

Managing and Manipulating Data Using R

Lecture 10, Looping

Ozan Jaquette

1. Introduction
2. Accessing elements of vectors and lists
3. Loop basics
4. Different ways to loop over a vector
5. Loop tips
6. Modifying vs. creating new object
7. More Practice

Introduction

Libraries

```
library(tidyverse)
#> -- Attaching packages -----
#> v ggplot2 3.0.0      v purrr  0.2.5
#> v tibble  1.4.2      v dplyr  0.7.6
#> v tidyr   0.8.1      v stringr 1.3.1
#> v readr   1.1.1      v forcats 0.3.0
#> -- Conflicts -----
#> x dplyr::filter() masks stats::filter()
#> x dplyr::lag()    masks stats::lag()
```

Accessing elements of vectors and lists

Types of vectors

Recall that there are two broad types of vectors, **atomic vectors** and **lists**

1. **Atomic vectors.** There are six types:
 - ▷ logical, integer, double, character, complex, and raw
2. **lists.** “sometimes called recursive vectors lists can contain other lists”

Main difference between atomic vectors and lists:

- atomic vectors are “homogenous,” meaning each element in vector must have same type (e.g., integer, logical, character)
- lists are “heterogeneous,” meaning that data type can differ across elements within a list

Link to figure of data structures overview [HERE](#)

Accessing elements of (atomic) vectors

Review: Types of atomic vectors

1. logical. each element can be three potential values: TRUE, FALSE, NA

```
typeof(c(TRUE,FALSE,NA))  
#> [1] "logical"  
typeof(c(1==1,1==2))  
#> [1] "logical"
```

2. Numeric (integer or double)

```
typeof(c(1.5,2,1))  
#> [1] "double"  
typeof(c(1,2,1))  
#> [1] "double"
```

- Numbers are doubles by default. To make integer, place L after number:

```
typeof(c(1L,2L,1L))  
#> [1] "integer"
```

3. character

```
typeof(c("element of character vector","another element"))  
#> [1] "character"  
length(c("element of character vector","another element"))  
#> [1] 2
```


Review: functions that identify type of vector

Function	logical	int	dbl	chr	list
is_logical()	X				
is_integer()		X			
is_double()			X		
is_numeric()		X	X		
is_character()				X	
is_atomic()	X	X	X	X	
is_list()					X
is_vector()	X	X	X	X	X

Recall that elements of a vector must have the same type

- if vector contains elements of different type, the vector type will be the most complex
- from simplest to most complex: logical, integer, double, character

```
is_logical(c(TRUE,TRUE,NA))  
is_logical(c(TRUE,TRUE,NA,1))
```

```
typeof(c(TRUE,1L))  
is_integer(c(TRUE,1L))
```

```
typeof(c(TRUE,1L,1.5,"b"))  
is_character(c(TRUE,1L,1.5,"b"))
```

Review: naming vectors

All vectors can be “named” (i.e., you name individual elements within the vector)

Unnamed vector

```
x <- c(1,2,3,"hi!")
x
#> [1] "1" "2" "3" "hi!"
str(x)
#> chr [1:4] "1" "2" "3" "hi!"
```

named vector

```
y <- c(a=1,b=2,3,c="hi!")
y
#>      a      b      c
#>  "1"  "2"  "3" "hi!"
str(y)
#> Named chr [1:4] "1" "2" "3" "hi!"
#> - attr(*, "names")= chr [1:4] "a" "b" "" "c"
```

Subsetting elements of vector, based on position number

“Subsetting” a vector, refers to isolating particular elements of a vector

- I sometimes refer to this as “accessing elements of a vector”
- subsetting elements of a vector is similar to “filtering” rows of a data-frame

`[]` is the subsetting function for vectors

- contents inside `[]` can refer to element number (also called “position”).
 - ▷ e.g., `[3]` refers to contents of 3rd element (or position 3)
- contents inside `[]` can also refer to name of element
 - ▷ e.g., `["a"]` refers to contents inside an element named “a”

Referring to elements based on position

```
x <- c("a", "b", "c", "d", "e")
```

```
x[1]
```

```
#> [1] "a"
```

```
x[5]
```

```
#> [1] "e"
```

```
c(x[1], x[2], x[2])
```

```
#> [1] "a" "b" "b"
```

```
x[c(1, 2, 2)]
```

```
#> [1] "a" "b" "b"
```

Subsetting elements of vector, based on position number

Referring to elements based on position, continued

```
y <- c(4,5,10,29,15,12)
length(y)
#> [1] 6
```

```
y[c(1,3,6)]
#> [1] 4 10 12
y[c(3,6,1)]
#> [1] 10 12 4
```

While subsetting with positive numbers keeps elements in those positions, subsetting with negative numbers drops elements at those positions

```
y
#> [1] 4 5 10 29 15 12
y[c(-3,-4,-5,-6)]
#> [1] 4 5
```

Subsetting elements of vector, with a logical vector

Grolemund and Wickham (chapter 20): "Subsetting with a logical vector keeps all values corresponding to a TRUE value"

```
x <- c(10, 3, NA, 5, 8, 1, NA, "Hi!")
```

```
typeof(x)
```

```
#> [1] "character"
```

```
x[is.na(x)]
```

```
#> [1] NA NA
```

```
x[!is.na(x)]
```

```
#> [1] "10" "3" "5" "8" "1" "Hi!"
```

```
x[is.numeric(x)]
```

```
#> character(0)
```

```
x[is.character(x)]
```

```
#> [1] "10" "3" NA "5" "8" "1" NA "Hi!"
```

```
y <- c(10, 3, NA, 5, 8, 1, NA)
```

```
typeof(y)
```

```
#> [1] "double"
```

```
y[is.numeric(y)]
```

```
#> [1] 10 3 NA 5 8 1 NA
```

Subsetting elements of named vector, by element name

If you have a named vector, you can subset it with a character vector:

```
x <- c(abc = 1, def = 2, xyz = 5)
x
#> abc def xyz
#>  1  2  5
x[c("xyz", "def")]
#> xyz def
#>  5  2
```

Accessing elements of lists

Review: Lists

Like atomic vectors, lists are objects that contain elements. However, the “type” of elements can vary within a list and elements of a list can contain another list

Examples:

```
rm(list = ls())
x1 <- list(c(1, 2), c(3, 4))
x2 <- list(list(1, 2), list(3, 4))
x3 <- list(1, list(2, list(3)))
```

`str()` function is helpful for understanding structure and contents of a list

```
str(x1)
#> List of 2
#> $ : num [1:2] 1 2
#> $ : num [1:2] 3 4
str(x2)
#> List of 2
#> $ :List of 2
#> ..$ : num 1
#> ..$ : num 2
#> $ :List of 2
#> ..$ : num 3
#> ..$ : num 4
```


Review: Data frames are lists

Recall the relationship between “lists” and “data frames”

- data frames have “type==list”
- data frames are lists with these additional structure requirements
 - ▷ each element of data frame must be a vector (not a list)
 - ▷ each element (i.e., vector) in data frame must have the same length
- data frames have additional attributes
 - ▷ e.g., each vector is named

```
(df <- tibble(x = 1:3, y = 3:1))
#> # A tibble: 3 x 2
#>       x     y
#>   <int> <int>
#> 1     1     3
#> 2     2     2
#> 3     3     1
typeof(df)
#> [1] "list"
str(df)
#> Classes 'tbl_df', 'tbl' and 'data.frame':   3 obs. of  2 variables:
#>  $ x: int  1 2 3
#>  $ y: int  3 2 1

load("../..data/recruiting/recruit_event_somevars.Rdata")
typeof(df_event)
#> [1] "list"
```

Subsetting/accessing elements of a list

Accessing elements of a list important for looping and many applications in R

Will demonstrate accessing elements of a list using two lists:

1. A list that has more complicated structure than a data frame (from Grolemund and Wickham example)

```
list_a <- list(a = 1:3, b = "a string", c = pi, d = list(-1, -5))  
typeof(list_a)  
str(list_a)
```

2. List that is 7 variables and first 5 obs of df_event, corresponding to University of Alabama

```
df_bama <- df_event %>% arrange(univ_id,event_date) %>%  
  select(instnm,univ_id,event_date,event_type,event_state,zip,med_inc) %>%  
  filter(row_number()<6)  
  
typeof(df_bama)  
str(df_bama)
```

Subsetting/accessing elements of a list

Three ways to “subset” (access elements of) a list (from <http://r4ds.had.co.nz/vectors.html#subsetting-1>):

1. `[]` “extracts a sub-list. The result will always be a list”
 - ▶ like subsetting vectors, you can subset with a logical, integer, or character vector
2. `[[]]` “extracts a single component from a list. It removes a level of hierarchy from the list”
3. `$` “shorthand for extracting named elements of a list. It works similarly to `[[]]` except that you don’t need to use quotes.”

Subset a list using []

[“extracts a sub-list”

- contents of `[]` can be position number, name of element in list, logical vect

```
str(list_a)
length(list_a)

str(list_a[1])
str(list_a["a"])

str(list_a[1:2])
str(list_a[c(1,2)])

str(list_a[c("a", "c")])
```

Key takeaway about subsetting a list using []: **The result will always be a list**

- that is, [does not remove a level of hierarchy
- structure and attributes of object you isolate using [will be the same as its structure and attributes in the list it is taken from

Subset a list using `[]`: Student task

Applying `[]` to the object `df_bama`:

- Isolate the 1st element of `df_bama`
- Isolate the 3rd through 5th element of `df_bama`
- Isolate the 3rd, 7th, and 1st element of `df_bama`
- Isolate the element named "event_type"
- Isolate the elements named "event_type" and "med_inc"

Subset a list using []: Student task [SOLUTIONS]

Applying [] to the object df_bama:

```
#- Isolate the 1st element of `df_bama`  
df_bama[1]  
str(df_bama[1])  
#- Isolate the 3rd through 5th element of `df_bama`  
df_bama[3:5]  
str(df_bama[3:5])  
#- Isolate the 3rd, 7th, and 1st element of `df_bama`  
df_bama[c(3,7,1)]  
#- Isolate the element named `"event_type"`  
df_bama["event_type"]  
str(df_bama["event_type"])  
#- Isolate the elements named `"event_type"` and `"med_inc"`  
df_bama[c("event_type", "med_inc")]
```

Subset a list using [[]]

[[]] extracts a single component from a list. It removes a level of hierarchy from the list.

```
str(list_a)

str(list_a[1]) # []
str(list_a[[1]]) # [[]]

str(list_a["a"])
str(list_a[["a"]])

str(list_a[4]) # []
str(list_a[[4]]) # [[]]
```

Subset a list using `[]`, data frames

Comparing `[]` to `[[[]]` when working with lists that are data frames, `df`

- Data frame object always has type=list and each element a vector
- If you subset using `[]` the result will always have type=list
- If you subset using `[[[]]` the result will always have type==vector

```
df_bama[1]
```

```
df_bama[[1]]
```

```
str(df_bama[1])
```

```
str(df_bama[[1]])
```

```
typeof(df_bama[1])
```

```
typeof(df_bama[[1]])
```

```
class(df_bama[1])
```

```
class(df_bama[[1]])
```

```
attributes(df_bama[3])
```

```
attributes(df_bama[[3]])
```

#can perform all of same above tasks with any element, accessed by position number

```
str(df_bama["event_type"])
```

```
str(df_bama[["event_type"]])
```

```
attributes(df_bama["event_type"])
```

```
attributes(df_bama[["event_type"]])
```


Subset a list using \$

\$ is a shorthand for extracting **named** elements of a list. It works similarly to [[]] except that you don't need to use quotes.

- note: we have been using this method of subsetting variables in a data frame all quarter!
- Like [[]], subsetting using \$ removes a level of hierarchy

```
str(list_a)
```

```
list_a$a
```

```
list_a[["a"]]
```

```
df_bama$med_inc
```

```
#> [1] 77380 39134 38272 89203 127972
```

```
df_bama[["med_inc"]]
```

```
#> [1] 77380 39134 38272 89203 127972
```

Key concepts for loops

Review Key concepts for loops

Sequences

(Loose) definition

- a sequence is a list of numbers in ascending or descending order

Creating sequences using colon operator

```
-5:5  
#> [1] -5 -4 -3 -2 -1 0 1 2 3 4 5  
5:-5  
#> [1] 5 4 3 2 1 0 -1 -2 -3 -4 -5
```

Creating sequences using seq() function

- basic syntax:

```
seq(from = 1, to = 1, by = ((to - from)/(length.out - 1)),  
    length.out = NULL, along.with = NULL, ...)
```

- examples:

```
seq(10,15)  
#> [1] 10 11 12 13 14 15  
seq(from=10,to=15,by=1)  
#> [1] 10 11 12 13 14 15  
seq(from=100,to=150,by=10)  
#> [1] 100 110 120 130 140 150
```

Length, vectors

The **length** of an object is its number of elements

Length of vectors, using `length()` function

```
x <- c(1,2,3,4,"ha ha"); length(x)
#> [1] 5
y <- seq(1,10); length(y)
#> [1] 10
z <- c(seq(1,10),"ho ho"); length(z)
#> [1] 11
```

Once vector length known, isolate element contents based on position number using `[`

```
x[5]
#> [1] "ha ha"
z[1]
#> [1] "1"
```

Applying `[[` to vector gives same result as applying `[`

```
x[[5]]
#> [1] "ha ha"
z[[1]]
#> [1] "1"
```

Length of lists

The **length** of an object is its number of elements

```
typeof(df_bama); length(df_bama)
#> [1] "list"
#> [1] 7
```

Once list length known, isolate element contents based on position number using [] or [[]]

- subset one element of list with [] yields list w/ length==1

```
typeof(df_bama[7]); length(df_bama[7])
#> [1] "list"
#> [1] 1
```

- subset one element of list with [[]] yields vector w length==# rows

```
df_bama[[7]]; typeof(df_bama[[7]]); length(df_bama[[7]])
#> [1] 77380 39134 38272 89203 127972
#> [1] "double"
#> [1] 5
```

subset one element of list with \$ is same as [[]]

```
df_bama$med_inc; typeof(df_bama$med_inc); length(df_bama$med_inc)
#> [1] 77380 39134 38272 89203 127972
#> [1] "double"
#> [1] 5
```

Combine sequences and length

When writing loops, very common to create a sequence from 1 to the length (i.e., number of elements) of an object

Here, we do this with a vector object

```
(x <- c("a", "b", "c", "d", "e"))  
#> [1] "a" "b" "c" "d" "e"  
length(x)  
#> [1] 5  
  
1:length(x)  
#> [1] 1 2 3 4 5  
seq(from=1, to=length(x), by=1)  
#> [1] 1 2 3 4 5
```

Can do same thing with list object

```
length(df_bama)  
#> [1] 7  
  
1:length(df_bama)  
#> [1] 1 2 3 4 5 6 7  
seq(2, length(df_bama))  
#> [1] 2 3 4 5 6 7
```

Loop basics

Simple loop example

Loops execute some set of commands multiple times

- we build loops using the `for()` function
- each time the loop executes the set of commands is an **iteration**
- the below loop iterates 4 times

Create loop that prints each value of vector `c(1,2,3,4)`, one at a time

```
c(1,2,3,4)
#> [1] 1 2 3 4

for(i in c(1,2,3,4)) { # Loop sequence
  print(i) # Loop body
}
#> [1] 1
#> [1] 2
#> [1] 3
#> [1] 4
```

Components of a loop

```
for(i in c(1,2,3,4)) { # Loop sequence
  print(i) # Loop body
}
#> [1] 1
#> [1] 2
#> [1] 3
#> [1] 4
```

Components of a loop

1. **Sequence.** Determines what to “loop over” (e.g., from 1 to 4 by 1)
 - ▷ sequence in above loop is `for(i in c(1,2,3,4))`
 - ▷ this creates a temporary object called `i`
 - ▷ each iteration of loop will assign a different value to `i`
 - ▷ `c(1,2,3,4)` is the set of values that will be assigned to `i`
 - in first iteration, value of `i` is 1
 - in second iteration, value of `i` is 2, etc.
2. **Body.** What commands to execute for each iteration through the loop
 - ▷ Body in above loop is `print(i)`
 - ▷ Each time (i.e., iteration) through the loop, body prints the value of object `i`

Components of a loop

These three loops all do the same thing

```
for(z in c(1,2,3,4)) { # Loop sequence
  cat("object z=",z, fill=TRUE) # Loop body
}
#> object z= 1
#> object z= 2
#> object z= 3
#> object z= 4
for(z in 1:4) { # Loop sequence
  cat("object z=",z, fill=TRUE) # Loop body
}
#> object z= 1
#> object z= 2
#> object z= 3
#> object z= 4
num_sequence <- 1:4
for(z in num_sequence) { # Loop sequence
  cat("object z=",z, fill=TRUE) # Loop body
}
#> object z= 1
#> object z= 2
#> object z= 3
#> object z= 4
```

When building loops, I always include a line like `cat("i=",i, fill=TRUE)` to help me understand what loop is doing

Student task

1. Create a numeric vector that has year of birth of members of your family (you decide who to include) e.g., for my mom, dad, wife, son: `birth_years <- c(1944, 1950, 1981, 2016)`
2. Write a loop that calculates current year minus birth year and prints this number for each member of your family

Student task [SOLUTION]

1. Create a numeric vector that has year of birth of members of your family (you decide who to include)
2. Write a loop that calculates current year minus birth year and prints this number for each member of your family

```
birth_years <- c(1944,1950,1981,2016)
birth_years
#> [1] 1944 1950 1981 2016

for(y in birth_years) { # Loop sequence
  cat("object y=",y, fill=TRUE) # Loop body
  z <- 2018-y
  cat("value of",y,"minus",2018,"is",z, fill=TRUE)
}
#> object y= 1944
#> value of 1944 minus 2018 is 74
#> object y= 1950
#> value of 1950 minus 2018 is 68
#> object y= 1981
#> value of 1981 minus 2018 is 37
#> object y= 2016
#> value of 2016 minus 2018 is 2
```

Different ways to loop over a vector

Plan for learning more about loops

Rest of lecture on loops will proceed as follows:

1. Describe the three different ways to “loop over” a vector
2. Describe the two broad sorts of tasks to accomplish within body of a loop
 - 2.1 Modify an existing object (e.g., vector or list/data frame)
 - 2.2 Create a new object

Throughout, I'll try to give you lots of examples and practice

Three ways to loop over an object

There are three ways to loop over elements of an object

1. **Loop over the elements** [approach we have used so far]
2. **Loop over names of the elements**
3. **Loop over numeric indices associated with element position** [approach recommended by Grolemond and Wickham]

Will demonstrate these approaches on a named vector and list/data frame

- o Create named vector

```
vec=c("a"=5, "b"=10, "c"=-5, "d"=30)
vec
#>   a    b    c    d
#>  5  10  -5   30
```

- o Create data frame, with 4 columns of fictitious data

```
set.seed(12345)
df <- tibble(a = rnorm(10), b = rnorm(10), c = rnorm(10), d = rnorm(10))
str(df)
#> Classes 'tbl_df', 'tbl' and 'data.frame':    10 obs. of  4 variables:
#> $ a: num  0.586 0.709 -0.109 -0.453 0.606 ...
#> $ b: num  -0.116 1.817 0.371 0.52 -0.751 ...
#> $ c: num  0.78 1.456 -0.644 -1.553 -1.598 ...
#> $ d: num  0.812 2.197 2.049 1.632 0.254 ...
```


Loop over elements

Approach 1: loop over elements of object

- **sequence** syntax: `for (i in object_name)`
 - ▷ Sequence iterates through each element of the object
 - ▷ that is, sequence iterates through *value* of each element, rather than *name* or *position* of element
- in **body**.
 - ▷ value of `i` is equal to the contents of the `i`th element of the object

Example, object is a vector

```
vec
for (i in vec) {
  cat("\n", "value of object i=", i, fill=TRUE)
  cat("object type=", typeof(i), "; length=", length(i), "; class=", class(i),
      "; attributes=", attributes(i), sep="", fill=TRUE)
}
```

Example, object is a list/data frame

```
for (i in df) {
  cat("\n", "value of object i=", i, fill=TRUE)
  cat("object type=", typeof(i), "; length=", length(i), "; class=", class(i),
      "; attributes=", attributes(i), sep="", fill=TRUE)
}
```

Approach 1: loop over elements of object

Example task:

- calculate mean value of each element of list object `df`

```
for (i in df) { # sequence

  cat("\n", "value of object i=", i, fill=TRUE)
  cat("mean value of object i=", mean(i), fill=TRUE)
}

#>
#> value of object i= 0.5855288 0.709466 -0.1093033 -0.4534972 0.6058875
#> -1.817956 0.6300986 -0.2761841 -0.2841597 -0.919322
#> mean value of object i= -0.1329441
#>
#> value of object i= -0.1162478 1.817312 0.3706279 0.5202165 -0.750532
#> 0.8168998 -0.8863575 -0.3315776 1.120713 0.2987237
#> mean value of object i= 0.2859778
#>
#> value of object i= 0.7796219 1.455785 -0.6443284 -1.553137 -1.59771
#> 1.805098 -0.4816474 0.6203798 0.6121235 -0.162311
#> mean value of object i= 0.08338741
#>
#> value of object i= 0.8118732 2.196834 2.04919 1.632446 0.2542712
#> 0.4911883 -0.3240866 -1.66205 1.767734 0.02580105
#> mean value of object i= 0.72432
```

Loop over element names

Approach 2: loop over names of elements in object

To use this approach, elements in object must have name attributes

- **sequence** syntax: `for (i in names(object_name))`
 - ▷ Sequence iterates through the *name* of each element in object
- in **body**, value of `i` is equal to *name* of *i*th element in object
 - ▷ Access element contents using `object_name[i]`
 - same object type as `object_name`; retains attributes (e.g., *name*)
 - ▷ Access element contents using `object_name[[i]]`
 - removes level of hierarchy, thereby removing attributes
 - Approach recommended by Wickham because isolates value of element

Example, object is a vector

```
vec
for (i in names(vec)) {
  cat("\n", "value of object i=", i, "; type=", typeof(i), sep="", fill=TRUE)
  print(str(vec[i])) # "Access element contents using []"
  print(str(vec[[i]])) # "Access element contents using [[]]"
}
```

Example, object is a list

```
for (i in names(df)) {
  cat("\n", "value of object i=", i, "; type=", typeof(i), sep="", fill=TRUE)
  print(str(df[i])) # "Access element contents using []"
  print(str(df[[i]])) # "Access element contents using [[]]"
}
```

Approach 2: loop over names of elements in object

Example task:

- calculate mean value of each element of list object `df`, using `[[]]` to access element contents

```
for (i in names(df)) {  
  cat("mean of element named",i,"is",mean(df[[i]]),fill=TRUE)  
}  
#> mean of element named a is -0.1329441  
#> mean of element named b is 0.2859778  
#> mean of element named c is 0.08338741  
#> mean of element named d is 0.72432
```

What happens if we try to complete task using `,` using `[]` to access element contents?

```
for (i in names(df)) {  
  cat("mean of element named",i,"is",mean(df[i]),fill=TRUE)  
  #print(typeof(df[i]))  
  #print(class(df[i]))  
}  
#?mean # mean function only works for particular *classes* of objects
```

Loop over element position number

Approach 3: Loop over numeric indices of element position

sequence syntax: `for (i in 1:length(object_name))`

```
length(vec)
1:length(vec)

for (i in 1:length(vec)) {
  cat("value of object i=", i, fill=TRUE)
}
```

Wickham's preferred **sequence** syntax: `for (i in seq_along(object_name))`

- `seq_along(x)` returns a sequence from 1 value of `length(x)`

```
length(vec)
seq_along(vec)

for (i in seq_along(vec)) {
  cat("value of object i=", i, fill=TRUE)
}
```


Approach 3: Loop over numeric indices [SKIP/SKIM]

Why Wickham prefers `seq_along(object_name)` over `1:length(object_name)`

- `seq_along` handles zero-length vectors correctly, and is therefore the “safe” version of `1:length(object_name)`

```
# create vector of length=0
y <- vector("double", 0)
length(y)
#> [1] 0

1:length(y)
#> [1] 1 0
for (i in 1:length(y)) {
  cat("value of object i=", i, fill=TRUE)
}
#> value of object i= 1
#> value of object i= 0

seq_along(y)
#> integer(0)
for (i in seq_along(y)) {
  cat("value of object i=", i, fill=TRUE)
}
```

Personally, I find `1:length(object_name)` much more intuitive

Approach 3: Loop over numeric indices of element position

- **sequence** syntax: `for (i in 1:length(object_name))` OR `for (i in seq_along(object_name))`
 - ▷ Sequence iterates through *position number* of each element in the object
- in **body**, value of `i` equals the *position number* of `i`th element in object
 - ▷ Access element contents using `object_name[i]`
 - same object type as `object_name`; retains attributes (e.g., `name`)
 - ▷ Access element contents using `object_name[[i]]` [RECOMMENDED]
 - removes level of hierarchy, thereby removing attributes

Example, object is a vector

```
vec
for (i in 1:length(vec)) {
  cat("\n", "value of object i=", i, "; type=", typeof(i), sep="", fill=TRUE)
  print(str(vec[i])) # "Access element contents using []"
  print(str(vec[[i]])) # "Access element contents using [[]]"
}
```

Example, object is a list

```
for (i in 1:length(df)) {
  cat("\n", "value of object i=", i, "; type=", typeof(i), sep="", fill=TRUE)
  print(str(df[i])) # "Access element contents using []"
  print(str(df[[i]])) # "Access element contents using [[]]"
}
```

Approach 3: Loop over numeric indices of element position

Example task:

- calculate mean value of each element of list object `df`, using `for (i in seq_along(df))` to create sequence and using `[[i]]` to access element contents

```
for (i in seq_along(df)) {  
  cat("mean of element named",i,"is",mean(df[[i]]),fill=TRUE)  
}  
#> mean of element named 1 is -0.1329441  
#> mean of element named 2 is 0.2859778  
#> mean of element named 3 is 0.08338741  
#> mean of element named 4 is 0.72432
```

What happens if we try to complete task using , using `[]` to access element contents?

```
for (i in seq_along(df)) {  
  cat("mean of element named",i,"is",mean(df[i]),fill=TRUE)  
  #print(typeof(df[i]))  
  #print(class(df[i]))  
}
```

Approach 3: Loop over numeric indices of element position

When looping over numeric indices, you can extract element names based on element position

- First, let's experiment w/ `names()` function

```
attributes(df)
attributes(df[1])
attributes(df[[1]]) # null because removing level of hierarchy removes attribute

names(df)
names(df[1])
names(df[[1]]) # null because object df[[1]] has no attributes

names(df)[[1]] # not null because we extract names of object df and then select
```

- Next, apply what we learned to the loop

```
for (i in seq_along(df)) {
  #print(names(df)[[i]])
  cat("i=", i, "; names=", names(df)[[i]], sep=" ", fill=TRUE)
}

#> i=1; names=a
#> i=2; names=b
#> i=3; names=c
#> i=4; names=d
```

Summary: Three ways to loop over object

1. Loop over elements
2. Loop over element names
3. Loop over numeric indices of element position

Grolemund and Wickham prefer “loop over numeric indices” approach. Why?

“Iteration over the numeric indices is the most general form, because given the position you can extract both the name [approach 2] and the value [approach 1]”

```
for (i in seq_along(df)) {  
  cat("\n", "i=", i, sep="", fill=TRUE)  
  
  name <- names(df)[[i]] # value of object "name" is what we loop over in approach 2  
  cat("name=", name, sep="", fill=TRUE)  
  
  value <- df[[i]] # value of object "value" is what we loop over in approach 1  
  cat("value=", value, sep="", fill=TRUE)  
}
```

Loop tips

When to write a loop

Broadly, the rationale for writing a loop is the same as rationale for writing a function:

- do not duplicate code
- can make changes to code in one place rather than many

When to write a loop:

- Grolemund and Wickham say **don't copy and paste more than twice**; if you find yourself doing this, consider writing a loop or a function

Don't worry about knowing all the situations you should write a loop

- rather, You'll be creating analysis dataset or analyzing data and you will notice there is some task that you are repeating over and over, and then you'll say "oh, I should write a loop or function for this"

When to write a loop vs a functions

Usually it is obvious when you are duplicating, but unclear whether you should write a loop or a whether you should write a function.

- Often, a repeated task can be completed with a loop or a function

In my experience, loops are better for repeated tasks when the individual tasks are **very** similar to one another

- e.g., a loop that reads in data sets from individual years; each dataset you read in differs only by directory and name
- e.g., a loop that converts negative values to NA for a set of variables

Because functions can have many arguments, functions are better when the individual tasks differ substantially from one another

- e.g., a function that runs a regression and spits out a nice table and allows you to specify the dependent variable, the independent variables, what model to run, etc. as function arguments

Note

- you can embed loops within functions; and you can call functions within loops
- But for now, just try to understand basics of functions and loops, and then these things will not seem overwhelming

Recipe for how to write loop

The general recipe for how to write a loop is very similar to the recipe for writing a function:

1. Complete the task for one instance outside a loop (this is akin to writing the **body** of the loop)
2. Write the **sequence**
3. Which parts of the body need to change with each iteration
4. *if* you are creating a new object store output of the loop, create this outside of the loop
5. Construct the loop

Modifying vs. creating new object

Modify object or create new object

Loops can be used to complete all sorts of data manipulation tasks (e.g., read-in data, tidy data, create variables) and data analysis tasks (e.g., run regressions, create plots)

Grolemund and Wickham differentiate between two types of tasks that loops accomplish: Modify an existing object; and create a new object

1. **Modify an existing object**

- ▶ examples: looping through a set of variables in a data frame and modifying these variables or creating new ones
- ▶ When writing loops in Stata/SAS/SPSS, we are usually modifying an existing object because Stata/SAS/SPSS typically only has one object - a dataset - open at a time)

2. Create a new object

- ▶ Example: Create an object that has summary statistics for each variable; this object will be the basis for a table or graph
- ▶ Often the new object will be a vector of results based on looping through elements of a data frame
- ▶ In R (as opposed to Stata/SAS/SPSS) creating a new object is very common because R can hold many objects at the same time

Creating a new object

So far our loops have two components: sequence and body

When we create a new object to store the results of a loop, our loops have three components

1. sequence
2. body
3. output

▸ this is the new object that will store results created from your loop

Grolemund and Wickham recommend creating this new object **prior** to writing the loop (rather than creating the new object within the loop)

“Before you start the loop, you must always allocate sufficient space for the output. This is very important for efficiency: if you grow the for loop at each iteration using `c()` (for example), your for loop will be very slow.”

Creating a new object

Task:

- Using the data frame `df`, which contained data on four numeric variables, create a new object that contains the mean value of each variable

```
set.seed(12345)
df <- tibble(a = rnorm(10), b = rnorm(10), c = rnorm(10), d = rnorm(10))
```

In a previous example, we calculated mean for each variable

```
for (i in seq_along(df)) {
  cat("mean of element named", i, "is", mean(df[[i]]), fill=TRUE)
}
#> mean of element named 1 is -0.1329441
#> mean of element named 2 is 0.2859778
#> mean of element named 3 is 0.08338741
#> mean of element named 4 is 0.72432
```

Now we just have to create an object to store these results

Creating a new object

Task: Create a new object that contains the mean value of each variable in `df`

Wickham recommends creating new object **prior** to creating loop

- You must specify type and length of new object
- New object will contain the mean for each variable, so it should be a numeric vector with number of elements (length) equal to number of variables in `df`

```
output <- vector("double", ncol(df)) # create object
typeof(output)
#> [1] "double"
length(output); length(df)
#> [1] 4
#> [1] 4
```

Now we can create a loop that uses position number to assign variable means to elements of the vector `output`

```
for (i in seq_along(df)) {
  cat("i=", i, fill=TRUE)
  output[[i]] <- mean(df[[i]]) # mean of df[[1]] assigned to output[[1]], etc.
}
#> i= 1
#> i= 2
#> i= 3
#> i= 4
output
#> [1] -0.13294415  0.28597776  0.08338741  0.72432003
```

Example of modifying an object: z-score loop

Task (from Christenson lecture):

- Write a loop that calculates z-score for a set of variables in a data frame and then replaces the original variables with the z-score variables

The z-score for observation i is number of standard deviations from mean:

$$Z_i = \frac{x_i - \bar{x}}{sd(x)}$$

We wrote a z-score function in the functions lecture; this can be basis of our z-score loop

```
z_score <- function(x) {  
  (x - mean(x, na.rm=TRUE))/sd(x, na.rm=TRUE)  
}  
z_score(df$a)  
#> [1] 0.88301870 1.03534018 0.02905509 -0.39396655 0.90803990  
#> [6] -2.07091570 0.93779586 -0.17604498 -0.18584722 -0.96647527  
z_score(df[["a"]]) # same  
#> [1] 0.88301870 1.03534018 0.02905509 -0.39396655 0.90803990  
#> [6] -2.07091570 0.93779586 -0.17604498 -0.18584722 -0.96647527  
z_score(df[[1]]) # same  
#> [1] 0.88301870 1.03534018 0.02905509 -0.39396655 0.90803990  
#> [6] -2.07091570 0.93779586 -0.17604498 -0.18584722 -0.96647527
```

Example of modifying an object: z-score loop

Task (from Christenson lecture):

- Write loop that replaces variables with z-scores of those variables

When modifying existing object, we only need to write **sequence** and **body**

- sequence.** Data frame `df` has 4 variables and all are quantitative, so write a sequence that loops across each element of `df`

▷ `for (i in seq_along(df))`

- Modify **body**.

▷ body of z-score function: $(x - \text{mean}(x, \text{na.rm=TRUE}))/\text{sd}(x, \text{na.rm=TRUE})$

▷ Substitute `df[[i]]` for `x`:

— $(\text{df}[[i]] - \text{mean}(\text{df}[[i]], \text{na.rm=TRUE}))/\text{sd}(\text{df}[[i]], \text{na.rm=TRUE})$

▷ Assign (replace) each observation the value of its z-score:

— `df[[i]] <- (df[[i]] - mean(df[[i]], na.rm=TRUE))/sd(df[[i]], na.rm=TRUE)`

```
set.seed(12345)
df <- tibble(a = rnorm(10), b = rnorm(10), c = rnorm(10), d = rnorm(10))
df
for (i in seq_along(df)) {
  cat("i=", i, "; mean=", mean(df[[i]], na.rm=TRUE), "; sd=", sd(df[[i]], na.rm=TRUE))
  #print((df[[i]] - mean(df[[i]], na.rm=TRUE))/sd(df[[i]], na.rm=TRUE))
  df[[i]] <- (df[[i]] - mean(df[[i]], na.rm=TRUE))/sd(df[[i]], na.rm=TRUE)
}
df
```


Example of modifying an object: z-score loop

What happens if we apply our loop to the data frame `df_bama`, which has both string and numeric variables?

```
for (i in seq_along(df_bama)) {  
  cat("i=", i, "; mean=", mean(df_bama[[i]], na.rm=TRUE), "; sd=", sd(df_bama[[i]], na.rm=TRUE), "\n")  
  #print((df_bama[[i]] - mean(df_bama[[i]], na.rm=TRUE))/sd(df_bama[[i]], na.rm=TRUE))  
  df_bama[[i]] <- (df_bama[[i]] - mean(df_bama[[i]], na.rm=TRUE))/sd(df_bama[[i]], na.rm=TRUE)  
}  
df_bama
```

Let's modify our loop so that it only calculates z-score if for non-integer, numeric variables

```
df_bama  
for (i in seq_along(df_bama)) {  
  cat("i=", i, "; var name=", names(df_bama)[[i]], "; type=", typeof(df_bama[[i]]),  
      "; class=", class(df_bama[[i]]), sep="", fill=TRUE)  
  
  if(is.numeric(df_bama[[i]]) & (!is.integer(df_bama[[i]]))) {  
    df_bama[[i]] <- (df_bama[[i]] - mean(df_bama[[i]], na.rm=TRUE))/sd(df_bama[[i]], na.rm=TRUE)  
  } else {  
    # do nothing  
  }  
}  
df_bama
```

Example of modifying an object: z-score loop

Can we embed this loop in a function that takes the data frame as an argument so we don't have to modify loop for each data frame?

```
z_score <- function(x) {  
  
  for (i in seq_along(x)) {  
    cat("i=", i, "; var name=", names(x)[[i]], "; type=", typeof(x[[i]]),  
        "; class=", class(x[[i]]), sep="", fill=TRUE)  
  
    if(is.numeric(x[[i]]) & (!is.integer(x[[i]]))) {  
      x[[i]] <- (x[[i]] - mean(x[[i]], na.rm=TRUE))/sd(x[[i]], na.rm=TRUE)  
    } else {  
      #do nothing  
    }  
  }  
}  
  
#apply to data frame df  
df_z <- z_score(df)  
df; df_z  
  
#apply to data frame df_bama  
df_bama_z <- z_score(df_bama)  
df_bama; df_bama_z
```

More Practice

Count of events for each university

Count of events for each university

EITHER CUT OR MAKE THIS LATER

This would use the `df_event` dataframe - loop that uses `df_event` and creates event table of number of events for each institution

How well do public universities cover in-state public high schools?

Load recruiting data

Load data frame with one observation per high school and variables for visits by each public research university in sample

- Note: this data frame has more vars than previous data frame we used

```
rm(list = ls()) # remove all objects
load("../data/recruiting/recruit_school_allvars.Rdata")
```

We are interested in creating measures of how good a job public universities are doing visiting in-state public high schools

- Create data frame with one observation for each public high school

```
#names(df_school_all)
df_school_all %>% str()
df_pubhs <- df_school_all %>% # Create data-frame that keeps public high schools
  filter(school_type=="public") %>% select(-school_type)
rm(df_school_all)
```

Create standalone objects (output and code omitted)

1. character vector containing ID for each public university
2. A named list containing university name

How well do public universities cover in-state public high schools

Task: for each public research university, calculate the number and percent of public high schools in the university's home state that received a visit

First, let's accomplish task outside of a loop for one university [Tidyverse]

- let's choose "U of South Carolina", ID==218663

```
#"state_code" is the 2-letter high school state code
```

```
df_pubhs %>% select(state_code) %>% str()
```

```
#variables starting with "inst_" identify state the university is located in
```

```
df_pubhs %>% select(inst_218663) %>% str()
```

```
df_pubhs %>% select(inst_218663) %>% count(inst_218663) # these vars don't vary
```

```
#variables starting with "visits_by_" indicate number of visits HS got in 2017
```

```
df_pubhs %>% select(visits_by_218663) %>% str()
```

```
df_pubhs %>% select(visits_by_218663) %>% count(visits_by_218663)
```

```
#filter only obs where HS state code equals home state of university
```

```
df_pubhs %>% filter(state_code==inst_218663) %>% count() # count pub HS in SC
```

```
#Create measures: number pub HS in SC; number w/ visit; pct w/ visit
```

```
df_pubhs %>% filter(state_code==inst_218663) %>% select(visits_by_218663) %>%
```

```
  mutate(got_visit=ifelse(visits_by_218663>0,1,0)) %>%
```

```
  summarise(n_hs=n(),n_visit=sum(got_visit),pct_visit=sum(got_visit)/n())
```


How well do public universities cover in-state public high schools

Task: for each public research university, calculate the number and percent of public high schools in the university's home state that received a visit

First, let's accomplish task outside of a loop for one university [Base R]

- let's choose "U of South Carolina", ID==218663

```
#"state_code" is the 2-letter high school state code
```

```
str(df_pubhs$state_code)
```

```
#variables starting with "inst_" identify state the university is located in
```

```
str(df_pubhs$inst_218663)
```

```
table(df_pubhs$inst_218663, useNA='ifany') # these vars don't vary
```

```
#variables starting with "visits_by_" indicate number of visits HS got in 2017
```

```
str(df_pubhs$visits_by_218663)
```

```
table(df_pubhs$visits_by_218663, useNA='ifany')
```

```
#filter only obs where HS state code equals home state of university
```

```
tempdf <- subset(df_pubhs,df_pubhs[["state_code"]]==df_pubhs[["inst_218663"]])
```

```
#tempdf <- subset(df_pubhs,df_pubhs$state_code==df_pubhs$inst_218663) # same as above
```

```
#tempdf <- subset(df_pubhs,state_code==inst_218663) # same as above
```

```
#Create 0/1 indicator of whether got visit
```

```
tempdf$got_visit <- ifelse(tempdf$visits_by_218663>0,1,0)
```

```
#frequency count of schools that got visits vs. not
```

```
table(tempdf$got_visit, useNA='ifany')
```

How well do public universities cover in-state public high schools

Task: for each public research university, calculate the number and percent of public high schools in the university's home state that received a visit

Build loop [Base R approach]

- first, loop through each value of list `instnm`

```
instnm
for (i in seq_along(instnm)) {
  cat("\n", "i=", i, sep="", fill=TRUE)

  name <- names(instnm)[[i]] # name of element
  cat("name=", name, sep="", fill=TRUE)

  value <- instnm[[i]] # value of element
  cat("value=", value, sep="", fill=TRUE)
}
```

How well do public universities cover in-state public high schools

Task: for each public research university, calculate the number and percent of public high schools in the university's home state that received a visit

Build loop

- next, create "inst_..." and "visits_by_..." vars for each id
- keep obs in same state as university
- create 0/1 variable of whether high school got a visit

```
for (i in seq_along(instnm)) {  
  cat("\n", "i=", i, "; ", names(instnm)[[i]], sep="", fill=TRUE)  
  
  #create object called inst_var; value is "inst_id" (e.g., "inst_166629")  
  cat("inst_", instnm[[i]], sep="", fill=TRUE)  
  inst_var <- paste("inst_", instnm[[i]], sep="")  
  print(inst_var)  
  
  #create object called visits_by_var; value is "visits_by_id" (e.g., "visits_by_  
  visits_by_var <- paste("visits_by_", instnm[[i]], sep="")  
  print(visits_by_var)  
  
  #create subset of data with high schools in same state as the university  
  tempdf <- subset(df_pubhs, df_pubhs[["state_code"]] == df_pubhs[[inst_var]])  
  # code df_pubhs[[inst_var]] evaluates to df_pubhs[["inst_166629"]] or wh  
  # this is same as instnm[[i]] evaluating to instnm[[16]] or whatever curre  
  #Create 0/1 indicator of whether got visit  
  tempdf$got_visit <- ifelse(tempdf[[visits_by_var]] > 0, 1, 0)  
}
```

How well do public universities cover in-state public high schools

Task: for each public research university, calculate the number and percent of public high schools in the university's home state that received a visit

Build loop

- next, create count of number of visited and non-visited in-state schools

```
for (i in seq_along(instnm)) {  
  cat("\n", "i=", i, "; ", names(instnm)[[i]], sep="", fill=TRUE)  
  
  inst_var <- paste("inst_", instnm[[i]], sep="")  
  visits_by_var <- paste("visits_by_", instnm[[i]], sep="")  
  
  tempdf <- subset(df_pubhs, df_pubhs[["state_code"]] == df_pubhs[[inst_var]]) # keep  
  tempdf$got_visit <- ifelse(tempdf[[visits_by_var]] > 0, 1, 0) # create 0/1 indicator  
  
  # create frequency table of number of schools with and without visits  
  print(table(tempdf$got_visit, useNA='ifany'))  
  ct_table <- table(tempdf$got_visit, useNA='ifany') # named vector with 2 elements  
  
  # create proportion table  
  print(prop.table(ct_table))  
  pr_table <- prop.table(ct_table) # named vector with 2 elements str(pr_table)  
}
```

How well do public universities cover in-state public high schools

Task: for each public research university, calculate the number and percent of public high schools in the university's home state that received a visit

Here is tidyverse approach to loop, which uses some programming concepts we haven't covered

```
for (i in seq_along(instnm)) {  
  cat("\n", "i=", i, "; ", names(instnm)[[i]], sep="", fill=TRUE)  
  
  #create object called inst_var; value is "inst_id" (e.g., "inst_166629")  
  inst_var <- paste("inst_", instnm[[i]], sep="")  
  
  #create object called visits_by_var; value is "visits_by_id" (e.g., "visits_by_  
  visits_by_var <- paste("visits_by_", instnm[[i]], sep="")  
  
  #Create measures: number pub HS in SC; number w/ visit; pct w/ visit  
  df_pubhs %>% filter_(glue::glue("state_code=={inst_var}")) %>%  
    select_(visits_by_var) %>%  
    mutate_(got_visit=glue::glue("ifelse({visits_by_var}>0,1,0)) %>%  
    summarise(n_hs=n(), n_visit=sum(got_visit), pct_visit=sum(got_visit)/n())  
}
```