

# 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.

table(x)

See counts of values.

rev(x)

Return x reversed.

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  
}  
else {  
  Do something different  
}
```

#### Example

```
if (i > 3){  
  print('Yes')  
} else {  
  print('No')  
}
```

#### If Statements

```
if (condition){  
  Do something  
}  
else {  
  Do something different  
}
```

#### 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.

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

## Maths Functions

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

## Variable Assignment

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

## The Environment

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

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

## Matrices

`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: $m^*x = n$

## 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>	

`nrow(df)`

Number of rows.

`ncol(df)`

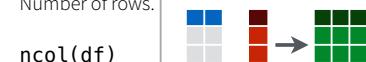
Number of columns.

`dim(df)`

Number of columns and rows.

`cbind` - Bind columns.

`rbind` - Bind rows.



## Strings

Also see the `string` package.

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

Join multiple vectors together.

`paste(x, collapse = ' ')`

Join elements of a vector together.

`grep(pattern, x)`

Find regular expression matches in x.

`gsub(pattern, replace, x)`

Replace matches in x with a string.

`toupper(x)`

Convert to uppercase.

`tolower(x)`

Convert to lowercase.

`nchar(x)`

Number of characters in a string.

## Factors

`factor(x)`

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

`cut(x, breaks = 4)`

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

## Statistics

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

Linear model.

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

Generalised linear model.

`summary`

Get more detailed information out a model.

`t.test(x, y)`

Perform a t-test for difference between means.

`pairwise.t.test`

Perform a t-test for paired data.

`prop.test`

Test for a difference between proportions.

`aov`

Analysis of variance.

## Distributions

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

## Plotting

Also see the `ggplot2` package.

`plot(x)`

Values of x in order.

`plot(x, y)`

Values of x against y.

`hist(x)`

Histogram of x.

## Dates

See the `lubridate` package.

# 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')`

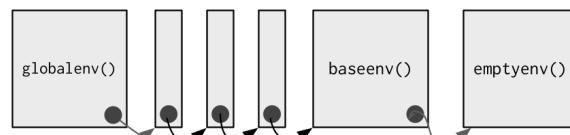


Figure 1 – The Search Path

- 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

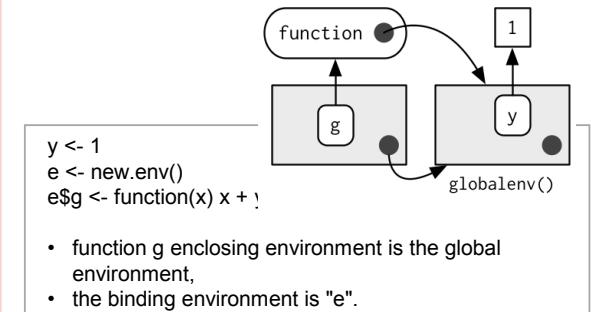
1. **Enclosing environment** - an environment where the function is created. It determines how function finds value.

- 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

- Mode:** represents how an object is stored in memory
  - 'type' of the object from R's point of view
  - Access with: **typeof()**
- 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

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

```
class(x) -> 'factor'
```

```
levels(x) # defines the set of allowed values
```

- 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:

- Factors look and often behave like character vectors, they are actually integers. Be careful when treating them like strings.
- 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 :

- 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#

- 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.

- 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) &amp; !isS4(x) or pryr::obj_type()</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 **simpliest** possible data structure that can represent the output

### 2. Preserving subsetting

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

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

Simplifying behavior varies slightly between different data types:

### 1. Atomic Vector

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

### 2. List

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

### 3. Factor

- Drops any unused levels but it remains a factor class

### 4. Matrix or Array

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

### 5. Data Frame

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

## Data Frame Subsetting

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

### 1. Subset with a single vector : Behave like lists

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

### 2. Subset with two vectors : Behave like matrices

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

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

### 1. Behave like matrices

```
str(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	<code>df1\$col3 &lt;- NULL</code>
Subset to return only columns you want	<code>df1[c('col1', 'col2')]</code>

### 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)

- Using integer subsetting is more effective when:

- You want to find the first (or last) TRUE.
- You have very few TRUEs and very many FALSEs; a set representation may be faster and require less storage.

- which() - conversion from boolean representation to integer representation

```
which(c(T, F, T F)) -> 1 3
```

- Integer representation length : is always  $\leq$  boolean representation length
- Common mistakes :
  - Use `x[which(y)]` instead of `x[y]`
  - `x[-which(y)]` is not equivalent to `x[!y]`

#### Recommendation:

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

## Subsetting with Assignment

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

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

- 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

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

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

- Avoid [ 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`

# R For Data Science Cheat Sheet

## data.table

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



### data.table

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



Load the package:

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

### Creating A data.table

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

Create a `data.table` and call it `DT`

### Subsetting Rows Using i

<code>&gt; DT[3:5]</code>	Select 3rd to 5th row
<code>&gt; DT[3:5]</code>	Select 3rd to 5th row
<code>&gt; DT[V2=="A"]</code>	Select all rows that have value A in column V2
<code>&gt; DT[V2 %in% c("A", "C")]</code>	Select all rows that have value A or C in column V2

### Manipulating on Columns in j

<code>&gt; DT[, V2]</code> [1] "A" "B" "C" "A" "B" "C" ... <code>&gt; DT[, .(V2, V3)]</code>	Return <code>V2</code> as a vector
<code>&gt; DT[, sum(V1)]</code>	Return <code>V2</code> and <code>V3</code> as a <code>data.table</code>
<code>&gt; DT[, .(sum(V1), sd(V3))]</code>	Return the sum of all elements of <code>V1</code> in a vector
<code>v1 v2</code> 1: 18 0.4546055	Return the sum of all elements of <code>V1</code> and the std. dev. of <code>V3</code> in a <code>data.table</code>
<code>&gt; DT[, .(Aggregate=sum(V1), Sd.V3=sd(V3))]</code>	The same as the above, with new names
<code>1: 18 0.4546055</code>	Select column <code>V2</code> and compute std. dev. of <code>V3</code> , which returns a single value and gets recycled
<code>&gt; DT[, .(print(V2), plot(V3), NULL)]</code>	Print column <code>V2</code> and plot <code>V3</code>

### Doing j by Group

<code>&gt; DT[, .(V4.Sum=sum(V4)), by=V1]</code>	Calculate sum of <code>V4</code> for every group in <code>V1</code>
<code>v1 V4.Sum</code> 1: 1 36 2: 2 42	Calculate sum of <code>V4</code> for every group in <code>V1</code> and <code>V2</code>
<code>&gt; DT[, .(V4.Sum=sum(V4)), by=(V1, V2)]</code>	Calculate sum of <code>V4</code> for every group in <code>sign(V1-1)</code>
<code>&gt; DT[, .(V4.Sum=sum(V4)), by=sign(V1-1)]</code>	The same as the above, with new name for the variable you're grouping by
<code>sign V4.Sum</code> 1: 0 36 2: 1 42	Calculate sum of <code>V4</code> for every group in <code>V1</code> after subsetting on the first 5 rows
<code>&gt; DT[1:5, .(V4.Sum=sum(V4)), by=V1]</code>	Count number of rows for every group in <code>V1</code>

### General form: `DT[i, j, by]`

"Take `DT`, subset rows using `i`, then calculate `j` grouped by `by`"

### Adding/Updating Columns By Reference in j Using :=

```
> DT[, V1:=round(exp(V1), 2)]
> DT
  V1 V2   V3 V4
1: 2.72 A -0.1107 1
2: 7.39 B -0.1427 2
3: 2.72 C -1.8893 3
4: 7.39 A -0.3571 4
...
> DT[, c("V1", "V2"):=list(round(exp(V1), 2),
LETTERS[4:6])]
> DT[, ' :='(V1:=round(exp(V1), 2),
V2=LETTERS[4:6])][]
  V1 V2   V3 V4
1: 15.18 D -0.1107 1
2: 1619.71 E -0.1427 2
3: 15.18 F -1.8893 3
4: 1619.71 D -0.3571 4
> DT[, V1:=NULL]
> DT[, c("V1", "V2"):=NULL]
> Cols.chosen=c("A", "B")
> DT[, Cols.Chosen:=NULL]
> DT[, (Cols.Chosen) :=NULL]
```

`V1` is updated by what is after `:=`  
Return the result by calling `DT`

Columns `V1` and `V2` are updated by what is after `:=`  
Alternative to the above one. With `[, ]`, you print the result to the screen

Remove `V1`  
Delete the column with column name `Cols.chosen`  
Delete the columns specified in the variable `Cols.chosen`

### Advanced Data Table Operations

<code>&gt; DT[., N-1]</code>	Return the penultimate row of the <code>DT</code>
<code>&gt; DT[., N]</code>	Return the number of rows
<code>&gt; DT[., .(V2, V3)]</code>	Return <code>V2</code> and <code>V3</code> as a <code>data.table</code>
<code>&gt; DT[, list(V2, V3)]</code>	Return <code>V2</code> and <code>V3</code> as a <code>data.table</code>
<code>&gt; DT[, mean(V3), by=.(V1, V2)]</code>	Return the result of <code>j</code> , grouped by all possible combinations of groups specified in <code>by</code>
<code>v1 V2   V1</code> 1: 1 A 0.4053 2: 1 B 0.4053 3: 1 C 0.4053 4: 2 A -0.6443 5: 2 B -0.6443 6: 2 C -0.6443	

### .SD & .SDcols

<code>&gt; DT[, print(.SD), by=V2]</code>	Look at what <code>.SD</code> contains
<code>&gt; DT[, .SD[c(1, .N)], by=V2]</code>	Select the first and last row grouped by <code>V2</code>
<code>&gt; DT[, lapply(.SD, sum), by=V2]</code>	Calculate sum of columns in <code>.SD</code> grouped by <code>V2</code>
<code>&gt; DT[, lapply(.SD, sum), by=V2, .SDcols=c("V3", "V4")]</code>	Calculate sum of <code>V3</code> and <code>V4</code> in <code>.SD</code> grouped by <code>V2</code>
<code>v2   V3 V4</code> 1: A -0.478 22 2: B -0.478 26 3: C -0.478 30	Calculate sum of <code>V3</code> and <code>V4</code> in <code>.SD</code> grouped by <code>V2</code>

### Chaining

<code>&gt; DT &lt;- DT[., .(V4.Sum=sum(V4)), by=V1]</code>	Calculate sum of <code>V4</code> , grouped by <code>V1</code>
<code>v1 V4.Sum</code> 1: 1 36 2: 2 42	Select that group of which the sum is >40
<code>&gt; DT[V4.Sum&gt;40]</code>	Select that group of which the sum is >40 (chaining)
<code>&gt; DT[., (V4.Sum=sum(V4)), by=V1][V4.Sum&gt;40]</code>	Calculate sum of <code>V4</code> , grouped by <code>V1</code> , ordered on <code>V1</code>
<code>v1 V4.Sum</code> 1: 2 42 2: 1 36	

### set() -Family

#### set()

Syntax: `for (i in from:to) set(DT, row, column, new value)`

<code>&gt; rows &lt;- list(3:4, 5:6)</code>	
<code>&gt; cols &lt;- 1:2</code>	
<code>&gt; for(i in seq_along(rows))</code>	Sequence along the values of <code>rows</code> , and for the values of <code>cols</code> , set the values of those elements equal to <code>NA</code> (invisible)
<code>  {set(DT,</code>	
<code>    i\$rows[[i]],</code>	
<code>    j=cols[i],</code>	
<code>    value=NA)}</code>	

#### setnames()

Syntax: `setnames(DT, "old", "new") []`

<code>&gt; setnames(DT, "V2", "Rating")</code>	Set name of <code>V2</code> to <code>Rating</code> (invisible)
<code>&gt; setnames(DT,</code>	Change 2 column names (invisible)
<code>  c("V2", "V3"),</code>	
<code>  c("V2.rating", "V3.DC"))</code>	

#### setnames()

Syntax: `setcolorder(DT, "neworder")`

<code>&gt; setcolorder(DT,</code>	Change column ordering to contents
<code>  c("V2", "V1", "V4", "V3"))</code>	of the specified vector (invisible)

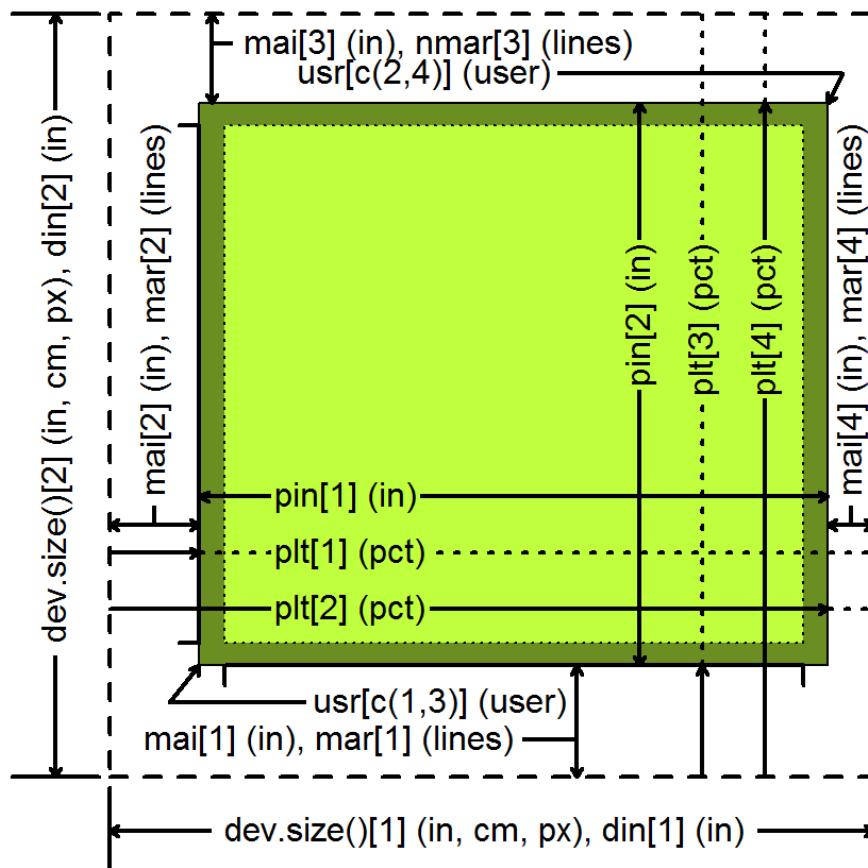


# How Big is Your Graph?

## An R Cheat Sheet

### Introduction

All functions that open a device for graphics will have **height** and **width** arguments to control the size of the graph and a **pointsize** argument to control the relative font size. In **knitr**, you control the size of the graph with the chunk options, **fig.width** and **fig.height**. This sheet will help you with calculating the size of the graph and various parts of the graph within R.



### Your graphics device

**dev.size()** (width, height)  
**par("din")** (r.o.) (width, height) in inches

Both the **dev.size** function and the **din** argument of **par** will tell you the size of the graphics device. The **dev.size** function will report the size in

1. inches (**units="in"**), the default
2. centimeters (**units="cm"**)
3. pixels (**units="px"**)

Like several other **par** arguments, **din** is read only (r.o.) meaning that you can ask its current value (**par("din")**) but you cannot change it (**par(din=c(5,7)** will fail).

### Your plot margins

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

Margins provide you space for your axes, axis, labels, and titles.

A "line" is the amount of vertical space needed for a line of text.

If your graph has no axes or titles, you can remove the margins (and maximize the plotting region) with

**par(mar=rep(0,4))**

### Your plotting region

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

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

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

### Your x-y coordinates

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

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

Run **par("usr")** to find the minimum X value, the maximum X value, the minimum Y value, and the maximum Y value. If you assign new values to **usr**, you will update the x-y coordinates to the new values.

### Getting a square graph

**par("pty")**

You can produce a square graph manually by setting the width and height to the same value and setting the margins so that the sum of the top and bottom margins equal the sum of the left and right margins. But a much easier way is to specify **pty="s"**, which adjusts the margins so that the size of the plotting region is always square, even if you resize the graphics window.

### Converting units

For many applications, you need to be able to translate user coordinates to pixels or inches. There are some cryptic shortcuts, but the simplest way is to get the range in user coordinates and measure the proportion of the graphics device devoted to the plotting region.

```
user.range <- par("usr")[c(2,4)] -  
           par("usr")[c(1,3)]
```

```
region.pct <- par("plt")[c(2,4)] -  
              par("plt")[c(1,3)]
```

```
region.px <-  
           dev.size(units="px") * region.pct
```

```
px.per.xy <- region.px / user.range
```

To convert a horizontal or distance from the x-coordinate value to pixels, multiply by **px.per.xy[1]**. To convert a vertical distance, multiply by **region.px.per.xy[2]**. To convert a diagonal distance, you need to invoke Pythagoras.

```
a.px <- x.dist*px.per.xy[1]  
b.px <- y.dist*px.per.xy[2]  
c.px <- sqrt(a.px^2+b.px^2)
```

To rotate a string to match the slope of a line segment, you need to convert the distances to pixels, calculate the arctangent, and convert from radians to degrees.

```
segments(x0, y0, x1, y1)  
delta.x <- (x1 - x0) * px.per.xy[1]  
delta.y <- (y1 - y0) * px.per.xy[2]  
angle.radians <- atan2(delta.y, delta.x)  
angle.degrees <- angle.radians * 180 / pi  
text(x1, y1, "TEXT", srt=angle.degrees)
```

## Panels

`par("fig")` (width, height) in pct  
`par("fin")` (width, height) in inches

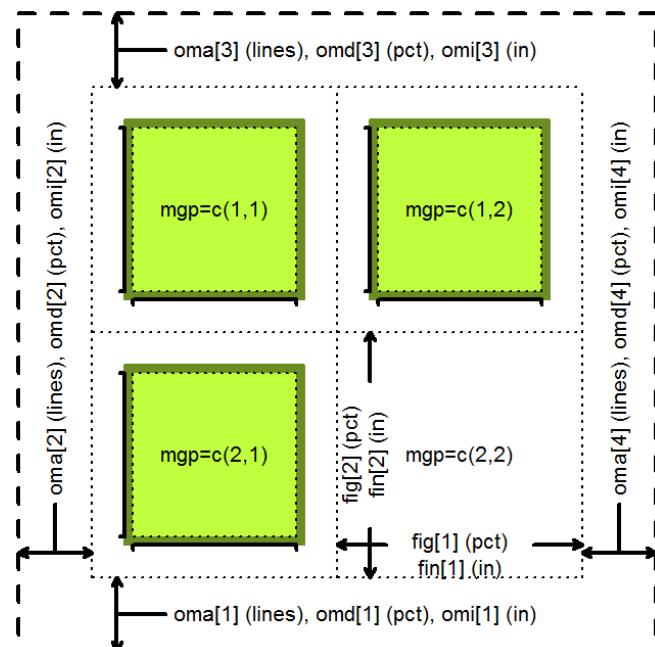
If you display multiple plots within a single graphics window (e.g., with the `mfrw` or `mfcol` arguments of `par` or with the `layout` function), then the `fig` and `fin` arguments will tell you the size of the current subplot window in percent or inches, respectively.

`par("oma")` (bottom, left, top, right) in lines  
`par("omd")` (bottom, left, top, right) in pct  
`par("omi")` (bottom, left, top, right) in inches

Each subplot will have margins specified by `mai` or `mar`, but no outer margin around the entire set of plots, unless you specify them using `oma`, `omd`, or `omi`. You can place text in the outer margins using the `mtext` function with the argument `outer=TRUE`.

`par("mpg")` (r, c) or (r, c, maxr, maxc)

The `mpg` argument of `par` will allow you to jump to a subplot in a particular row and column. If you query with `par("mpg")`, you will get the current row and column followed by the maximum row and column.



## Character and string sizes

### `strheight()`

The `strheight` functions will tell you the height of a specified string in inches (`units="inches"`), x-y user coordinates (`units="user"`) or as a percentage of the graphics device (`units="figure"`).

For a single line of text, `strheight` will give you the height of the letter "M". If you have a string with one or more linebreaks ("\\n"), the `strheight` function will measure the height of the letter "M" plus the height of one or more additional lines. The height of a line is dependent on the line spacing, set by the `lheight` argument of `par`. The default line height (`lheight=1`), corresponding to single spaced lines, produces a line height roughly 1.5 times the height of "M".

### `strwidth()`

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

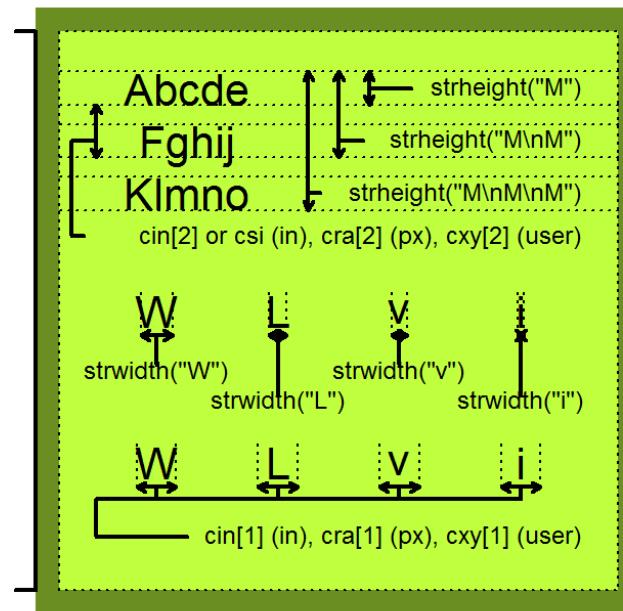
`par("cin")` (r.o.) (width, height) in inches  
`par("csi")` (r.o.) height in inches  
`par("cra")` (r.o.) (width, height) in pixels  
`par("cxy")` (r.o.) (width, height) in xy coordinates

The single value returned by the `csi` argument of `par` gives you the height of a line of text in inches. The second of the two values returned by `cin`, `cra`, and `cxy` gives you the height of a line, in inches, pixels, or xy (user) coordinates.

The first of the two values returned by the `cin`, `cra`, and `cxy` arguments to `par` gives you the approximate width of a single character, in inches, pixels, or xy (user) coordinates. The width, very slightly smaller than the actual width of the letter "W", is a rough estimate at best and ignores the variable width of individual letters.

These values are useful, however, in providing fast ratios of the relative sizes of the differing units of measure

`px.per.in <- par("cra") / par("cin")`  
`px.per.xy <- par("cra") / par("cxy")`  
`xy.per.in <- par("cxy") / par("cin")`



## If your fonts are too big or too small

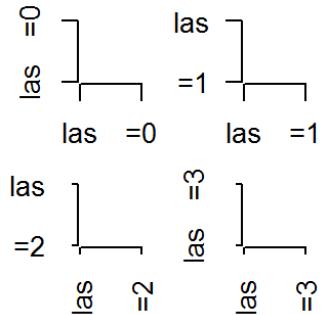
Fixing this takes a bit of trial and error.

- Specify a larger/smaller value for the `pointsize` argument when you open your graphics device.
- Try opening your graphics device with different values for `height` and `width`. Fonts that look too big might be better proportioned in a larger graphics window.
- Use the `cex` argument to increase or decrease the relative size of your fonts.

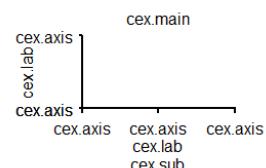
## If your axes don't fit

There are several possible solutions.

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



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



# Basic Regular Expressions in R

## Cheat Sheet

Character Classes	
<code>[:digit:]</code> or <code>\d</code>	Digits; [0-9]
<code>\D</code>	Non-digits; [^0-9]
<code>[:lower:]</code>	Lower-case letters; [a-z]
<code>[:upper:]</code>	Upper-case letters; [A-Z]
<code>[:alpha:]</code>	Alphabetic characters; [A-z]
<code>[:alnum:]</code>	Alphanumeric characters [A-z0-9]
<code>\w</code>	Word characters; [A-z0-9_]
<code>\W</code>	Non-word characters
<code>[:xdigit:]</code> or <code>\x</code>	Hexadec. digits; [0-9A-Fa-f]
<code>[:blank:]</code>	Space and tab
<code>[:space:]</code> or <code>\s</code>	Space, tab, vertical tab, newline, form feed, carriage return
<code>\S</code>	Not space; [^[:space:]]
<code>[:punct:]</code>	Punctuation characters; !#\$%&'()*, -./;=>?[@]^_`{ ~
<code>[:graph:]</code>	Graphical char.; [:alnum:][:punct:]\s
<code>[:print:]</code>	Printable characters; [:alnum:][:punct:]\s
<code>[:cntrl:]</code> or <code>\c</code>	Control characters; \n, \r etc.

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

Lookaheads and Conditionals*	
<code>(?=)</code>	Lookahead (requires PERL = TRUE), e.g. (?=yx): position followed by 'yx'
<code>(?!)</code>	Negative lookahead (PERL = TRUE); position NOT followed by pattern
<code>(?&lt;=)</code>	Lookbehind (PERL = TRUE), e.g. (?<=yx): position following 'xy'
<code>(?&lt;=)</code>	Negative lookbehind (PERL = TRUE); position NOT following pattern
<code>?(if)then</code>	If-then-condition (PERL = TRUE); use lookaheads, optional char. etc in if-clause
<code>?(if)then else</code>	If-then-else-condition (PERL = TRUE)

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

The diagram illustrates four functions for pattern matching:

- Detect pattern:** Shows a magnifying glass icon over a string represented by a sequence of colored squares (blue, red, blue, red, blue).
- Locate pattern:** Shows two arrows pointing to specific positions in a string represented by a sequence of colored squares (blue, red, blue, red, blue).
- Extract pattern:** Shows a red box highlighting a portion of a string represented by a sequence of colored squares (red, red, blue, blue).
- Replace pattern:** Shows a sequence of colored squares where some are blue and others are yellow, representing a string where some parts have been replaced.

**Functions for Pattern Matching**

**Detect pattern**

**Locate pattern**

**Extract pattern**

**Replace pattern**

```
> string <- c("Hipopotamus", "Rhymoceros", "time for bottomless lyrics")
> pattern <- "t.m"
```

**Detect Patterns**

`grep(pattern, string)`  
[1] 1 3

`grep(pattern, string, value = TRUE)`  
[1] "Hipopotamus"  
[2] "time for bottomless lyrics"

`grepl(pattern, string)`  
[1] TRUE FALSE TRUE

`stringr::str_detect(string, pattern)`  
[1] TRUE FALSE TRUE

**Locate Patterns**

`regexpr(pattern, string)`  
find starting position and length of first match

`gregexpr(pattern, string)`  
find starting position and length of all matches

`stringr::str_locate(string, pattern)`  
find starting and end position of first match

`stringr::str_locate_all(string, pattern)`  
find starting and end position of all matches

**Replace Patterns**

`sub(pattern, replacement, string)`  
replace first match

`gsub(pattern, replacement, string)`  
replace all matches

`stringr::str_replace(string, pattern, replacement)`  
replace first match

`stringr::str_replace_all(string, pattern, replacement)`  
replace all matches

**Character Classes and Groups**

- `.` Any character except \n
- `|` Or, e.g. (a|b)
- `[...]` List permitted characters, e.g. [ abc ]
- `[a-z]` Specify character ranges
- `[^...]` List excluded characters
- `(...)` Grouping, enables back referencing using \\N where N is an integer

**Anchors**

- `^` Start of the string
- `$` End of the string
- `\b` Empty string at either edge of a word
- `\B` NOT the edge of a word
- `\B` Beginning of a word
- `\E` End of a word

**Quantifiers**

- `*` Matches at least 0 times
- `+` Matches at least 1 time
- `?` Matches at most 1 time; optional string
- `{n}` Matches exactly n times
- `{n,}` Matches at least n times
- `{,n}` Matches at most n times
- `{n,m}` Matches between n and m times

**General Modes**

By default R uses *POSIX extended regular expressions*. You can switch to *PCRE regular expressions* using PERL = TRUE for base or by wrapping patterns with perl() for stringr.

All functions can be used with literal searches using fixed = TRUE for base or by wrapping patterns with fixed() for stringr.

All base functions can be made case insensitive by specifying ignore.cases = TRUE.

**Escaping Characters**

Metacharacters ( . \* + etc.) can be used as literal characters by escaping them. Characters can be escaped using \\ or by enclosing them in \\Q...\\E.

**Case Conversions**

Regular expressions can be made case insensitive using (?i). In backreferences, the strings can be converted to lower or upper case using \\L or \\U (e.g. \\L\\1). This requires PERL = TRUE.

**Greedy Matching**

By default the asterisk \* is greedy, i.e. it always matches the longest possible string. It can be used in lazy mode by adding ?, i.e. ?\*.

Greedy mode can be turned off using (?U). This switches the syntax, so that (?U)a\* is lazy and (?U)a\*? is greedy.

**Note**

Regular expressions can conveniently be created using `rex::rex()`.

# R For Data Science Cheat Sheet

## xts

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



### xts

**eXtensible Time Series (xts)** is a powerful package that provides an extensible time series class, enabling uniform handling of many R time series classes by extending zoo.

Load the package as follows:

```
> library(xts)
```

### xts Objects

xts objects have three main components:

- **coredata**: always a matrix for xts objects, while it could also be a vector for zoo objects
- **index**: vector of any Date, POSIXct, chron, yearmon, yearqtr, or DateTime classes
- **xtsAttributes**: arbitrary attributes

### Creating xts Objects

```
> xts1 <- xts(x=1:10, order.by=Sys.Date()-1:10)
> data <- rnorm(5)
> dates <- seq(as.Date("2017-05-01"), length=5, by="days")
> xts2 <- xts(x=data, order.by=dates)
> xts3 <- xts(x=rnorm(10),
+               order.by=as.POSIXct(Sys.Date()+1:10),
+               born=as.POSIXct("1899-05-08"))
> xts4 <- xts(x=1:10, order.by=Sys.Date()+1:10)
```

### Convert To And From xts

```
> data(AirPassengers)
> xts5 <- as.xts(AirPassengers)
```

### Import From Files

```
> dat <- read.csv(tmp_file)
> xts(dat, order.by=as.Date(rownames(dat), "%m/%d/%Y"))
> dat_zoo <- read.zoo(tmp_file,
+                      index.column=0,
+                      sep=",",
+                      format="%m/%d/%Y")
> dat_zoo <- read.zoo(tmp, sep=",", FUN=as.yearmon)
> dat_xts <- as.xts(dat_zoo)
```

### Inspect Your Data

```
> core_data <- coredata(xts2)
> index(xts1)
```

Extract core data of objects  
Extract index of objects

### Class Attributes

```
> indexClass(xts2)
> indexClass(convertIndex(xts, 'POSIXct'))
> indexTZ(xts5)
> indexFormat(xts5) <- "%Y-%m-%d"
```

Get index class  
Replacing index class  
Get index class  
Change format of time display

### Time Zones

```
> tzzone(xts1) <- "Asia/Hong_Kong"
> tzzone(xts1)
```

Change the time zone  
Extract the current time zone

## Export xts Objects

```
> data_xts <- as.xts(matrix)
> tmp <- tempfile()
> write.zoo(data_xts, sep=",", file=tmp)
```

## Replace & Update

```
> xts2[dates] <- 0
> xts5["1961"] <- NA
> xts2["2016-05-02"] <- NA
```

Replace values in xts2 on dates with 0  
Replace dates from 1961 with NA  
Replace the value at 1 specific index with NA

## Applying Functions

```
> ep1 <- endpoints(xts4, on="weeks", k=2)
[1] 0 5 10
> ep2 <- endpoints(xts5, on="years")
[1] 0 12 24 36 48 60 72 84 96 108 120 132 144
> period.apply(xts5, INDEX=ep2, FUN=mean)
> xts5_yearly <- split(xts5, f="years")
> lapply(xts5_yearly, FUN=mean)
> do.call(rbind,
+          lapply(split(xts5, "years"),
+                function(w) last(w, n="1 month")))
> do.call(rbind,
+          lapply(split(xts5, "years"),
+                cumsum))
> rollapply(xts5, 3, sd)
```

Take index values by time  
Calculate the yearly mean  
Split xts5 by year  
Create a list of yearly means  
Find the last observation in each year in xts5  
Calculate cumulative annual passengers  
Apply sd to rolling margins of xts5

## Selecting, Subsetting & Indexing

### Select

```
> mar55 <- xts5["1955-03"]
```

Get value for March 1955

### Subset

```
> xts5_1954 <- xts5["1954"]
> xts5_janmarch <- xts5["1954/1954-03"]
> xts5_janmarch <- xts5["/1954-03"]
> xts4[ep1]
```

Get all data from 1954  
Extract data from Jan to March '54  
Get all data until March '54  
Subset xts4 using ep2

### first() and last()

```
> first(xts4, '1 week')
> first(last(xts4, '1 week'), '3 days')
```

Extract first 1 week  
Get first 3 days of the last week of data

### Indexing

```
> xts2[index(xts3)]
> days <- c("2017-05-03", "2017-05-23")
> xts3[days]
> xts2[as.POSIXct(days, tz="UTC")]
> index <- which(.indexwday(xts1)==0).indexwday(xts1)==6)
> xts1[index]
```

Extract rows with the index of xts3  
Extract rows using the vector days  
Extract rows using days as POSIXct  
Index of weekend days  
Extract weekend days of xts1

## Missing Values

```
> na.omit(xts5)
> xts_last <- na.locf(xts2)
> xts_last <- na.locf(xts2,
+                      fromLast=TRUE)
> na.approx(xts2)
```

Omit NA values in xts5  
Fill missing values in xts2 using last observation  
Fill missing values in xts2 using next observation  
Interpolate NAs using linear approximation

## Arithmetic Operations

### coredata() or as.numeric()

```
> xts3 + as.numeric(xts2)
> xts3 * as.numeric(xts4)
> coredata(xts4) - xts3
> coredata(xts4) / xts3
```

Addition  
Multiplication  
Subtraction  
Division

### Shifting Index Values

```
> xts5 - lag(xts5)
> diff(xts5, lag=12, differences=1)
```

Period-over-period differences  
Lagged differences

### Reindexing

```
> xts1 + merge(xts2, index(xts1), fill=0)
[1] 0.231538
[2] 0.829257
[3] 4.000000
[4] 3.000000
[5] 2.000000
[6] 1.000000
> xts1 - merge(xts2, index(xts1), fill=na.locf)
[1] 0.231538
[2] 0.829257
[3] 4.829257
[4] 3.829257
[5] 2.829257
[6] 1.829257
```

Addition

Subtraction

## Merging

```
> merge(xts2, xts1, join='inner')
[1] 0.231538
[2] 0.829257
[3] 4.000000
[4] 3.000000
[5] 2.000000
[6] 1.000000
> merge(xts2, xts1, join='left', fill=0)
[1] 0.231538
[2] 0.829257
[3] 4.829257
[4] 3.829257
[5] 2.829257
[6] 1.829257
```

Inner join of xts2 and xts1

Left join of xts2 and xts1, fill empty spots with 0

Combine xts1 and xts4 by rows

## Other Useful Functions

```
> .index(xts4)
> .indexwday(xts3)
> .indexhour(xts3)
> start(xts3)
> end(xts4)
> str(xts3)
> time(xts1)
> head(xts2)
> tail(xts2)
```

Extract raw numeric index of xts1  
Value of week(day), starting on Sunday, in index of xts3  
Value of hour in index of xts3  
Extract first observation of xts3  
Extract last observation of xts4  
Display structure of xts3  
Extract raw numeric index of xts1  
First part of xts2  
Last part of xts2



# The eurostat package

## R tools to access open data from Eurostat database

### Search and download

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

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

### Find the table code

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

```
library("eurostat")
query <- search_eurostat("road", type = "table")
query[1:3,1:2]
##          title      code
## 1 Goods transport by road ttr00005
## 2 People killed in road accidents tsdtr420
## 3 Enterprises with broadband access tin00090
```

### Download the table

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

```
dat <- get_eurostat(id="tsdtr420", time_format="num")
head(dat)
##   unit sex geo time values
## 1 NR T AT 1999 1079
## 2 NR T BE 1999 1397
## 3 NR T CZ 1999 1455
## 4 NR T DK 1999 514
## 5 NR T EL 1999 2116
## 6 NR T ES 1999 5738
```

### Add labels

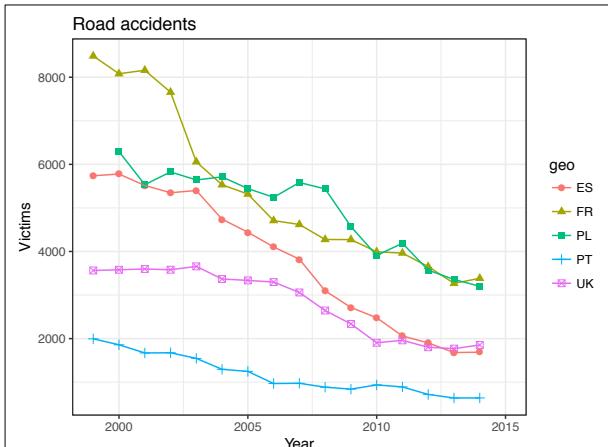
The `label_eurostat(x, lang = "en", ...)` gets definitions for Eurostat codes and replace them with labels in given language ("en", "fr" or "de").

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

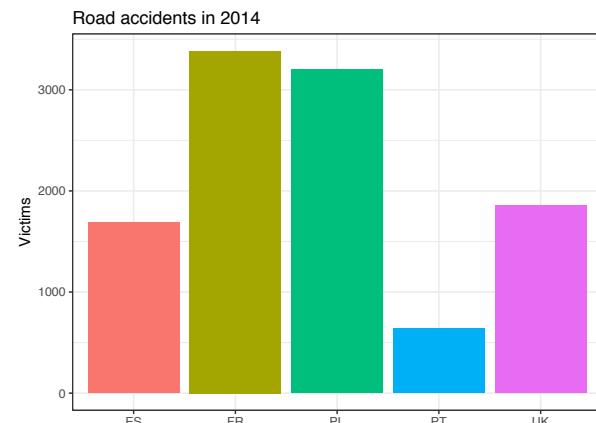
### eurostat and plots

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

```
t1 <- get_eurostat("tsdtr420", filters =
  list(geo = c("UK", "FR", "PL", "ES", "PT")))
library("ggplot2")
ggplot(t1, aes(x = time, y = values, color = geo,
  group = geo, shape = geo)) +
  geom_point(size = 2) +
  geom_line() + theme_bw() +
  labs(title = "Road accidents", x = "Year", y = "Victims")
```



```
library("dplyr")
t2 <- t1 %>% filter(time == "2014-01-01")
ggplot(t2, aes(geo, values, fill = geo)) +
  geom_bar(stat = "identity") + theme_bw() +
  theme(legend.position = "none") +
  labs(title = "Road accidents in 2014", x = "", y = "Victims")
```



### eurostat and maps

#### Fetch and process data

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

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

```
mapdata <- merge_eurostat_geodata(fertility,
  resolution = "20")
```

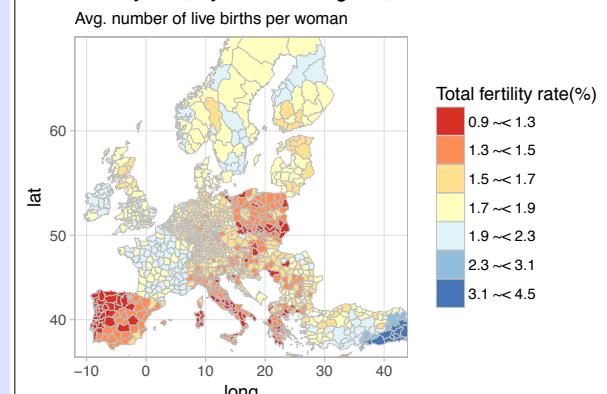
```
head(select(mapdata, geo, values, cat, long, lat, order, id))
##   geo values cat long lat order id
## 1 AT124 1.3 ~< 1.5 15.54245 48.90770 214 10
## 2 AT124 1.3 ~< 1.5 15.75363 48.85218 215 10
## 3 AT124 1.3 ~< 1.5 15.88763 48.78511 216 10
## 4 AT124 1.3 ~< 1.5 15.81535 48.69270 217 10
## 5 AT124 1.3 ~< 1.5 15.94094 48.67173 218 10
## 6 AT124 1.3 ~< 1.5 15.90833 48.59815 219 10
```

#### Draw a cartogram

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

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

#### Fertility rate, by NUTS-3 regions, 2014

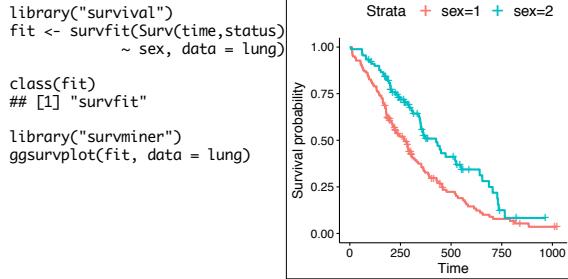


# Creating Survival Plots

## Informative and Elegant with survminer

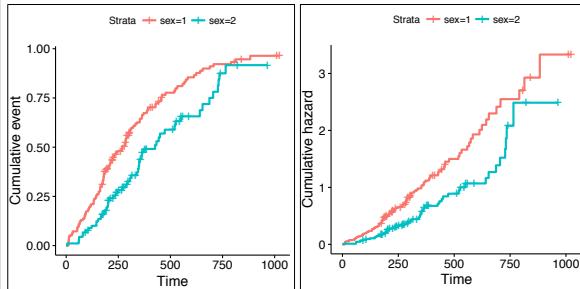
### Survival Curves

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



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

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



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

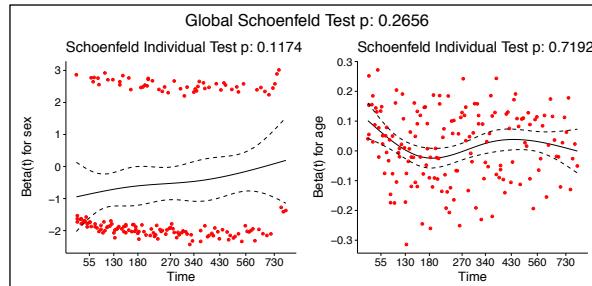
```
ggsurvplot(fit, data = lung,
conf.int = TRUE,
pval = TRUE,
fun = "pct",
risk.table = TRUE,
size = 1,
linetype = "strata",
palette = c("#E78000",
 "#2E9FDF"),
legend = "bottom",
legend.title = "Sex",
legend.labs = c("Male",
 "Female"))
```

Time	Male	Female
0	138	90
250	62	53
500	20	21
750	7	3
1000	2	0

### Diagnostics of Cox Model

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

```
library("survival")
fit <- coxph(Surv(time, status) ~ sex + age, data = lung)
ftest <- cox.zph(fit)
ftest
##          rho chisq      p
## sex     0.1236 2.452 0.117
## age    -0.0275 0.129 0.719
## GLOBAL   NA 2.651 0.266
library("survminer")
gcoxzph(ftest)
```



The function `ggcoxdiagnostics()` plots different types of residuals as a function of time, linear predictor or observation id. The type of residual is selected with `type` argument. Possible values are "martingale", "deviance", "score", "schoenfeld", "dfbeta", "dfbetas", and "scaledsch".

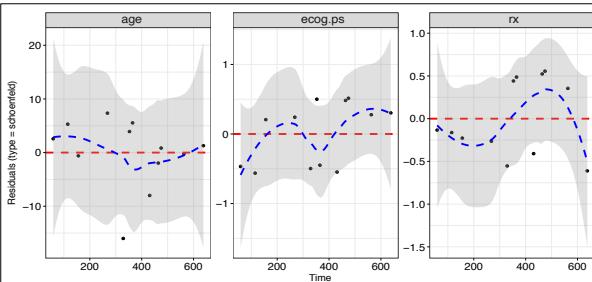
The `ox.scale` argument defines what shall be plotted on the OX axis. Possible values are "linear.predictions", "observation.id", "time".

Logical arguments `hline` and `sline` may be used to add horizontal line or smooth line to the plot.

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

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

```
ggcoxdiagnostics(fit,
type = "schoenfeld",
ox.scale = "time")
```

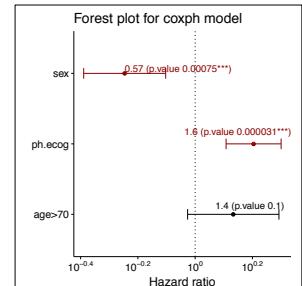


### Summary of Cox Model

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

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

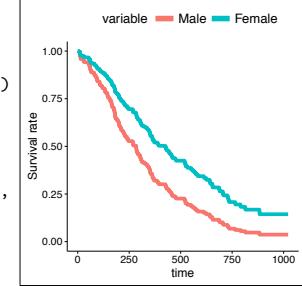
## Call:
## coxph(formula = Surv(time, status) ~ sex+ph.ecog+age, data=lung)
##
##            coef exp(coef) se(coef)      z      p
## sex      -0.567   0.567  0.168 -3.37 0.00075
## ph.ecog   0.470   1.600  0.113  4.16 3.1e-05
## age>70   0.307   1.359  0.187  1.64 0.10175
##
## Likelihood ratio test=31.6 on
## n= 227, number of events= 164
ggforest(fit)
```



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

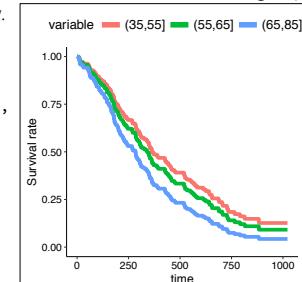
Note that these curves differ from Kaplan Meier estimates since they present expected survival based on given Cox model.

```
library("survival")
library("survminer")
lung$sex <- ifelse(lung$sex == 1,
"Male", "Female")
fit <- coxph(Surv(time, status) ~ sex + ph.ecog + age,
data = lung)
ggcoxadjustedcurves(fit, data=lung,
variable=lung$sex)
```



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

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



# Leaflet Cheat Sheet



for 

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

## Quick Start

### Installation

Use `install.packages("leaflet")` to install the package or directly from Github [devtools::install\\_github\("rstudio/leaflet"\)](#).

### First Map

```
m <- leaflet() %>%
  addTiles() %>%
  addMarkers(lng = 174.768, lat = -36.852, popup = "The birthplace of R")
# add a single point layer
```



## Map Widget

### Initialization

<code>m &lt;- leaflet(options = leafletOptions(...))</code>	
<code>center</code>	Initial geographic center of the map
<code>zoom</code>	Initial map zoom level
<code>minZoom</code>	Minimum zoom level of the map
<code>maxZoom</code>	Maximum zoom level of the map

### Map Methods

```
m %>% setView(lng, lat, zoom, options = list())
Set the view of the map (center and zoom level)
m %>% fitBounds(lng1, lat1, lng2, lat2)
Fit the view into the rectangle [lng1, lat1] - [lng2, lat2]
m %>% clearBounds()
Clear the bound, automatically determine from the map elements
```

### Data Object

Both `leaflet()` and the `map` layers have an optional data parameter that is designed to receive spatial data with the following formats:

#### Base R

The arguments of all layers take normal R objects:

```
df <- data.frame(lat = ..., lng = ...)
leaflet(df) %>% addTiles() %>% addCircles()
library(sp) Useful functions:
SpatialPoints, SpatialLines, SpatialPolygons, ...
library(maps) Build a map of states with colors:
mapStates <- map("state", fill = TRUE, plot = FALSE)
leaflet(mapStates) %>% addTiles() %>%
  addPolygons(fillColor = topo.colors(10, alpha =
    NULL), stroke = FALSE)
```

#### sp package

`library(sp)` Useful functions:

```
SpatialPoints, SpatialLines, SpatialPolygons, ...
library(maps) Build a map of states with colors:
mapStates <- map("state", fill = TRUE, plot = FALSE)
leaflet(mapStates) %>% addTiles() %>%
  addPolygons(fillColor = topo.colors(10, alpha =
    NULL), stroke = FALSE)
```

#### maps package

## Markers

Use markers to call out points, express locations with latitude/longitude coordinates, appear as icons or as circles.

Data come from vectors or assigned data frame, or `sp` package objects.

### Icon Markers

Regular Icons: default and simple

`addMarkers(lng, lat, popup, label)` [add basic icon markers](#)

`makeIcon(Icons)` (`iconUrl`, `iconWidth`, `iconHeight`, `iconAnchorX`, `iconAnchorY`, `shadowUrl`, `shadowWidth`, `shadowHeight`, ...) [customize marker icons](#)  
`iconList()` [create a list of icons](#)

Awesome Icons: customizable with colors and icons

`addAwesomeMarkers, makeAwesomeIcon, awesomeIcons, awesomeIconList`

Marker Clusters: option of `addMarkers()`

`clusterOptions = markerClusterOptions()`

`freezeAtZoom` Freeze the cluster at assigned zoom level

### Circle Markers

`addCircleMarkers(color, radius, stroke, opacity, ...)`

Customize their color, radius, stroke, opacity

## Popups and Labels

`addPopups(lng, lat, ...content..., options)` [Add standalone popups](#)

`options = popupOptions(closeButton=FALSE)`

`addMarkers(..., popup, ...)` [Show popups with markers or shapes](#)

`addMarkers(..., label, labelOptions...)` [Show labels with markers or shapes](#)

`labelOptions = labelOptions(noHide, textOnly, textSize, direction, style)`

`addLabelOnlyMarkers()` [Add labels without markers](#)

## Lines and Shapes

### Polygons and Polylines

`addPolygons(color, weight=1, smoothFactor=0.5, opacity=1.0, fillOpacity=0.5, fillColor=~colorQuantile("YIOrRd", ALAND)(ALAND), highlightOptions, ...)`

`highlightOptions(color, weight=2, bringToFront=TRUE)` [highlight shapes](#)

Use `rmapshaper::ms_simplify` to simplify complex shapes

`Circles` `addCircles(lng, lat, weight=1, radius, ...)`

`Rectangles` `addRectangles(lng1, lat1, lng2, lat2, fillColor="transparent", ...)`

## Basemaps

`addTiles()`



`providers$Stamen.Toner, CartoDB.Positron, Esri.NatGeoWorldMap`



Default Tiles

Use `addTiles()` to add a custom map tile URL template, use `addWMSTiles()` to add WMS (Web Map Service) tiles

## GeoJSON and TopoJSON

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

\* To read into `sp` objects with the `geojsonio` or `rgdal` package:  
`geojsonio::geojson_read(..., what="sp") rgdal::readOGR(..., "OGRGeoJSON")`

\* Or to use the `addGeoJSON()` and `addTopoJSON()` functions:  
`addTopoJSON/addGeoJSON(... weight, color, fill, opacity, fillOpacity, ...)`

Styles can also be tuned separately with a `style: [...]` object.

Other packages including `RJSONIO` and `jsonlite` can help fast parse or generate the data needed.

## Shiny Integration

To integrate a Leaflet map into an app:

\* In the UI, call `leafletOutput("name")`

\* On the server side, assign a `renderLeaflet(...)` call to the output

\* Inside the `renderLeaflet` expression, return a Leaflet map object

### Modification

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

Other useful functions to edit your map:

`fitBounds(o, 0, 11, 11)` [similar to setView](#)

[fit the view to within these bounds](#)

`addCircles(1:10, 1:10, layerId = LETTERS[1:10])`

[create circles with layerIds of "A", "B", "C"...](#)

`removeShape(c("B", "F"))` [remove some of the circles](#)

`clearShapes()` [clear all circles \(and other shapes\)](#)

### Inputs/Events

#### Object Events

Object event names generally use this pattern:

`input$MAPID_OBJECTCATEGORY_EVENTNAME`

Trigger an event changes the value of the Shiny input at this variable.

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

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

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

\* `lat` The latitude of the object, if available; otherwise, the mouse cursor

\* `lng` The longitude of the object, if available; otherwise, the mouse cursor

\* `id` The `layerId`, if any

GeoJSON events also include additional properties:

\* `featureId` The feature ID, if any

\* `properties` The feature properties

#### Map Events

`input$MAPID_click` [when the map background or basemap is clicked](#)  
value -- a list with `lat` and `lng`

`input$MAPID_bounds` [provide the lat{lng bounds of the visible map area](#)  
value -- a list with `north`, `east`, `south` and `west`

`input$MAPID_zoom` [an integer indicates the zoom level](#)