(dateInput("a", ""),
  submitButton())
)
```



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

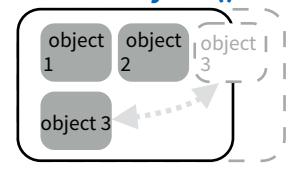
Organize panels and elements into a layout with a layout function. Add elements as arguments of the layout functions.

fluidRow()



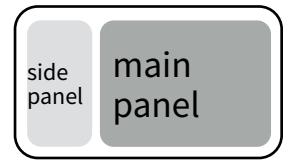
```
ui <- fluidPage(
  fluidRow(column(width = 4),
    column(width = 2, offset = 3)),
  fluidRow(column(width = 12))
)
```

flowLayout()



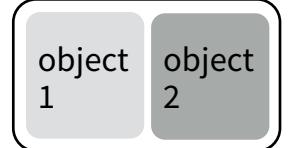
```
ui <- fluidPage(
  flowLayout(# object 1,
            # object 2,
            # object 3
  )
)
```

sidebarLayout()



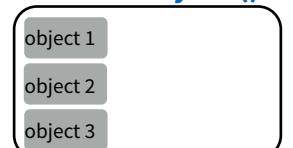
```
ui <- fluidPage(
  sidebarLayout(
    sidebarPanel(),
    mainPanel()
  )
)
```

splitLayout()

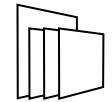


```
ui <- fluidPage(
  splitLayout(# object 1,
             # object 2
  )
)
```

verticalLayout()



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



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

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



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



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

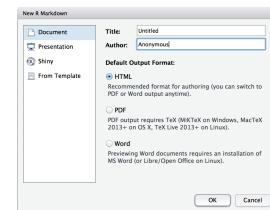


R Markdown :: CHEAT SHEET

What is R Markdown?

- .Rmd files** • An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.
- Reproducible Research** • At the click of a button, or the type of a command, you can rerun the code in an R Markdown file to reproduce your work and export the results as a finished report.
- Dynamic Documents** • You can choose to export the finished report in a variety of formats, including html, pdf, MS Word, or RTF documents; html or pdf based slides, Notebooks, and more.

Workflow



- ① Open a new .Rmd file at File ► New File ► R Markdown. Use the wizard that opens to pre-populate the file with a template
- ② Write document by editing template
- ③ Knit document to create report; use knit button or render() to knit
- ④ Preview Output in IDE window
- ⑤ Publish (optional) to web server
- ⑥ Examine build log in R Markdown console
- ⑦ Use output file that is saved along side .Rmd

The screenshot illustrates the R Markdown workflow within RStudio. It shows the R Markdown editor with code like `summary(cars)` and its resulting output. The preview window shows the rendered HTML page. The file browser shows the generated `report.html` file. The terminal window shows the command used to render the document.

render

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

input - file to render
output_format

output_options - List of render options (as in YAML)

output_file
output_dir

params - list of params to use

envir - environment to evaluate code chunks in

encoding - of input file

Embed code with knitr syntax

INLINE CODE

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

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

CODE CHUNKS

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

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

GLOBAL OPTIONS

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

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

IMPORTANT CHUNK OPTIONS

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

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

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

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

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

dependson - chunk dependencies for caching (default = NULL)

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

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

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

eval - Run code in chunk (default = TRUE)

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

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

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

highlight - highlight source code (default = TRUE)

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

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

results (default = 'markup')

'asis' - passthrough results

'hide' - do not display results

'hold' - put all results below all code

tidy - tidy code for display (default = FALSE)

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

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



.rmd Structure

YAML Header

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

At start of file

Between lines of ---

Text

Narration formatted with markdown, mixed with:

Code Chunks

Chunks of embedded code. Each chunk:

Begins with `{{r}}`

ends with `{{ }}`

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

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

Parameters

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

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

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

2. **Call parameters** • Call parameter values in code as `params$<name>`

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

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


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

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

Interactive Documents

Turn your report into an interactive Shiny document in 4 steps

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

```
---
```

`output: html_document`
`runtime: shiny`

```
---
```

````{r, echo = FALSE}`

`numericInput("n",`

`"How many cars?", 5)`

`renderTable({`

`head(cars, input$n)`

`)`

How many cars?  
5  

speed	dist
4.00	2.00
4.00	10.00
7.00	4.00
7.00	22.00
8.00	16.00

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

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



# Pandoc's Markdown

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

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

```
equation block:
```

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

```
> block quote
```

```
Header1 {#anchor}
```

```
Header 2 {#css_id}
```

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

```
Header 4
```

```
Header 5
```

```
Header 6
```

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

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

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

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

```
E = mc2
```

```
block quote
```

## Header1

## Header 2

### Header 3

#### Header 4

##### Header 5

##### Header 6

HTML ignored in pdfs

http://www.rstudio.com

link

Jump to Header 1

image:



Caption

- unordered list
  - sub-item 1
  - sub-item 2
  - sub-sub-item 1
- item 2

Continued (indent 4 spaces)

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

(@) A list whose numbering

continues after

2. an interruption

Term 1

Definition 1

Right	Left	Default	Center
12	12	12	12
123	123	123	123
1	1	1	1

- slide bullet 1
- slide bullet 2

(>- to have bullets appear on click)

horizontal rule/slide break:

\*\*\*

A footnote<sup>[1]</sup>

[^1]: Here is the footnote.

- (>- to have bullets appear on click)
- horizontal rule/slide break:
- 
1. Here is the footnote.<sup>1</sup>

# Set render options with YAML

When you render, R Markdown

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



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

```

output: html_document

Body
```

**output value**

**creates**

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

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

```

output: html_document:
 code_folding: hide
 toc_float: TRUE

Body
```

**html tabs**

Use tablet css class to place sub-headers into tabs

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



## Create a Reusable Template

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

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

**template.yaml** (see below)

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

any supporting files

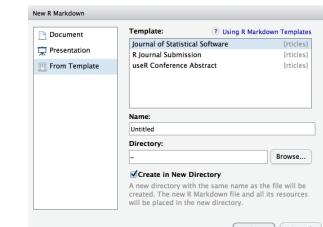
3. **Install the package**

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

```

name: My Template

```



**sub-option**      **description**

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

## Table Suggestions

Several functions format R data into tables

Table with kable	
<code>eruptions waiting</code>	
1	3.60 79.00
2	1.80 54.00
3	3.33 74.00
4	2.28 62.00
2.283	62

eruptions waiting	
1	3.60 79
2	1.80 54
3	3.33 74
4	

# Data Science in Spark with sparklyr :: CHEAT SHEET

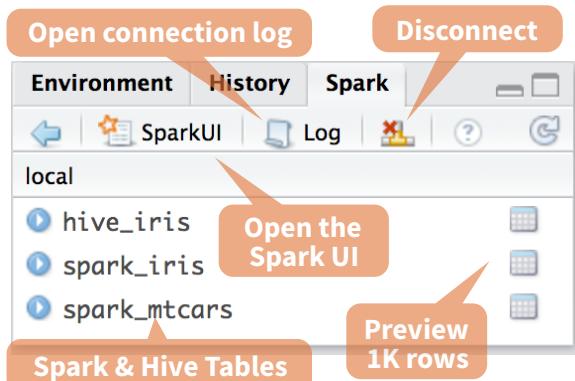


## Intro

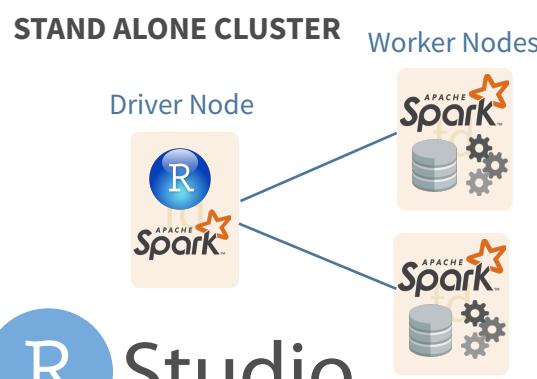
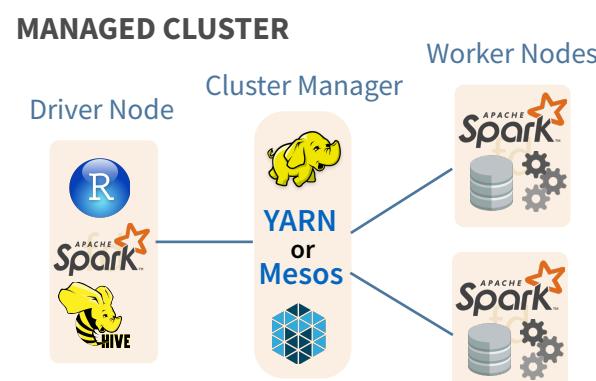
**sparklyr** is an R interface for Apache Spark™, it provides a complete **dplyr** backend and the option to query directly using **Spark SQL** statement. With sparklyr, you can orchestrate distributed machine learning using either **Spark's MLLib** or **H2O** Sparkling Water.

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

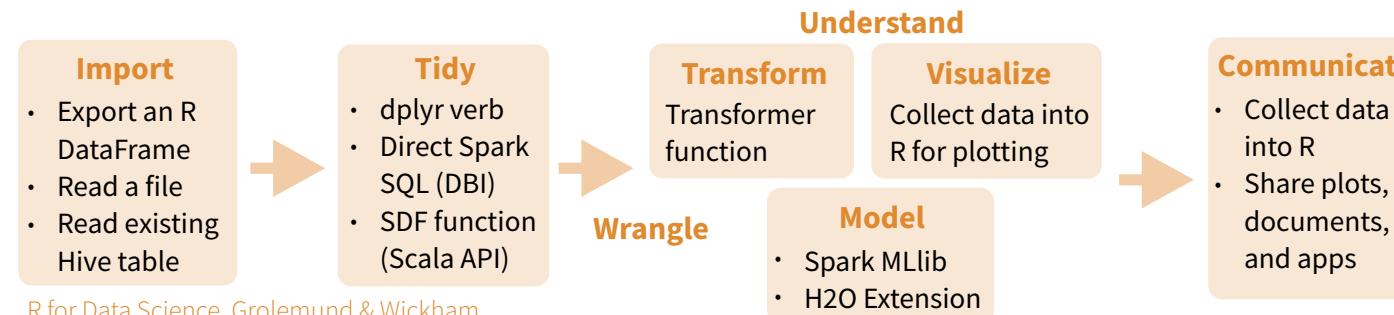
### RStudio Integrates with sparklyr



## Cluster Deployment



## Data Science Toolchain with Spark + sparklyr



## Getting Started

### LOCAL MODE (No cluster required)

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

### ON A MESOS MANAGED CLUSTER

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

### ON A YARN MANAGED CLUSTER

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

### ON A SPARK STANDALONE CLUSTER

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

### USING LIVY (Experimental)

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

## Tuning Spark

### EXAMPLE CONFIGURATION

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

### IMPORTANT TUNING PARAMETERS with defaults

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

## Using sparklyr

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

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

**Install Spark locally**

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

**Connect to local version**

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

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

**Copy data to Spark memory**

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

**Partition data**

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

**Create a hive metadata for each partition**

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

**Spark ML Decision Tree Model**

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

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

**Create reference to Spark table**

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

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

**Bring data back into R memory for plotting**

```
spark_disconnect(sc)
```

**Disconnect**



# Reactivity

## COPY A DATA FRAME INTO SPARK

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

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

## IMPORT INTO SPARK FROM A FILE

Arguments that apply to all functions:

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

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

**JSON** `spark_read_json()`

**PARQUET** `spark_read_parquet()`

## SPARK SQL COMMANDS

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

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

## FROM A TABLE IN HIVE

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

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

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

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

# Wrangle

## SPARK SQL VIA DPLYR VERBS

Translates into Spark SQL statements

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

## DIRECT SPARK SQL COMMANDS

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

`DBI::dbGetQuery(conn, statement)`

## SCALA API VIA SDF FUNCTIONS

`sdf_mutate(.data)`

Works like dplyr mutate function

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

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

`sdf_register(x, name = NULL)`

Gives a Spark DataFrame a table name

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

`sdf_sort(x, columns)`  
Sorts by >=1 columns in ascending order

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

`sdf_predict(object, newdata)`  
Spark DataFrame with predicted values

## ML TRANSFORMERS

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

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

`ft_binarizer(threshold = 0.5)`  
Assigned values based on threshold

`ft_bucketizer(splits)`  
Numeric column to discretized column

`ft_discrete_cosine_transform(inverse = FALSE)`  
Time domain to frequency domain

`ft_elementwise_product(scaling.col)`  
Element-wise product between 2 cols

`ft_index_to_string()`  
Index labels back to label as strings

`ft_one_hot_encoder()`  
Continuous to binary vectors

`ft_quantile_discretizer(n.buckets = 5L)`  
Continuous to binned categorical values

`ft_sql_transformer(sql)`  
Column of labels into a column of label indices.

`ft_string_indexer(params = NULL)`  
Combine vectors into single row-vector

# Visualize & Communicate

## DOWNLOAD DATA TO R MEMORY

`r_table <- collect(my_table)`  
`plot(Petal_Width ~ Petal_Length, data = r_table)`

`dplyr::collect(x)`  
Download a Spark DataFrame to an R DataFrame

`sdf_read_column(x, column)`

Returns contents of a single column to R

## SAVE FROM SPARK TO FILE SYSTEM

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

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

**JSON** `spark_read_json(mode = NULL)`

**PARQUET** `spark_read_parquet(mode = NULL)`

# Model (MLlib)

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

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

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

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

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

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

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

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

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

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

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

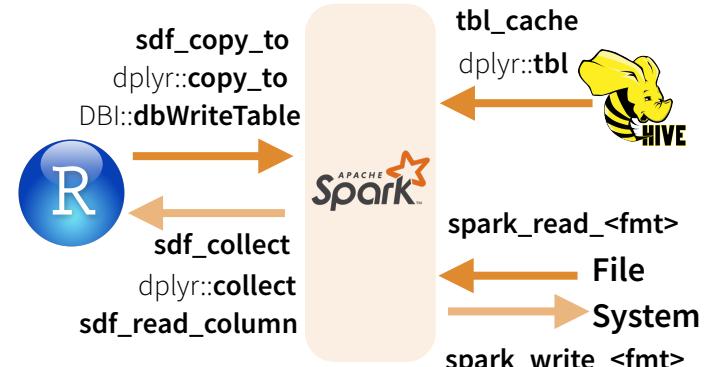
`ml_random_forest(x, response, features, max.bins = 32L, max.depth = 5L, num.trees = 20L, type = c("auto", "regression", "classification"), ml.options = ml_options())`

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

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

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

`ml_tree_feature_importance(sc, model)`



# Extensions

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

## CORE TYPES

`spark_connection()` Connection between R and the Spark shell process

`spark_jobj()` Instance of a remote Spark object

`spark_dataframe()` Instance of a remote Spark DataFrame object

## CALL SPARK FROM R

`invoke()` Call a method on a Java object

`invoke_new()` Create a new object by invoking a constructor

`invoke_static()` Call a static method on an object

## MACHINE LEARNING EXTENSIONS

`ml_create_dummy_variables()` `ml_options()`

`ml_prepare_dataframe()` `ml_model()`

`ml_prepare_response_features_intercept()`

## sparklyr

is an R interface  
for

