

# R Programming Cheat Sheet

## JUST THE BASICS

CREATED BY: ARIANNE COLTON AND SEAN CHEN

### GENERAL

- R version 3.0 and greater adds support for 64 bit integers
- R is case sensitive
- R index starts from 1

### HELP

`help(functionName)` OR `?functionName`

Help Home Page	<code>help.start()</code>
Special Character Help	<code>help('[')</code>
Search Help	<code>help.search(..)</code> or <code>??..</code>
Search Function - with Partial Name	<code>apropos('mea')</code>
See Example(s)	<code>example(topic)</code>

### OBJECTS in current environment

Display Object Name	<code>objects()</code> OR <code>ls()</code>
Remove Object	<code>rm(object1, object2,..)</code>

#### Notes:

- .name starting with a period are accessible but invisible, so they will not be found by 'ls'
- To guarantee memory removal, use 'gc', releasing unused memory to the OS. R performs automatic 'gc' periodically

### SYMBOL NAME ENVIRONMENT

- If multiple packages use the same function name the function that the package loaded the last will get called.
- To avoid this precede the function with the name of the package. e.g. `packageName::functionName(..)`

### LIBRARY

Only trust reliable R packages i.e., 'ggplot2' for plotting, 'sp' for dealing spatial data, 'reshape2', 'survival', etc.

Load Package	<code>library(packageName)</code> OR <code>require(packageName)</code>
Unload Package	<code>detach(packageName)</code>

**Note:** `require()` returns the status(True/False)

### MANIPULATING STRINGS

Putting Together Strings	<code>paste('string1', 'string2', sep = '/')</code> # separator ('sep') is a space by default <code>paste(c('1', '2'), collapse = '/')</code> # returns '1/2'
Split String	<code>stringr::str_split(string = v1, pattern = '-')</code> # returns a list
Get Substring	<code>stringr::str_sub(string = v1, start = 1, end = 3)</code>
Match String	<code>isJohnFound &lt;- stringr::str_detect(string = df1\$col1, pattern = ignore.case('John'))</code> # returns True/False if John was found <code>df1[isJohnFound, c('col1', ...)]</code>

### DATA TYPES

**Check data type:** `class(variable)`

#### FOUR BASIC DATA TYPES

- Numeric** - includes float/double, int, etc.

`is.numeric(variable)`

- Character(string)**

`nchar(variable)` # length of a character or numeric

- Date/POSIXct**

- Date:** stores just a date. In numeric form, number of days since 1/1/1970 (see below).

`date1 <- as.Date('2012-06-28'), as.numeric(date1)`

- POSIXct:** stores a date and time. In numeric form, number of seconds since 1/1/1970.

`date2 <- as.POSIXct('2012-06-28 18:00')`

**Note:** Use 'lubridate' and 'chron' packages to work with Dates

- Logical**

- (TRUE = 1, FALSE = 0)
- Use `==/!=` to test equality and inequality

`as.numeric(TRUE) => 1`

## DATA STRUCTURES

### VECTOR

- Group of elements of the SAME type
- R is a vectorized language, operations are applied to each element of the vector automatically
- R has no concept of column vectors or row vectors
- Special vectors: letters and LETTERS, that contain lower-case and upper-case letters

Create Vector	<code>v1 &lt;- c(1, 2, 3)</code>
Get Length	<code>length(v1)</code>
Check if All or Any is True	<code>all(v1); any(v1)</code>
Integer Indexing	<code>v1[1:3]; v1[c(1,6)]</code>
Boolean Indexing	<code>v1[is.na(v1)] &lt;- 0</code>
Naming	<code>c(first = 'a', ..) or names(v1) &lt;- c('first', ..)</code>

### FACTOR

- `as.factor(v1)` gets you the levels which is the number of unique values
- Factors can reduce the size of a variable because they only store unique values, but could be buggy if not used properly

### LIST

Store any number of items of ANY type

Create List	<code>list1 &lt;- list(first = 'a', ...)</code>
Create Empty List	<code>vector(mode = 'list', length = 3)</code>
Get Element	<code>list1[[1]]</code> or <code>list1[['first']]</code>
Append Using Numeric Index	<code>list1[[6]] &lt;- 2</code>
Append Using Name	<code>list1[['newElement']] &lt;- 2</code>

**Note:** repeatedly appending to list, vector, data.frame etc. is expensive, it is best to create a list of a certain size, then fill it.

### DATA.FRAME

- Each column is a variable, each row is an observation
- Internally, each column is a vector
- idata.frame is a data structure that creates a reference to a data.frame, therefore, no copying is performed

Create Data Frame	<code>df1 &lt;- data.frame(col1 = v1, col2 = v2, v3)</code>
Dimension	<code>nrow(df1); ncol(df1); dim(df1)</code>
Get/Set Column Names	<code>names(df1)</code> <code>names(df1) &lt;- c(...)</code>
Get/Set Row Names	<code>rownames(df1)</code> <code>rownames(df1) &lt;- c(...)</code>
Preview	<code>head(df1, n = 10); tail(...)</code>
Get Data Type	<code>class(df1)</code> # is data.frame
Index by Column(s)	<code>df1['col1']</code> or <code>df1[1]</code> ; <code>df1[c('col1', 'col3')]</code> or <code>df1[c(1, 3)]</code>
Index by Rows and Columns	<code>df1[c(1, 3), 2:3]</code> # returns data from row 1 & 3, columns 2 to 3

† Index method: `df1$col1` OR `df1[, 'col1']` OR `df1[, 1]` returns as a vector. To return single column

data.frame while using single-square brackets, use 'drop': `df1[, 'col1', drop = FALSE]`

### DATA.TABLE

#### What is a data.table

- Extends and enhances the functionality of data.frames

#### Differences: data.table vs. data.frame

- By default data.frame turns character data into factors, while data.table does not
- When you print data.frame data, all data prints to the console, with a data.table, it intelligently prints the first and last five rows
- Key Difference:** Data.tables are fast because they have an index like a database.  
i.e., this search, `dt1$col1 > number`, does a sequential scan (vector scan). After you create a key for this, it will be much faster via binary search.

Create data.table from data.frame	<code>data.table(df1)</code>
Index by Column(s)*	<code>dt1[, 'col1', with = FALSE]</code> or <code>dt1[, list(col1)]</code>
Show info for each data.table in memory (i.e. size, ...)	<code>tables()</code>
Show Keys in data.table	<code>key(dt1)</code>
Create index for col1 and reorder data according to col1	<code>setkey(dt1, col1)</code>
Use Key to Select Data	<code>dt1[c('col1Value1', 'col1Value2'), ]</code>
Multiple Key Select	<code>dt1[J('1', c('2', '3')), ]</code>
Aggregation**	<code>dt1[, list(col1 = mean(col1), by = col2)]</code> <code>dt1[, list(col1 = mean(col1), col2Sum = sum(col2)), by = list(col3, col4)]</code>

\* Accessing columns must be done via list of actual names, not as characters. If column names are characters, then "with" argument should be set to FALSE.

\*\* Aggregate and d'ply functions will work, but built-in aggregation functionality of data table is faster

### MATRIX

- Similar to data.frame except every element must be the SAME type, most commonly all numerics
- Functions that work with data.frame should work with matrix as well

Create Matrix	<code>matrix1 &lt;- matrix(1:10, nrow = 5), # fills rows 1 to 5, column 1 with 1:5, and column 2 with 6:10</code>
Matrix Multiplication	<code>matrix1 %*% t(matrix2)</code> # where t() is transpose

### ARRAY

- Multidimensional vector of the SAME type
- `array1 <- array(1:12, dim = c(2, 3, 2))`
- Using arrays is not recommended
- Matrices are restricted to two dimensions while array can have any dimension

# DATA MUNGING

## APPLY (apply, tapply, lapply, mapply)

- Apply - most restrictive. Must be used on a matrix, all elements must be the same type
- If used on some other object, such as a data.frame, it will be converted to a matrix first

```
apply(matrix1, 1 - rows or 2 - columns,
function to apply)
# if rows, then pass each row as input to the function
```

- By default, computation on NA (missing data) always returns NA, so if a matrix contains NAs, you can ignore them (use `na.rm = TRUE` in the `apply(...)` which doesn't pass NAs to your function)

## lapply

Applies a function to each element of a list and returns the results as a list

## sapply

Same as lapply except return the results as a vector

**Note:** lapply & sapply can both take a vector as input, a vector is technically a form of list

## AGGREGATE (SQL GROUPBY)

- `aggregate(formulas, data, function)`
- Formulas:  $y \sim x$ ,  $y$  represents a variable that we want to make a calculation on,  $x$  represents one or more variables we want to group the calculation by
- Can only use one function in `aggregate()`. To apply more than one function, use the `plyr()` package

In the example below diamonds is a data.frame; price, cut, color etc. are columns of diamonds.

```
aggregate(price ~ cut, diamonds, mean)
# get the average price of different cuts for the diamonds
aggregate(price ~ cut + color, diamonds,
mean) # group by cut and color
aggregate(cbind(price, carat) ~ cut,
diamonds, mean) # get the average price and average
carat of different cuts
```

## PLYR ('split-apply-combine')

- `ddply()`, `lply()`, `ldply()`, etc. (1st letter = the type of input, 2nd = the type of output)
- `plyr` can be slow, most of the functionality in `plyr` can be accomplished using base function or other packages, but `plyr` is easier to use

## ddply

Takes a data.frame, splits it according to some variable(s), performs a desired action on it and returns a data.frame

## lply

- Can use this instead of lapply
- For sapply, can use `lapply` ('a' is array/vector/matrix), however, `lply` result does not include the names.

## DPLYR (for data.frame ONLY)

- Basic functions: `filter()`, `slice()`, `arrange()`, `select()`, `rename()`, `distinct()`, `mutate()`, `summarise()`,

`group_by()`, `sample_n()`

## Chain functions

```
df1 %>% group_by(year, month) %>%
select(col1, col2) %>% summarise(collmean
= mean(col1))
```

- Much faster than `plyr`, with four types of easy-to-use joins (inner, left, semi, anti)
- Abstracts the way data is stored so you can work with data frames, data tables, and remote databases with the same set of functions

## HELPER FUNCTIONS

`each()` - supply multiple functions to a function like `aggregate`

```
aggregate(price ~ cut, diamonds, each(mean,
median))
```

# DATA

## LOAD DATA FROM CSV

### Read csv

```
read.table(file = url or filepath, header =
TRUE, sep = ',')
```

- "stringAsFactors" argument defaults to TRUE, set it to FALSE to prevent converting columns to factors. This saves computation time and maintains character data
- Other useful arguments are "quote" and "colClasses", specifying the character used for enclosing cells and the data type for each column.
- If cell separator has been used inside a cell, then use `read.csv2()` or `read.delim2()` instead of `read.table()`

## DATABASE

Connect to Database	<code>db1 &lt;- RODB::odbcConnect('conStr')</code>
Query Database	<code>df1 &lt;- RODB::sqlQuery(db1, 'SELECT ..', stringAsFactors = FALSE)</code>
Close Connection	<code>RODB::odbcClose(db1)</code>

- Only one connection may be open at a time. The connection automatically closes if R closes or another connection is opened.
- If table name has space, use `[ ]` to surround the table name in the SQL string.
- `which()` in R is similar to 'where' in SQL

## INCLUDED DATA

R and some packages come with data included.

List Available Datasets	<code>data()</code>
List Available Datasets in a Specific Package	<code>data(package = 'ggplot2')</code>

## MISSING DATA (NA and NULL)

NULL is not missing, it's nothingness. NULL is atomical and cannot exist within a vector. If used inside a vector, it simply disappears.

Check Missing Data	<code>is.na()</code>
Avoid Using	<code>is.null()</code>

# FUNCTIONS AND CONTROLS

Create Function	<code>say_hello &lt;- function(first, last = 'hola') { }</code>
Call Function	<code>say_hello(first = 'hello')</code>

- R automatically returns the value of the last line of code in a function. This is bad practice. Use `return()` explicitly instead.
- `do.call()` - specify the name of a function either as string (i.e. 'mean') or as object (i.e. mean) and provide arguments as a list.

```
do.call(mean, args = list(first = '1st'))
```

## IF /ELSE /ELSE IF /SWITCH

	if { } else	ifelse
Works with Vectorized Argument	No	Yes
Most Efficient for Non-Vectorized Argument	Yes	No
Works with NA *	No	Yes
Use <code>&amp;&amp;</code> , <code>  </code> **†	Yes	No
Use <code>&amp;</code> , <code> </code> ***†	No	Yes

\* NA == 1 result is NA, thus if won't work, it'll be an error. For `ifelse`, NA will return instead

\*\* `&&`, `||` is best used in `if`, since it only compares the first element of vector from each side

\*\*\* `&`, `|` is necessary for `ifelse`, as it compares every element of vector from each side

† `&&`, `||` are similar to `if` in that they don't work with vectors, where `ifelse`, `&`, `|` work with vectors

- Similar to C++/Java, for `&`, `|`, both sides of operator are always checked. For `&&`, `||`, if left side fails, no need to check the right side.
- `} else, else` must be on the same line as }

# GRAPHICS

## DEFAULT BASIC GRAPHIC

```
hist(df1$col1, main = 'title', xlab = 'x
axis label')
plot(col2 ~ col1, data = df1),
aka y ~ x or plot(x, y)
```

## LATTICE AND GGLOT2 (more popular)

- Initialize the object and add layers (points, lines, histograms) using `+`, map variable in the data to an axis or aesthetic using 'aes'

```
ggplot(data = df1) + geom_histogram(aes(x
= col1))
```

- Normalized histogram (pdf, not relative frequency histogram)

```
ggplot(data = df1) + geom_density(aes(x =
col1), fill = 'grey50')
```

# DATA RESHAPING

## REARRANGE

Melt Data - from column to row	<code>reshape2.melt(df1, id.vars = c('col1', 'col2'), variable.name = 'newCol1', value.name = 'newCol2')</code>
Cast Data - from row to column	<code>reshape2.dcast(df1, col1 + col2 ~ newCol1, value.var = 'newCol2')</code>

If `df1` has 3 more columns, `col3` to `col5`, 'melting' creates a new df that has 3 rows for each combination of `col1` and `col2`, with the values coming from the respective `col3` to `col5`.

## COMBINE (multiple sets into one)

### 1. cbind - bind by columns

data.frame from two vectors	<code>cbind(v1, v2)</code>
data.frame combining df1 and df2 columns	<code>cbind(df1, df2)</code>

### 2. rbind - similar to cbind but for rows, you can assign new column names to vectors in cbind

```
cbind(col1 = v1, ...)
```

### 3. Joins - (merge, join, data.table) using common keys

#### 3.1 Merge

- `by.x` and `by.y` specify the key columns use in the `join()` operation

- Merge can be much slower than the alternatives

```
merge(x = df1, y = df2, by.x = c('col1',
'col3'), by.y = c('col3', 'col6'))
```

#### 3.2 Join

- Join in `plyr()` package works similar to `merge` but much faster, drawback is key columns in each table must have the same name

- `join()` has an argument for specifying left, right, inner joins

```
join(x = df1, y = df2, by = c('col1',
'col3'))
```

#### 3.3 data.table

```
dt1 <- data.table(df1, key = c('1',
'2')), dt2 <- ... ‡
```

- Left Join

```
dt1[dt2]
```

‡ Data table join requires specifying the keys for the data tables

Created by Arianne Colton and Sean Chen  
[data.scientist.info@gmail.com](mailto:data.scientist.info@gmail.com)

Based on content from  
 'R for Everyone' by Jared Lander

Updated: December 2, 2015

# R Programming Cheat Sheet

ADVANCED

CREATED BY: ARIANNE COLTON AND SEAN CHEN

## ENVIRONMENTS

### ENVIRONMENT BASICS

#### What is an Environment?

Data structure (that powers lexical scoping) is made up of two components, the frame, which contains the name-object bindings (and behaves much like a named list), and the parent environment.

#### Named List

- You can think of an environment as a 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.

#### Parent Environment

- Every environment has a parent, another environment. Only one environment doesn't have a parent: the empty environment.
- The parent is used to implement lexical scoping: if a name is not found in an environment, then R will look in its parent (and so on).

Environments can also be useful data structures in their own right because they have reference semantics.

### FOUR SPECIAL ENVIRONMENTS

- Global environment**, access with `globalenv()`, is the interactive workspace. This is the environment in which you normally work.

The parent of the global environment is the last package that you attached with `library()` or `require()`.

- Base environment**, access with `baseenv()`, is the environment of the base package. Its parent is the empty environment.

- Empty environment**, access with `emptyenv()`, is the ultimate ancestor of all environments, and the only environment without a parent. Empty environments contain nothing.

- Current environment**, access with `environment()`

### SEARCH PATH

#### What is the Search Path?

An R internal mechanism to look up objects, specifically, functions.

- Access with `search()`, which lists all parents of the global environment. (See Figure 1)
- It contains one environment for each attached package.
- Objects in the search path environments can be found from the top-level interactive workspace.

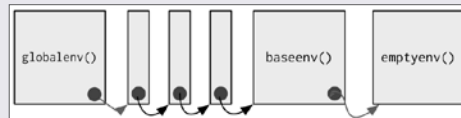


Figure 1. The Search Path

- If you look for a name in a search, it will always start from global environment first, then inside the latest attached package.

If there are functions with the same name in two different packages, the latest package will get called.

- Each time you load a new package with `library()` or `require()` it is inserted between the global environment and the package that was previously at the top of the search path.

```
search() :  
'_GlobalEnv' ... 'Autoloads' 'package:base'  
library(reshape2); search()  
'_GlobalEnv' 'package:reshape2' ...  
'Autoloads' 'package:base'
```

**Note:** There is a special environment called Autoloads which is used to save memory by only loading package objects (like big datasets) when needed.

## ENVIRONMENTS

Access any environment on the search list	<code>as.environment('package:base')</code>
Find the environment where a name is defined	<code>pryr::where('func1')</code>

### FUNCTION ENVIRONMENTS

There are 4 environments for functions.

#### 1. Enclosing environment (used for lexical scoping)

- When a function is created, it gains a reference to the environment where it was made. This is the enclosing environment.
- The enclosing environment belongs to the function, and never changes, even if the function is moved to a different environment.
- Every function has one and only one enclosing environment. For the three other types of environments, there may be 0, 1, or many environments associated with each function.
- You can determine the enclosing environment of a function by calling i.e. `environment(func1)`

#### 2. Binding environment

- The binding environments of a function are all the environments which have a binding to it.
- The enclosing environment determines how the function finds values; the binding environments determine how we find the function.

Example for enclosing and binding environment

```
y <- 1  
e <- new.env()  
e$g <- function(x) x + y  
# function g enclosing environment is the global  
# environment, and the binding environment is 'e'.
```

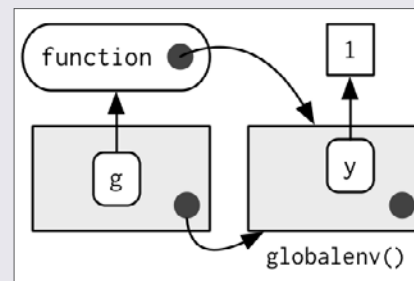


Figure 2. Function Environment

**Note:** Every R package has two environments associated with it (package and namespace). Every exported function is bound into the package environment, but enclosed by the namespace environment.

#### 3. Execution environment

- Each time a function is called, a new environment is created to host execution. The parent of the execution environment is the enclosing environment of the function.
- Once the function has completed, this environment is thrown away.

**Note:** Each execution environment has two parents: a calling environment and an enclosing environment.

- R's regular scoping rules only use the enclosing parent; `parent.frame()` allows you to access the calling parent.

#### 4. Calling environment

- This is the environment where the function was called.
- Looking up variables in the calling environment rather than in the enclosing environment is called dynamic scoping.
- Dynamic scoping is primarily useful for developing functions that aid interactive data analysis.

## BINDING NAMES TO VALUES

### Assignment

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

### Name rules

- A complete list of reserved words can be found in `?Reserved`.

### Regular assignment arrow, <-

- The regular assignment arrow always creates a variable in the current environment.

### Deep assignment arrow, <<-

- The deep assignment arrow modifies an existing variable found by walking up the parent environments. If `<<-` doesn't find an existing variable, it will create one in the global environment. This is usually undesirable, because global variables introduce non-obvious dependencies between functions.

## ENVIRONMENT CREATION

- To create an environment manually, use `new.env()`. You can list the bindings in the environment's frame with `ls()` and see its parent with `parent.env()`.
- When creating your own environment, note that you should set its parent environment to be the empty environment. This ensures you don't accidentally inherit objects from somewhere else.

# FUNCTIONS

## FUNCTION BASICS

The most important thing to understand about R is that **functions are objects** in their own right.

All R functions have three parts:

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

- When you print (`func1`) a function in R, it shows you these three important components. If the environment isn't displayed, it means that the function was created in the global environment.
- Like all objects in R, functions can also possess any number of additional attributes().

Every operation is a function call

- Everything that exists is an object
- Everything that happens in R is a function call, even if it doesn't look like it. (i.e. +, for, if, [, \$, { ...})

**Note:** the backtick (`), lets you refer to functions or variables that have otherwise reserved or illegal names: e.g. `x + y` is the same as ``+`(x, y)`

## LEXICAL SCOPING

What is Lexical Scoping?

Looks up value of a symbol. (See "Enclosing Environment" in the "Environment" section.)

- `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 "+"
```

\* This doesn't work because R relies on lexical scoping to find everything, even the + operator. It's never possible to make a function completely self-contained because you must always rely on functions defined in base R or other packages.

## FUNCTION ARGUMENTS

When calling a function you can specify arguments by position, by complete name, or by partial name. Arguments are matched first by exact name (perfect matching), then by prefix matching, and finally by position.

- Function arguments are passed by **reference** and **copied on modify**.
- You can determine if an argument was supplied or not with the `missing()` function.

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

- By default, R function arguments are lazy -- they're only evaluated if they're actually used

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

However, since x is not used. `stop("This is an error!")` never get evaluated.

- Default arguments are evaluated inside the function. This means that if the expression depends on the current environment the results will differ depending on whether you use the default value or explicitly provide one:

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

<code>f()</code>	<code>'a' 'x'</code>	<code>ls()</code> evaluated inside f
<code>f(ls())</code>		<code>ls()</code> evaluated in global environment

## RETURN VALUES

- The last expression evaluated in a function becomes the return value, the result of invoking the function.
- Only use explicit `return()` for when you are returning early, such as for an error.
- Functions can return only a single object. But this is not a limitation because you can return a list containing any number of objects.
- Functions can return invisible values, which are not printed out by default when you call the function.

```
f1 <- function() 1
f2 <- function() invisible(1)
```

- The most common function that returns invisibly is `<-`

## PRIMITIVE FUNCTIONS

- There is one exception to the rule that functions have three components.
- Primitive functions, like `sum()`, call C code directly with `.Primitive()` and contain no R code.
- Therefore their `formals()`, `body()`, and `environment()` are all NULL:

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

- Primitive functions are only found in the base package, and since they operate at a low level, they can be more efficient.

## INFIX FUNCTIONS

- Most functions in R are 'prefix' operators: the name of the function comes before the arguments.

- You can also create infix functions where the function name comes in between its arguments, like + or -. All user-created infix functions must start and end with `%`.

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

- 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

- Act like they modify their arguments in place, and have the special name `xxx<-`
- They typically have two arguments (x and value), although they can have more, and they must return the modified object.

```
`second<-` <- function(x, value) {
  x[2] <- value
  x
}
x <- 1:10
second(x) <- 5L
```

- I say they "act" like they modify their arguments in place, because they actually create a modified copy.
- We can see that by using `pryr::address()` to find the memory address of the underlying object.

# DATA STRUCTURES

	Homogeneous	Heterogeneous
<b>1d</b>	Atomic vector	List
<b>2d</b>	Matrix	Data frame
<b>nd</b>	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()`

# DATA STRUCTURES

- 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

\* Internally, data.frame is a list of equal-length vectors.

## 1D (VECTORS: ATOMIC VECTOR AND LIST)

- Use `is.atomic()` || `is.list()` to test if an object is actually a vector, not `is.vector()`.

Type	typeof()	what it is
Length	<code>length()</code>	how many elements
Attributes	<code>attributes()</code>	additional arbitrary metadata

## FACTORS

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

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

- While factors look (and often behave) like character vectors, they are actually integers. Be careful when treating them like strings.
- Factors are useful when you know the possible values a variable may take, even if you don't see all values in a given dataset.
- Most data loading functions in R automatically convert character vectors to factors, use the argument `stringsAsFactors = FALSE` to suppress this behavior.

## ATTRIBUTES

- All objects can have arbitrary additional attributes.
- Attributes can be accessed individually with `attr()` or all at once (as a list) with `attributes()`.

```
attr(v1, 'attr1') <- 'my vector'
```

- By default, most attributes are lost when modifying a vector. The only attributes not lost are the three most important:

Names	a character vector giving each element a name	<code>names(x)</code>
Dimensions	used to turn vectors into matrices and arrays	<code>dim(x)</code>
Class	used to implement the S3 object system	<code>class(x)</code>



# SUBSETTING (OPERATORS: [, [[, \$)

## SIMPLIFYING VS. PRESERVING SUBSETTING

- **Simplifying** subsetting returns the **simplest** possible data structure that can represent the output.
- **Preserving** subsetting keeps the structure of the output the **same** as the input.

	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, ]</code> or <code>x[, 1]</code>	<code>x[1, , drop = F]</code> or <code>x[, 1, drop = F]</code>
Data frame	<code>x[, 1]</code> or <code>x[[1]]</code>	<code>x[, 1, drop = F]</code> or <code>x[1]</code>

- When you use `drop = FALSE`, it's preserving.
- Omitting `drop = FALSE` when subsetting matrices and data frames is one of the most common sources of programming errors.
- `[[` is similar to `[`, except it can only return a single value and it allows you to pull pieces out of a list.

\* Simplifying behavior varies slightly between different data types:

- **Atomic Vector:** `x[[1]]` is the same as `x[1]`.
- **List:** `[ ]` always returns a list, to get the contents use `[[ ]]`.
- **Factor:** drops any unused levels but it remains a factor class.
- **Matrix or array:** if any of the dimensions has length 1, drops that dimension.
- **Data frame** is similar, if output is a single column, it returns a vector instead of a data frame.

## DATA.FRAMES SUBSETTING

- Data frames possess 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.

List Subsetting	<code>df1[c('col1', 'col2')]</code>
Matrix Subsetting	<code>df1[, c('col1', 'col2')]</code>

The subsetting results are the **same** in this example.

- **Single column subsetting:** matrix subsetting simplifies by default, list subsetting does not.

<code>str(df1[, 'col1']) -&gt; int [1:3]</code> # the result is a vector
<code>str(df1['col1']) -&gt; 'data.frame'</code> # the result remains a data frame of 1 column

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

## OUT OF BOUNDS

- `[` and `[[` differ slightly in their behavior when the index is out of bounds (OOB).
- For example, when you try to extract the fifth element of a length four vector, aka OOB `x[5] -> NA`, or subset a vector with NA or NULL: `x[NULL] -> x[0]`

Operator	Index	Atomic	List
<code>[</code>	OOB	NA	list(NULL)
<code>[</code>	NA_real_	NA	list(NULL)
<code>[</code>	NULL	<code>x[0]</code>	list(NULL)
<code>[[</code>	OOB	Error	Error
<code>[[</code>	NA_real_	Error	NULL
<code>[[</code>	NULL	Error	Error

- If the input vector is named, then the names of OOB, missing, or NULL components will be "`<NA>`".

## \$ SUBSETTING OPERATOR

- `$` is a useful shorthand for `[[` combined with character subsetting:

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

- One common mistake with `$` is to try and use 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]]
```

- There's one important difference between `$` and `[[`, `$` does partial matching, `[[` does not:

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

## SUBSETTING WITH ASSIGNMENT

- All subsetting operators can be combined with assignment to modify selected values of the input vector.
- Subsetting with nothing can be useful in conjunction with assignment because it will preserve the original object class and structure.

```
df1[] <- lapply(df1, as.integer)
# df1 will remain as a data frame
```

```
df1 <- lapply(df1, as.integer)
# df1 will become a list
```

## EXAMPLES

### 1. Lookup tables (character subsetting)

Character matching provides a powerful way to make lookup tables.

```
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), ]
```

- If you have multiple columns to match on, you'll need to first collapse them to a single column (with `interaction()`, `paste()`, or `plyr::id()`).
- You can also use `merge()` or `plyr::join()`, which do the same thing for you.

### 3. Expanding aggregated counts (integer subsetting)

- Sometimes you get a data frame where identical rows have been collapsed into one and a count column has been added.

- `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 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 the columns you want	<code>df1[c('col1', 'col2')]</code>

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

- Logical subsetting is probably the most commonly used technique for extracting rows out of a data frame.

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

- Remember to use the vector boolean operators `&` and `|`, not the short-circuiting scalar operators `&&` and `||` which are more useful inside if statements.
- `subset()` is a specialised shorthand function for subsetting data frames, and saves some typing because you don't need to repeat the name of the data frame.

```
subset(df1, col1 == 5 & col2 == 4)
```

## BOOLEAN ALGEBRA VS. SETS (LOGICAL & INTEGER SUBSETTING)

- It's useful to be aware of the natural equivalence between set operations (integer subsetting) and boolean algebra (logical subsetting).
- Using set operations 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()` allows you to convert a boolean representation to an integer representation. There's no reverse operation in base R.

```
which(c(T, F, T, F)) -> 1 3
# returns the index of the true*
```

\* The integer representation length is always `<=` boolean representation length.

- When first learning subsetting, a common mistake is to use `x[which(y)]` instead of `x[y]`.
- Here the `which()` achieves nothing, it switches from logical to integer subsetting but the result will be exactly the same.
- Also beware that `x[-which(y)]` is not equivalent to `x[!y]`. If `y` is all FALSE, `which(y)` will be `integer(0)` and `-integer(0)` is still `integer(0)`, so you'll get no values, instead of all values.
- In general, avoid switching from logical to integer subsetting unless you want, for example, the first or last TRUE value.

# DEBUGGING, CONDITION HANDLING, & DEFENSIVE PROGRAMMING

## DEBUGGING

Use `traceback()` and `browser()`, and interactive tools in RStudio:

- RStudio's error inspector or `traceback()` which list the sequence of calls that lead to the error.
- RStudio's breakpoints or `browser()` which open an interactive debug session at an arbitrary location in the code.
- RStudio's "Rerun with Debug" tool or `options(error = browser)*` which open an interactive debug session where the error occurred.

\* There are two other useful functions that you can use with the error option:

- Recover** is a step up from `browser`, as it allows you to enter the environment of any of the calls in the call stack.

This is useful because often the root cause of the error is a number of calls back.

- dump.frames** is an equivalent to recover for non-interactive code. It creates a `last.dump.rda` file in the current working directory.

Then, in a later interactive R session, you load that file, and use `debugger()` to enter an interactive debugger with the same interface as `recover()`. This allows interactive debugging of batch code.

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)

- Communicating potential problems** to the user is the job of **conditions**: errors, warnings, and messages:

- Fatal errors are raised by `stop()` and force all execution to terminate. Errors are used when there is no way for a function to continue.
- Warnings are generated by `warning()` and are used to display potential problems, such as when some elements of a vectorised input are invalid.
- Messages are generated by `message()` and are used to give informative output in a way

that can easily be suppressed by the user using `?suppressMessages()`.

### 2. Handling conditions programmatically:

- `try()` gives you the ability to continue execution even when an error occurs.
- `tryCatch()` 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.

- You can also capture the output of the `try()` and `tryCatch()` functions.

If successful, it will be the last result evaluated in the block, just like a function.

If unsuccessful it will be an invisible object of class "try-error".

### 3. Custom signal classes:

- One of the challenges of error handling in R is that most functions just call `stop()` with a string.
- Since conditions are S3 classes, the solution is to define your own classes if you want to distinguish different types of error.
- Each condition signalling function, `stop()`, `warning()`, and `message()`, can be given either a list of strings, or a custom S3 condition object.

## DEFENSIVE PROGRAMMING

The basic principle of defensive programming is to **"fail fast"**, to raise an error as soon as something goes wrong.

In R, this takes three particular forms:

- Checking that inputs are correct using `stopifnot()`, the `'assertthat'` package, or simple if statements and `stop()`
- Avoiding non-standard evaluation like `subset()`, `transform()`, and `with()`.  
These functions save time when used interactively, but because they make assumptions to reduce typing, when they fail, they often fail with uninformative error messages.
- Avoiding functions that can return different types of output. The two biggest offenders are `[]` and `apply()`.

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

# OBJECT ORIENTED (OO) FIELD GUIDE

## OBJECT ORIENTED SYSTEMS

R has three object oriented systems (plus the base types)

- S3** is a very casual system. It has no formal definition of classes. S3 implements a style of OO programming called generic-function OO.

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

Example:	<code>drawRect(canvas, 'blue')</code>
Language:	R

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

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

- S4** works similarly to S3, but is more formal. There are two major differences to S3.

- S4 has formal class definitions, which describe the representation and inheritance for each class, and has special helper functions for defining generics and methods.
- S4 also has multiple dispatch, which means that generic functions can pick methods based on the class of any number of arguments, not just one.

- Reference classes**, called RC for short, are quite different from S3 and S4.

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

## C STRUCTURE

- Underlying every R object is a C structure (or struct) that describes how that object is stored in memory.**

- The struct includes the contents of the object, the information needed for memory management and a **type**.

`typeof()` # determines an object's base type

- The **"Data structures"** section explains the most common base types (atomic vectors and lists), but base types also encompass functions, environments, and other more exotic objects like names, calls, and promises.

- To see if an object is a pure base type, (i.e., it doesn't also have S3, S4, or RC behavior), check that `is.object(x)` returns FALSE.

## S3

- S3 is R's first and simplest OO system. It is the only OO system used in the base and stats package.

- In S3, methods belong to functions, called generic functions, or generics for short. S3 methods do not belong to objects or classes.

- Given a class, the job of an S3 generic is to call the right S3 method. You can recognise S3 methods by their names, which look like `generic.class()`.

For example, the `Date` method for the `mean()` generic is called `mean.Date()`

This is the reason that most modern style guides discourage the use of `.` in function names, it makes them look like S3 methods.

- See all methods that belong to a generic :

```
methods('mean')
#> mean.Date
#> mean.default
#> mean.difftime
```

- List all generics that have a method for a given class :

```
methods(class = 'Date')
```

- S3 objects are usually built on top of lists, or atomic vectors with attributes. Factor and data frame are S3 class.

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

Created by Arianne Colton and Sean Chen  
[data.scientist.info@gmail.com](mailto:data.scientist.info@gmail.com)

Based on content from  
"Advanced R" by Hadley Wickham  
Updated: January 15, 2016