

# Lab 1, Fall 2020

Your Friendly GEs!

10/5/2020

## Introduction:

Hi everybody, I'm your friendly neighborhood GE and I will be helping you get through 421!

I have outlined my goals for our 90min lab today below

- Introduce myself, take a few polls, iron out zoom setting, talk about the process of learning and remembering syntax
- Get R/RStudio downloaded and working for most. (I realize that some may already have it installed while others may have some issues - I have reserved the last 30min of the class to help individuals)
- Walkthrough of R in RStudio

## Lesson 0: Basics

In general:

- We will be working with what is called a *script*. This is similar to a do file in Stata. It's basically your workspace.
- To execute a script, hit control+enter (cmd+enter if on Mac.) To save the script, control+s. There are other shortcuts as well. If you want to run a specific line, then you can move your cursor to that line and hit control+return and the R script will only run that one line.
- R uses *functions*, which we apply to *objects*. More on this shortly, but if you aren't sure what a function does, or how it works, you can use `?`  before the function to get the documentation. Ex:

```
?mean
```

```
## starting httpd help server ... done
```

In rstudio, this will pull up the function documentation in the help window. Learning to read the R documentation is a brutally necessary skill to becoming proficient in R. It can be confusing and frustrating but the more you do it the better R programmer you will become. Also, google is your best friend. I learn by split screening my script and a web-browser.

There are a ton of different types of objects (numeric (numbers), character (letters) and logical (true false statements) are the most common types), and not all functions will work on all objects. Let's talk a bit about objects.

## Lesson 1: All things are objects

An object is an assignment between a name and a value. You assign values to names using `<-` or `=`. The first assignment symbol consists of a `<` next to a dash `-` to make it look like an arrow.

If we want to make an object name 'a' refer to the number '2', we can do that by:

```
# assign the value '2' to the name 'a'  
a <- 2
```

We can do the same thing using the 'equals' sign

```
a = 2  
a
```

```
## [1] 2
```

When I run this code, we see an output that looks like `[1] 2`. The `[1]` refers to the output line. There is only one line here, as we only called one object. The 2 is the value associated with our object. In this case, a, which we set equal to 2.

You can combine objects together as well which lets us do some basic math operations

If you want to print intermediate steps of your code (ie, see them show up in your terminal) you can put parentheses around your code to get them to display

Let's find the value of `2*3` (equal to two times three), which should be equal to 6. Since a is already equal to 2, we can use that and add it to another variable. Of course, we could simply type `2*3`, but this way let's me show things off.

```
#assign the value of 3 to the name b  
b <- 3  
#assign the value of b (3) times the value of a (2), to a new name, c (now 6). Remember, parentheses wi  
(c <- a * b)
```

```
## [1] 6
```

```
#display c  
c
```

```
## [1] 6
```

objects however can hold more than one thing inside of them. For instance, we can make vectors. R also has a cool method to create a vector of integers with the colon operator.

```
#let's create a vector with the integers from 1 through 10.  
tmp <- 1:10  
  
tmp
```

```
## [1] 1 2 3 4 5 6 7 8 9 10
```

We can also create vectors with the following syntax:

```
avect <- c(1,4,6,-10,-5.25)
```

```
(avect)
```

```
## [1] 1.00 4.00 6.00 -10.00 -5.25
```

Note: In this syntax, `c()` stands for concatenate. Definition of concatenate: operation of joining character strings

we can even apply mathematical operators to vectors. Let's square this vector (this will square all values)

```
(avect^2)
```

```
## [1] 1.0000 16.0000 36.0000 100.0000 27.5625
```

one caveat though... R has two tendencies/rules to keep in mind. One, vectors have to all be the same 'type' of object, and two, R tries to help a little too much.

Let's make a 'bad' vector: we'll add to our temporary vector a few 'character' values, and one 'Null' value. We will work with this guy later.

```
#make a vector that doesn't behave itself  
bad <- c(avect, 11, "red", "dinosaurs", NULL)
```

before we look into what goes wrong in 'bad' let's learn a bit about functions. We've actually worked with a function already! The 'c' is actually a function. Let's do a quick overview of how functions work in R.

## Lesson 2: Functions

Functions are operations that can transform your created **object** in a ton of different ways. Let's look at some convenient ways to get a snapshot of the data, and summary statistics

Examples: `head`, `tail`, `mean`, `median`, `sd`, `summary`, `min`, `max`

These functions are good at summarizing data in a variety of ways. Let's see how they work

```
#print the first few objects in 'tmp'  
head(tmp)
```

```
## [1] 1 2 3 4 5 6
```

```
#print the first 3 objects in 'tmp'  
head(tmp, 3)
```

```
## [1] 1 2 3
```

```
#print the last few objects in 'tmp'  
tail(tmp)
```

```
## [1] 5 6 7 8 9 10
```

Now let's look at what happened to our 'bad' vector.

```
#let's print the last 6 objects in 'bad'  
tail(bad,6)
```

```
## [1] "6"          "-10"          "-5.25"        "11"          "red"          "dinosaurs"
```

Interesting. What is going on here? Why does our vector look different?

We can also use these to perform some basic or commonly used statistics, without the hassle of typing in the formula explicitly. These functions are built into R.

```
#mean of our vector tmp  
mean(tmp)
```

```
## [1] 5.5
```

```
#median of our vector  
median(tmp)
```

```
## [1] 5.5
```

```
#standard deviation of our vector  
sd(tmp)
```

```
## [1] 3.02765
```

*IMPORTANT* We can also print a **summary** of our object.

```
#This can work on many object types and is useful to get an overview of the object in general  
summary(tmp)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##      1.00   3.25   5.50   5.50   7.75   10.00
```

## Lesson 3: types of Objects

So far, we've mostly worked with numeric objects, which can be integer or numeric. you can detect what sort of object you're working with using the `class()` function.

```
(class(a))
```

```
## [1] "numeric"
```

```
(class(tmp))
```

```
## [1] "integer"
```

```
(class(avect))
```

```
## [1] "numeric"
```

another common class is called `character`

Let's create a character object. These are surrounded by either `"` or `'`. This distinguishes them from *object names*.

```
(some_text <- "calvin go do something you hate, it builds character")
```

```
## [1] "calvin go do something you hate, it builds character"
```

```
class(some_text)
```

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

with that in mind, let's look at `'bad'` again. Remember, it has the numbers 1-10 at first, so we'd expect those values to be stored as integers.

```
#we can use the 'head' function to look at the first value of our 'bad' vector. Remember, it should be  
class(head(bad,1))
```

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

huh, it seems like R tried to be helpful and make all of your objects into character values. This is important to keep in mind as you start to work with data.

The last type of object you'll work with frequently is a dataframe. However, to help you get going on the homework, I'm going to do this primarily with your own homework dataset, which requires us to start loading packages.

## Lesson 4: Load Packages & HW help

Base R (what comes installed on your computer) is an incredibly powerful programming language, but one of the best features of R are its packages, which are remotely stored functions written by anybody. You could even write a package if you wanted! This open source nature allows R to be extremely flexible. For now, we will load the `pacman` package management package, and then the `broom` and `tidyverse` packages which gives us access to the `tidy()` command (useful for summarizing regression objects) and a number of useful data manipulation tools.

Let's start by loading packages. Uncomment the `install.packages` function to get the `pacman` package to install. If you already have some of these packages, feel free to delete lines. The `install.packages` function can take a vector of package names, as characters, to install all of the above.

For the most part, we will use the super helpful package `'pacman'` to load new packages into our workspace

```
install.packages(c("pacman"), dependencies=T, repos = "http://cran.us.r-project.org")
```

Congrats! The `pacman` package (package management... get it?) will allow us to install, update and load packages in such a way that you won't have to worry about conflicts. Load `pacman` with the `library()` function.

```
library(pacman)
#p_load is pacman's 'library' in that it lets you load packages. and features a number of improvements.
p_load(ISLR, tidyverse, broom)
```

We're also going to use a new function that lets us read CSV data into our workspace. `read_csv`, which comes from a package we just loaded called 'tidyverse'.

In order to use this function effectively, we should create a name for our dataset. To keep things easy for you, let's call this `your_df`, but really you could call this anything. In order to load a CSV, you need to find the **filepath** for your dataset. Go to Canvas to download the file, and then find it in your downloads folder.

If you're on **Mac**, you can find your filepath by right clicking on the file, and then holding control and selecting the 'copy as path' option.

If you're on **PC**, you can find your filepath by right clicking the file, going to 'properties' and copying the filepath from there. Unfortunately, the slashes on PC are also what are known as 'escape characters' in R. What this means for you is either you must replace every "`\`" with "`\\`" or with "`/`".

Turn this filepath into a string by placing quotes around it, and then pass it to the `read_csv()` command.

```
#remember, your filepath needs to be passed to read_csv() as a string, meaning "/Users/You/..." not /Us
your_df <-read_csv(filepath)
```

```
## Parsed with column specification:
## cols(
##   first_name = col_character(),
##   sex = col_character(),
##   i_callback = col_double(),
##   n_jobs = col_double(),
##   n_expr = col_double(),
##   i_military = col_double(),
##   i_computer = col_double(),
##   i_female = col_double(),
##   i_male = col_double(),
##   race = col_character(),
##   i_black = col_double(),
##   i_white = col_double()
## )
```

you'll see that R helpfully tells you what each class is for the variables in your dataset. However, the first thing to do is always to get a good summary- so let's use the summary command

```
summary(your_df)
```

```
##   first_name          sex          i_callback          n_jobs
## Length:4864      Length:4864      Min.   :0.00000      Min.   :1.000
## Class :character  Class :character  1st Qu.:0.00000      1st Qu.:3.000
## Mode  :character  Mode  :character  Median :0.00000      Median :4.000
##                                     Mean  :0.08059      Mean   :3.663
##                                     3rd Qu.:0.00000      3rd Qu.:4.000
##                                     Max.   :1.00000      Max.   :7.000
##   n_expr          i_military          i_computer          i_female
## Min.   : 1.000      Min.   :0.00000      Min.   :0.00000      Min.   :0.00000
## 1st Qu.: 5.000      1st Qu.:0.00000      1st Qu.:1.00000      1st Qu.:1.00000
```

```
## Median : 6.000 Median :0.00000 Median :1.0000 Median :1.0000
## Mean : 7.842 Mean :0.09704 Mean :0.8205 Mean :0.7691
## 3rd Qu.: 9.000 3rd Qu.:0.00000 3rd Qu.:1.0000 3rd Qu.:1.0000
## Max. :44.000 Max. :1.00000 Max. :1.0000 Max. :1.0000
## i_male race i_black i_white
## Min. :0.0000 Length:4864 Min. :0.0000 Min. :0.0000
## 1st Qu.:0.0000 Class :character 1st Qu.:0.0000 1st Qu.:0.0000
## Median :0.0000 Mode :character Median :1.0000 Median :0.0000
## Mean :0.2309 Mean :0.5002 Mean :0.4998
## 3rd Qu.:0.0000 3rd Qu.:1.0000 3rd Qu.:1.0000
## Max. :1.0000 Max. :1.0000 Max. :1.0000
```

you'll see that for the numeric columns, you'll be looking at summary statistics, but for character (word) columns you only get the length and the Class/Mode.

we can use our snapshot tools here in the same way we did with vectors to take a peek at this df

```
head(your_df,8)
```

```
## # A tibble: 8 x 12
## first_name sex i_callback n_jobs n_expr i_military i_computer i_female
## <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Jill f 0 2 5 0 1 1
## 2 Kenya f 0 4 21 0 1 1
## 3 Latonya f 0 3 3 0 1 1
## 4 Tyrone m 0 2 6 0 0 0
## 5 Aisha f 0 4 8 0 1 1
## 6 Allison f 0 4 8 0 1 1
## 7 Aisha f 0 4 4 0 1 1
## 8 Carrie f 0 2 4 0 1 1
## # ... with 4 more variables: i_male <dbl>, race <chr>, i_black <dbl>,
## # i_white <dbl>
```

```
tail(your_df, 18)
```

```
## # A tibble: 18 x 12
## first_name sex i_callback n_jobs n_expr i_military i_computer i_female
## <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Meredith f 0 5 13 0 1 1
## 2 Tanisha f 0 7 18 0 1 1
## 3 Geoffrey m 0 6 8 0 1 0
## 4 Greg m 0 3 7 0 0 0
## 5 Jamal m 0 4 6 0 1 0
## 6 Tamika f 0 4 2 1 1 1
## 7 Jamal m 0 4 2 1 1 0
## 8 Latonya f 1 4 6 0 1 1
## 9 Matthew m 0 6 8 0 1 0
## 10 Sarah f 1 3 7 0 0 1
## 11 Allison f 0 5 13 0 1 1
## 12 Jill f 0 4 16 0 1 1
## 13 Lakisha f 0 5 26 0 1 1
## 14 Tamika f 0 2 1 0 1 1
## 15 Ebony f 0 4 6 0 1 1
```

```
## 16 Jay      m      0      6      8      0      1      0
## 17 Latonya  f      0      4      2      1      1      1
## 18 Laurie   f      0      3      7      0      0      1
## # ... with 4 more variables: i_male <dbl>, race <chr>, i_black <dbl>,
## #   i_white <dbl>
```

```
glimpse(your_df)
```

```
## Rows: 4,864
## Columns: 12
## $ first_name <chr> "Jill", "Kenya", "Latonya", "Tyrone", "Aisha", "Allison"...
## $ sex <chr> "f", "f", "f", "m", "f", "f", "f", "f", "f", "m", "m", "...
## $ i_callback <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
## $ n_jobs <dbl> 2, 4, 3, 2, 4, 4, 4, 2, 2, 3, 3, 3, 2, 2, 3, 3, 2, 3, 3,...
## $ n_expr <dbl> 5, 21, 3, 6, 8, 8, 4, 4, 5, 4, 5, 6, 6, 8, 4, 3, 2, 7, 3...
## $ i_military <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
## $ i_computer <dbl> 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 1,...
## $ i_female <dbl> 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0,...
## $ i_male <dbl> 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1,...
## $ race <chr> "w", "b", "b", "b", "b", "w", "b", "w", "b", "w", "w", "w", "...
## $ i_black <dbl> 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0,...
## $ i_white <dbl> 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1,...
```

We can also check the size of the dataframe out by using the `dim()` command

```
dim(your_df)
```

```
## [1] 4864 12
```

If we're interested in looking at a specific subset of the data, we can do this in one of two ways. The first is by referencing the column name. Let's look at column 'race'

```
head(your_df$race,4)
```

```
## [1] "w" "b" "b" "b"
```

if you're interested in many columns, you can find them by using what's called "indexing". Indexing is passed to a set of square brackets, and is labeled row, column

```
# let's look at columns 2 & 3
```

```
head(your_df[,c(2,3)])
```

```
## # A tibble: 6 x 2
##   sex  i_callback
##   <chr>      <dbl>
## 1 f      0
## 2 f      0
## 3 f      0
## 4 m      0
## 5 f      0
## 6 f      0
```



```
# Or
```

```
your_df %>% select(sex, i_callback) %>% head(10)
```

```
## # A tibble: 10 x 2
##   sex    i_callback
##   <chr>      <dbl>
## 1 f          0
## 2 f          0
## 3 f          0
## 4 m          0
## 5 f          0
## 6 f          0
## 7 f          0
## 8 f          0
## 9 f          0
## 10 m         0
```

```
#this creates a mini-dataframe with what you selected. You can also do this for a single column like be
```

```
head(your_df[,3])
```

```
## # A tibble: 6 x 1
##   i_callback
##   <dbl>
## 1          0
## 2          0
## 3          0
## 4          0
## 5          0
## 6          0
```

```
# Or
```

```
i_callback_vec <- your_df %>% select(i_callback) %>% head()
i_callback_vec
```

```
## # A tibble: 6 x 1
##   i_callback
##   <dbl>
## 1          0
## 2          0
## 3          0
## 4          0
## 5          0
## 6          0
```

If you want to see the names of the columns you were looking at, you can use

```
names(your_df[,3])
```

```
## [1] "i_callback"
```

To print them out.

you can also select rows with indexing. Let's look at rows 5-10 of columns 2-6

```
your_df[5:10,6:12]
```

```
## # A tibble: 6 x 7
##   i_military i_computer i_female i_male race  i_black i_white
##   <dbl>      <dbl>    <dbl> <dbl> <chr>   <dbl>   <dbl>
## 1         0         1        1     0 b         1         0
## 2         0         1        1     0 w         0         1
## 3         0         1        1     0 b         1         0
## 4         0         1        1     0 w         0         1
## 5         0         0        1     0 b         1         0
## 6         0         1        0     1 w         0         1
```

or rows 2,7,10 and 12 of those same columns

```
your_df[c(2,7,10,12),2:6]
```

```
## # A tibble: 4 x 5
##   sex  i_callback n_jobs n_expr i_military
##   <chr>      <dbl>  <dbl>  <dbl>    <dbl>
## 1 f         0      4     21         0
## 2 f         0      4      4         0
## 3 m         0      3      4         0
## 4 f         0      3      6         0
```

We can also use other features of the tidyverse to look at specific portions of our data. Let's use the filter command. The way this works is:

```
#let's filter on whether or not race is equal to white. In this dataset, that is represented by 'w'
filter(your_df, race == 'w')
```

```
## # A tibble: 2,431 x 12
##   first_name sex  i_callback n_jobs n_expr i_military i_computer i_female
##   <chr>      <chr>      <dbl>  <dbl>  <dbl>    <dbl>      <dbl>    <dbl>
## 1 Jill      f         0      2      5         0         1         1
## 2 Allison   f         0      4      8         0         1         1
## 3 Carrie    f         0      2      4         0         1         1
## 4 Geoffrey  m         0      3      4         0         1         0
## 5 Matthew   m         0      3      5         0         0         0
## 6 Jill      f         0      2      6         0         0         1
## 7 Todd      m         0      3      3         0         1         0
## 8 Allison   f         0      2      2         0         0         1
## 9 Carrie    f         0      3      7         0         1         1
## 10 Greg      m         0      3      3         0         1         0
## # ... with 2,421 more rows, and 4 more variables: i_male <dbl>, race <chr>,
## #   i_black <dbl>, i_white <dbl>
```

This lets us find things like 'what percentage of white respondents were women'

```
#the mean of a binary variable is the percentage of times that binary variable is equal to '1'  
(mean(filter(your_df, race == 'w')$i_female))
```

```
## [1] 0.7638832
```

looks like it's about 76% of the white respondents are women.

we can also do this type of operation in steps using pipes (which come from the tidyverse package), using this funny carrot symbol %>%

let's look at the first 8 rows and find out the first names of our respondents using a tidyverse command called `select`

```
head(your_df) %>% select(first_name)
```

```
## # A tibble: 6 x 1  
##   first_name  
##   <chr>  
## 1 Jill  
## 2 Kenya  
## 3 Latonya  
## 4 Tyrone  
## 5 Aisha  
## 6 Allison
```

Now let's look and see if there is any gender disparity between our two groups by doing a **z-test**. What do we need to generate a z-score again? We need the n for each of our groups, a mean for each group, and a mean for females in the sample overall.

- number of black people in sample

we can use the `nrow` command to count the number of rows in the dataset after you filter by race = 'b'

```
n_b <- filter(your_df, race == "b") %>% nrow()  
n_b
```

```
## [1] 2433
```

- number of white people

```
n_w <- filter(your_df, race == "w") %>% nrow()  
n_w
```

```
## [1] 2431
```

- the means

```
#percentage of females who are black
mean_b <- filter(your_df, race == "b")$i_female %>% mean()
#percentage of females who are white
mean_w <- filter(your_df, race == "w")$i_female %>% mean()
#percentage of females in the sample overall
mean_all <- your_df$i_female %>% mean()
```

and we need to put them all together. We're going to be using a few mathematical operators here, primarily, `sqrt()` (square root) and `pnorm()` (distribution function of normal).

`pnorm` is a weird one, and it takes a quantile (so a z-score), mean, standard deviation, and a few other things. We're interested in a two sided z-test so we'll use `lower.tail` equal to `false`.

```
#use the formula for a z-score to calculate the z-stat
z_stat <- (mean_b - mean_w) / sqrt(mean_all*(1-mean_all)*(1/n_b + 1/n_w))
z_stat
```

```
## [1] 0.8663662
```

```
#plug it into the probability of seeing a value as large as our z-stat
2*pnorm(abs(z_stat), lower.tail = FALSE)
```

```
## [1] 0.3862894
```

this doesn't look too bad. There's a difference here but not statistically significant. Let's look at this problem another way.

## Lesson 5: Regression in R

we can also run regressions in R using the 'lm' command. Let's run a regression testing if female is significant by race.

```
#equations are a NEW object type and are denoted usually by the tilde (~) object
your_reg <- lm(i_female ~ i_black, data = your_df)
```

You can see here that `lm` is a function that takes an "equation" object of the following form:

`y ~ x1 + x2 + x3 + etc.`

It also needs to be able to reference the dataframe in question, so we need to tell it the name our dataframe has. In this case, `your_df`.

we can pass this to 'tidy', which comes from the broom package, to display the results nicely, we also can use `summary()`

```
your_reg %>% tidy()
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)    0.764    0.00855   89.4      0
## 2 i_black        0.0105   0.0121    0.866    0.386
```

Now for summary:

```
summary(your_reg)
```

```
##
## Call:
## lm(formula = i_female ~ i_black, data = your_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.7743  0.2256  0.2256  0.2361  0.2361
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.763883   0.008548  89.366  <2e-16 ***
## i_black      0.010469   0.012086   0.866   0.386
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4214 on 4862 degrees of freedom
## Multiple R-squared:  0.0001543, Adjusted R-squared: -5.133e-05
## F-statistic: 0.7504 on 1 and 4862 DF, p-value: 0.3864
```

let's look at a more complicated equation involving interaction terms. We can call out interaction terms in R by putting a colon between the two interacted variables

```
your_reg_comp <- lm(i_female ~ i_black + i_callback + i_black:i_military, data = your_df)
```

```
summary(your_reg_comp)
```

```
##
## Call:
## lm(formula = i_female ~ i_black + i_callback + i_black:i_military,
##     data = your_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.8077  0.1923  0.2118  0.2380  0.3597
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.76199   0.00879  86.687  < 2e-16 ***
## i_black        0.02617   0.01241   2.109   0.035 *
## i_callback     0.01956   0.02218   0.882   0.378
## i_black:i_military -0.14782  0.02817 -5.247 1.61e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4203 on 4860 degrees of freedom
## Multiple R-squared:  0.005988, Adjusted R-squared:  0.005374
## F-statistic: 9.759 on 3 and 4860 DF, p-value: 2.041e-06
```

PHEW. That's it. I will see you guys next week! Thank you!