

Statistics 360: Advanced R for Data Science

Lecture 1

Brad McNeney

Course Objectives

R objects

Course Objectives

Course objectives

- ▶ Work through most, but not all of the book Advanced R by Hadley Wickham: <https://adv-r.hadley.nz/index.html>
- ▶ Topics:
 - ▶ R objects: names and values
 - ▶ Basic data structures and programming.
 - ▶ vectors, subsetting, control flow, functions, environments
 - ▶ No tidyverse this time
 - ▶ Object-oriented programming in R
 - ▶ Code performance: debugging, profiling, memory, calling Python, or C++ from R
 - ▶ Parallelizing R code (if time permits)

Getting started with R, RStudio and git

- ▶ Follow the “getting started” instructions on the class canvas page to get set up with R, RStudio and git.
 - ▶ R and RStudio will be familiar, but you may not have used git before, so leave some time for that.
- ▶ Please try to get R and RStudio installed and create an RStudio project linked to the class GitHub repository (or a forked copy) as soon as possible.
- ▶ Those still having trouble after the weekend should ask our TA, Pulindu, for help during the first lab sessions in week 2.
 - ▶ **Note: No lab first week.**

Reading

- ▶ Welcome, Preface and Chapters 1,2 of the text.

R objects

R objects

- ▶ In R, data structures and functions are all referred to as “objects”.
- ▶ Objects are created with the assignment operator `<-`; e.g.,
`x <- c(1,2,3)`.
 - ▶ The objects a user creates from the R console are contained in the user’s workspace, called the global environment.
 - ▶ Use `ls()` to see a list of all objects in the workspace.
 - ▶ Use `rm(x)` to remove object `x` from the workspace.

Digging deeper

- ▶ The above understanding is an over-simplification that is usually OK, but will sometimes lead to misunderstandings about memory usage and when R makes copies of objects
- ▶ Object copying is a **major** source of computational overhead in R, so it pays to understand what will trigger it.
- ▶ Reference: text, chapter 2

Binding names and to objects

- ▶ The R code `x <- c(1,2,3)` does two things: (i) creates an object in computer memory that contains the values 1, 2, 3 and (ii) “binds” that object to the “name” x.

```
# install.packages("lobstr")  
library(lobstr)  
x <- c(1,2,3)  
ls()
```

```
## [1] "x"
```

```
obj_addr(x) # changes every time this code chunk is run
```

```
## [1] "0x13584f328"
```

Binding multiple names to the same object

- ▶ The following binds the name `y` to the same object that `x` is bound to.

```
y <- x  
obj_addr(y)
```

```
## [1] "0x13584f328"
```

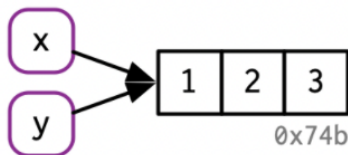


Figure 1: Bind two names to the same object

Aside: Syntactic vs non-syntactic

- ▶ Valid, or “syntactic” names in R can consist of letters, digits, . and _ but should start with a letter.
- ▶ Names that start with . are hidden from directory listing with `ls()`.
- ▶ Names that start with _ or a digit are non-syntactic and will cause an error.
- ▶ If you need to create or access a non-syntactic name, use backward single-quotes (“backticks”).

```
x <- 1  
.x <- 1  
`_x` <- 1  
ls()
```

```
## [1] "_x" "x"  "y"
```

Modifying causes copying

- ▶ Modifying a variable causes a copy to be made, with the modified variable name bound to the copy.

```
c(obj_addr(x),obj_addr(y))
```

```
## [1] "0x1249c35f0" "0x13584f328"
```

```
x[[2]] <- 10 # Note: x[2] <- 10 has the same effect  
c(obj_addr(x),obj_addr(y))
```

```
## [1] "0x1242e4608" "0x13584f328"
```

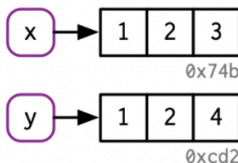


Figure 2: Bindings after copying

Tracing copying

- ▶ The `tracemem()` function marks an object so that a message is printed whenever a copy is made.

```
x <- c(1,2,3)
tracemem(x)
```

```
## [1] "<0x134aa01f8>"
```

```
x[[2]] <- 10
```

```
## tracemem[0x134aa01f8 -> 0x134aaa8e8]: eval eval withVisible w
```

```
untracemem(x)  # remove the trace
x[[1]] <- 10
```

More on tracemem()

- ▶ As the output of `tracemem()` suggests, the trace is on the object, not the name:

```
x <- c(1,2,3)
tracemem(x)
```

```
## [1] "<0x134beea48>"
```

```
y <- x
c(obj_addr(x),obj_addr(y))
```

```
## [1] "0x134beea48" "0x134beea48"
```

```
y[[2]] <- 10
```

```
## tracemem[0x134beea48 -> 0x1258b3318]: eval eval withVisible w
c(obj_addr(x),obj_addr(y))
```

```
## [1] "0x134beea48" "0x1258b3318"
```

Function calls

- R has a reputation for passing copies to functions, but in fact the copy-on-modify applies to functions too:

```
f <- function(arg) { return(arg) }  
x <- c(1,2,3)  
y <- f(x) # no copy made, so x and y bound to same obj  
c(obj_addr(x),obj_addr(y))
```

```
## [1] "0x134a9e998" "0x134a9e998"
```

```
f <- function(arg) { arg <- 2*arg; return(arg) }  
y <- f(x) # copy made  
c(obj_addr(x),obj_addr(y))
```

```
## [1] "0x134a9e998" "0x134d8c6b8"
```


Lists

- List elements point to objects too:

```
l1 <- list(1, 2, 3)
c(obj_addr(l1),obj_addr(l1[[1]]),obj_addr(l1[[2]]),obj_addr(l1[[3]]))
```

```
## [1] "0x125379998" "0x1249dd5b8" "0x1249dd580" "0x1249dd548"
```

```
# Note: ref(l1) will print a nicely formatted version of the above,  
# but doesn't work with my slides
```

Copy-on-modify in lists

```
l1 <- l2 <- list(1,2,3)
```

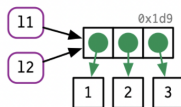


Figure 3: Bindings before

```
l2[[3]] <- 4
```

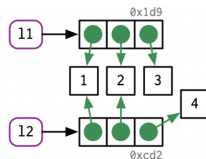


Figure 4: Bindings after modify

Copies of lists are said to be “shallow”

- ▶ As shown above, we copy the list itself and any list **elements** that are modified. This is called a “shallow” copy.
- ▶ By contrast, a “deep” copy would be a copy of all elements.

```
l1 <- list(1,2,3)
c(obj_addr(l1),obj_addr(l1[[1]]),obj_addr(l1[[2]]),obj_addr(l1[[3]]))
```

```
## [1] "0x1259d22b8" "0x1257579b8" "0x125757980" "0x125757948"
```

```
l1[[3]] <- 4
c(obj_addr(l1),obj_addr(l1[[1]]),obj_addr(l1[[2]]),obj_addr(l1[[3]]))
```

```
## [1] "0x134aab068" "0x1257579b8" "0x125757980" "0x1257577c0"
```

Data frames are lists with columns as list items

```
dd <- data.frame(x=1:3,y=4:6)
c(obj_addr(dd[[1]]),obj_addr(dd[[2]]))
```

```
## [1] "0x124842578" "0x124842498"
```

```
dd[,2] <- 7:9
c(obj_addr(dd[[1]]),obj_addr(dd[[2]])) # only changes second element
```

```
## [1] "0x124842578" "0x12540ac88"
```

```
dd[1,] <- c(11,22)
c(obj_addr(dd[[1]]),obj_addr(dd[[2]])) # changes to both elements
```

```
## [1] "0x13485a138" "0x13485a0e8"
```

```
dd[1,1] <- 111
c(obj_addr(dd[[1]]),obj_addr(dd[[2]])) # only changes first element
```

```
## [1] "0x1252142a8" "0x13485a0e8"
```

Beware of data frame overhead

- ▶ Data frames are convenient, but the convenience comes at a cost.

- ▶ For example, coercion to/from lists

```
dd <- data.frame(x=rnorm(100)) # try yourself with rnorm(1e7)  
tracemem(dd); tracemem(dd[[1]])
```

```
## [1] "<0x124918278>"
```

```
## [1] "<0x133fbf690>"
```

```
dmed <- lapply(dd,median) # makes a list copy of dd
```

```
## tracemem[0x124918278 -> 0x12499a2e8]: as.list.data.frame as.list lap
```

```
## tracemem[0x133fbf690 -> 0x133fbe910]: sort.int sort.default sort mea
```

```
dd[[1]] <- dd[[1]] - dmed[[1]] #
```

```
## tracemem[0x124918278 -> 0x1249bb548]: eval eval withVisible withCall
```

```
## tracemem[0x1249bb548 -> 0x1249be0a8]: [[<-.data.frame [[<- eval eval
```

- Fewer copies if we do the same with a list.

```
ll <- list(x=rnorm(100))  
tracemem(ll); tracemem(ll[[1]])
```

```
## [1] "<0x1249ff628>"
```

```
## [1] "<0x133f92af0>"
```

```
lmed <- lapply(ll,median) # no need for a list copy
```

```
## tracemem[0x133f92af0 -> 0x133f92e40]: sort.int sort.default sort.mean
```

```
ll[[1]] <- ll[[1]] - lmed[[1]]
```

```
## tracemem[0x1249ff628 -> 0x1251289b0]: eval eval.withVisible withCall
```

Modify-in-place

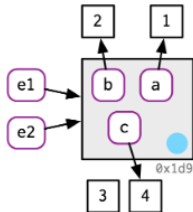
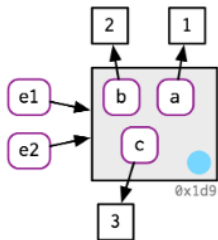
- ▶ The text says there are two exceptions to the copy-on-modify:
 1. Modify an element of an object with one binding, or
 2. Modify an environment. but in my experiments, only the second applies.

```
v <- c(1,2,3) # creates object (1,2,3) and binds v to it
tracemem(v)
```

```
## [1] "<0x1249b50f8>"
```

```
v[[3]] <- 4 # for me, this triggers a copy
```

```
## tracemem[0x1249b50f8 -> 0x12590e408]: eval eval withVisible withCall
```



```
e1 <- rlang::env(a = 1, b = 2, c = 3)
e2 <- e1
# note: can't use tracemem() on an environment
e1$c <- 4
e2$c
```

```
## [1] 4
```


Object size

- Use `lobstr::obj_size()` to find the size of objects.

```
obj_size(dd)
```

```
## 1,528 B
```

```
obj_size(ll)
```

```
## 1,128 B
```

```
obj_size(e1)
```

```
## 840 B
```

```
obj_size(e2)
```

```
## 840 B
```