

# Introduction to R II

## R for Advanced Stata Users

---

DIME Analytics

The World Bank | [WB Github](#)

November 2020



# Table of contents



1. Introduction
2. Initial settings
3. File paths
4. Exploring a dataset
5. Commenting
6. Creating a document outline in RStudio
7. Using packages
8. Functions inception
9. Looping
10. Custom functions
11. Indentation
12. Appendix

# Introduction

---



## Why are we here today?

- In the last session, you learned the basic concepts to work in R
- You are probably eager to get your hands into some data using R by now, and you would figure out **what should be in your code** for it to work
- But you would probably not know right away **how to write that**, so that in the end you might have code that is only intelligible for yourself -- and not for a very long time



## Why are we here today?

- In this session, we will cover common coding practices in R so that you can make **the most efficient use** for it
- We will also discuss some styling conventions to make your code **readable and reproducible**
- This will give you a solid foundation to code in R, and hopefully you'll be able to skip some painful steps of the "getting-your-hands-dirty" learning approach

# Initial settings

---

# Initial settings



- Ok, let's start by opening RStudio, or by reopening it if you had it opened already
- What do you see in your environment?
- If you saved the last RStudio session in `.Rhistory` (and that is the default), the objects that were in RStudio's memory last time you closed it will still be there whenever you open it again

# Initial settings



Raise your hand if you have ever seen these lines of code before:

A screenshot of a Stata Do-file Editor window titled "Do-file Editor - MyStataDoFile.do". The window has a menu bar with "File", "Edit", "View", "Project", and "Tools". Below the menu bar is a toolbar with various icons for file operations and editing. The main editing area shows a file named "MyStataDoFile.do" with a list of Stata commands:

```
1  
2      clear  
3      set memory 200m  
4      set maxvar 120000  
5      set more off  
6      set varabbrev off  
7
```



# Initial settings



- We **don't need to set the memory or the maximum number of variables in R**, and the equivalent of `set more off` is the default in R
- You can see all the objects currently in you memory in the *Environment* pane



## Exercise 1: Clear workspace

- Make sure the `Environment` window is open. Is anything there?
- Create an object called `foo` with any content you pick
- Type `rm(foo)` to remove the `foo` object from you memory
- Type `ls()` to print the names of the object in memory
- To remove all objects, use `rm(list=ls())`



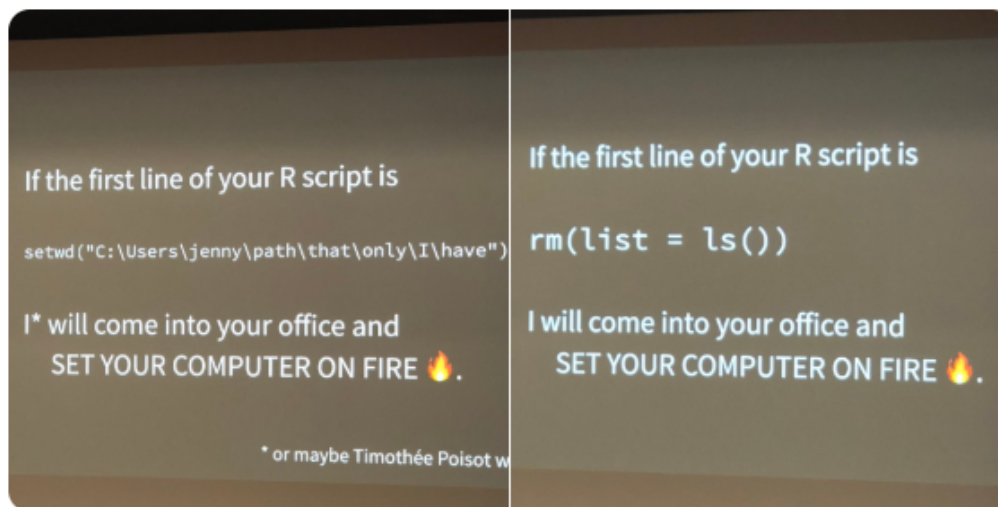
**Hadley Wickham** ✓

@hadleywickham

Following



The only two things that make @JennyBryan  
😞😡😱. Instead use projects + here::here()  
[#rstats](#)



4:50 PM - 10 Dec 2017

303 Retweets 992 Likes



64



303



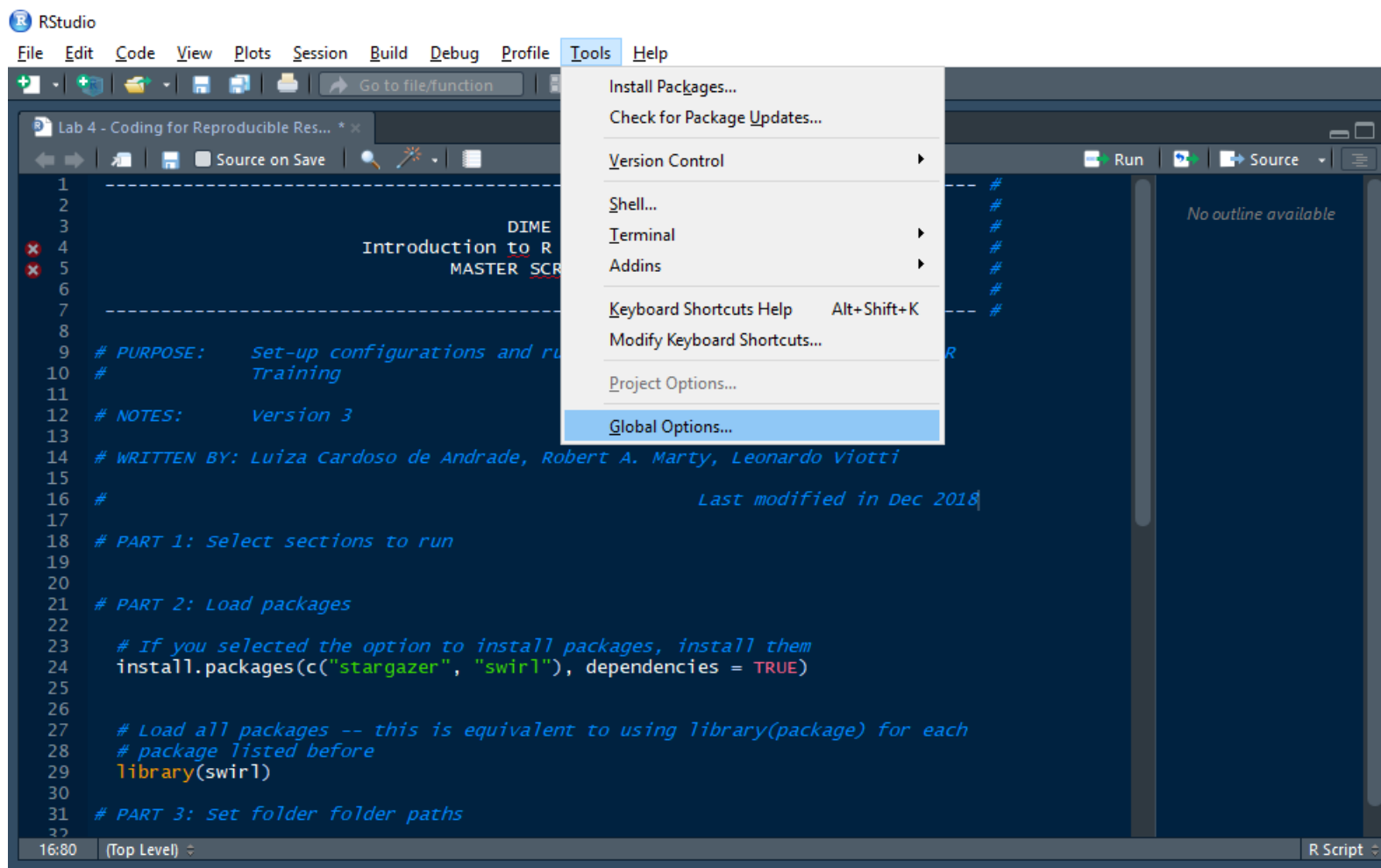
992



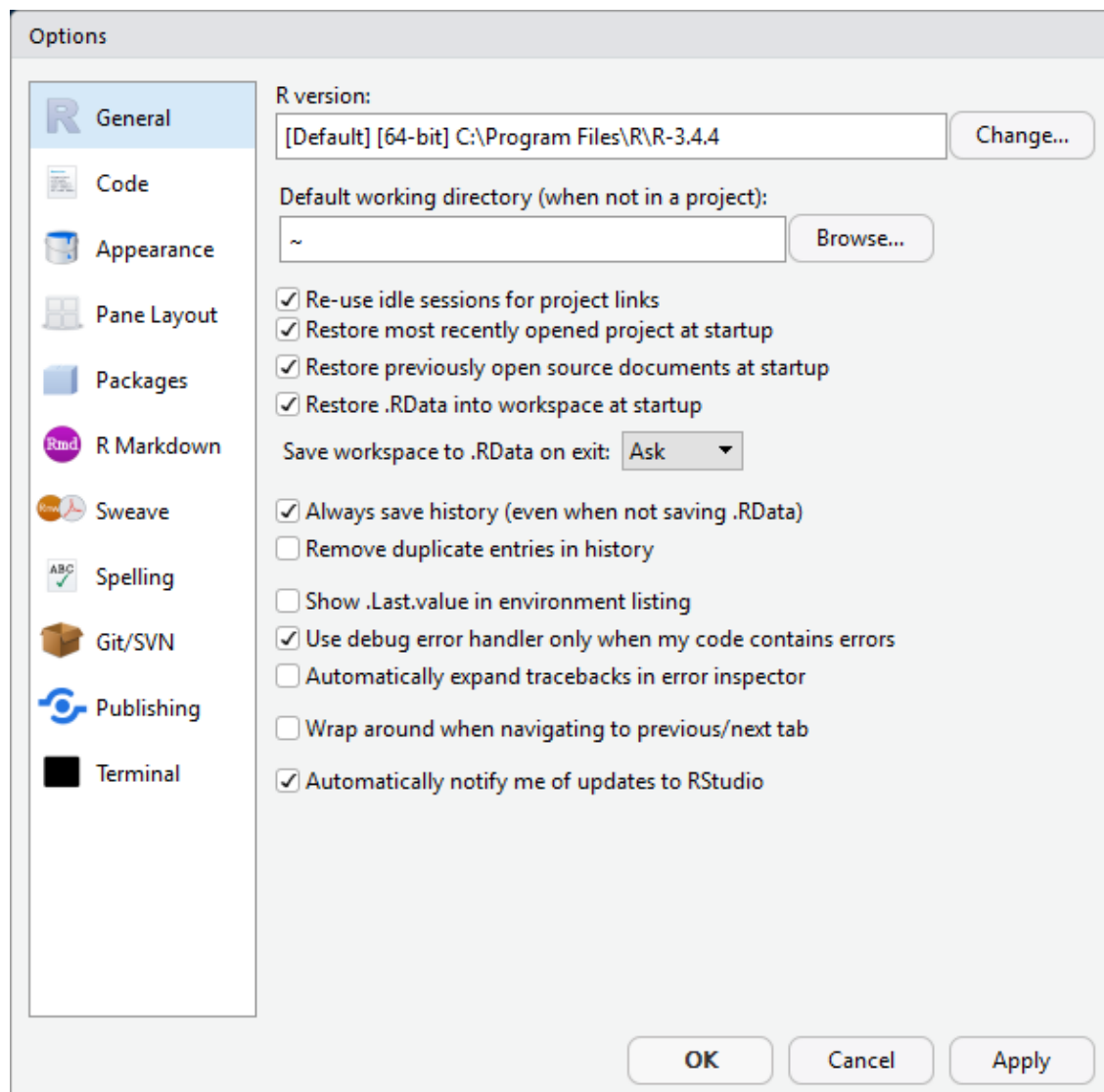


Exercise 2: To make sure that no one will burn your computer, here's how you change these settings:

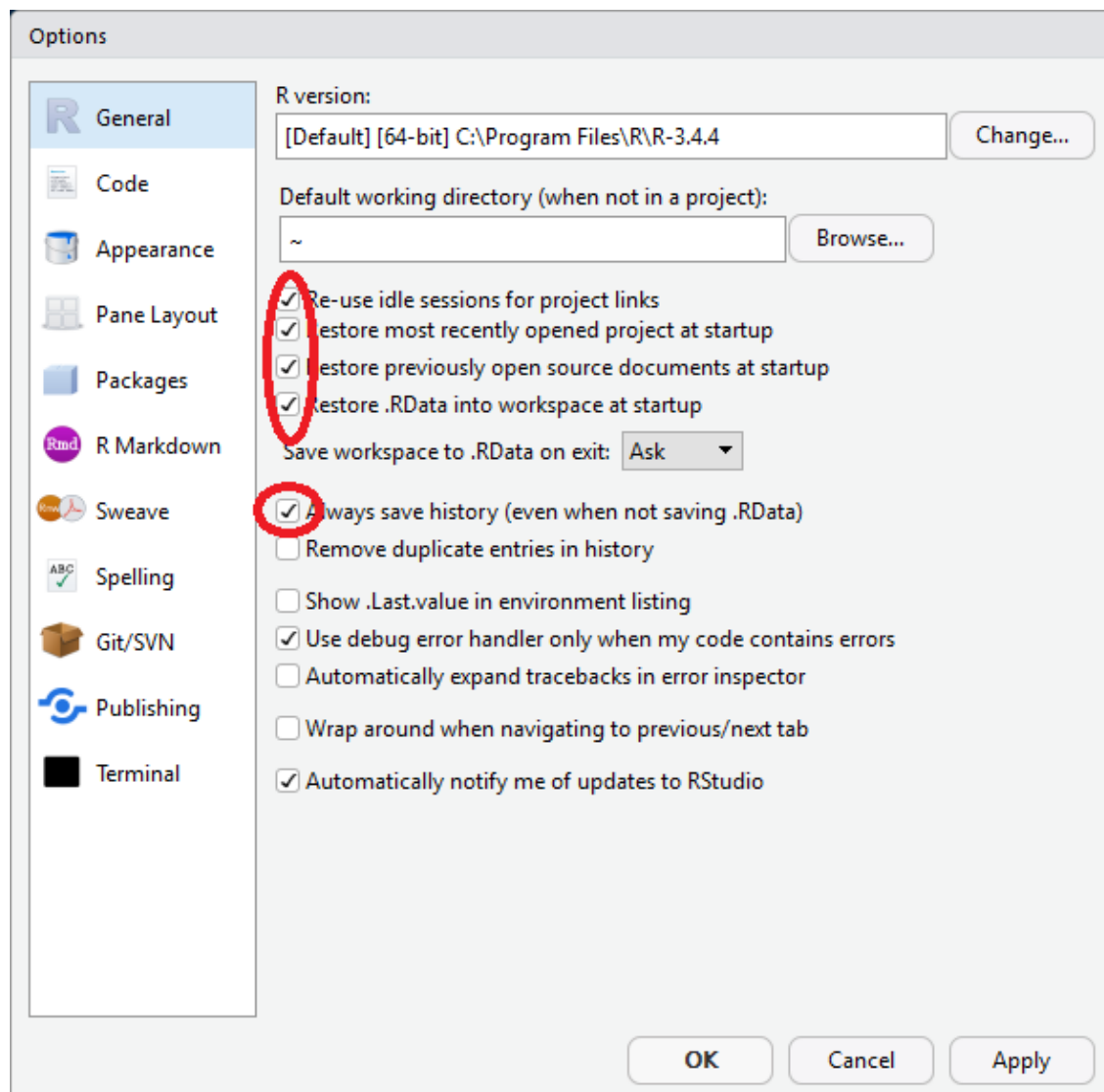
# Initial settings



# Initial settings



# Initial settings



# File paths

---





- For the purpose of this training, we will assume that you are dealing with a specific folder structure.
- Folder organization, though an important part of data work, is outside the scope of this course.
- You can find resources about it in the appendix, and we have shared with you a folder that is organized as we want it to be.
- To follow today's session, go to the `DataWork/Code` folder and open the file called `Lab 2 - Intro II.R`.
- We will use this script as a basis for the exercises, and you should modify it during this session to complete them.



- In the last session, we used the menu bar to load a data set into R
- Today, we will do that using code and referring a file's path
- File paths in R, as in Stata, are basically just strings
- Note, however, that in R we can only use forward slashes ( / ) to separate folder names

## Exercise 3: File path to your folder

Let's start by adding the folder path to the training's folder in your computer to the beginning of PART 2



- You can set file paths in your script using the `file.path()` function
- This function concatenates strings using `/` as a separator to create file paths
- This is similar to using globals to define folder paths in Stata

# File paths



Let's test if that worked:

```
# Project folder
projectFolder ←
  "/home/luis_eduardo/cs_projects/dime-r-training"

# Data work folder
dataWorkFolder ← file.path(projectFolder, "DataWork")

# Print data work folder
dataWorkFolder
```

```
## [1] "/home/luis_eduardo/cs_projects/dime-r-training/DataWork"
```



## Loading a dataset from CSV

Use the `read.csv()` command:

```
read.csv(file, header = FALSE, stringsAsFactors = FALSE)
```

- **file:** is the path (string) to the file you want to open, including its name and format (`.csv`).
- **header:** if `TRUE`, will read the first row as variable names (default is `header = FALSE`).
- **stringsAsFactors:** logical. See next slide for more.



## Loading a dataset from CSV

- Since R 4.0.0 and beyond, `stringsAsFactors = FALSE` is the default. In every previous version, the default is `TRUE`.
- This means that if your R version is 3.X.X, R will turn any string variables into factors by default when reading a `csv` file.
- This format **saves memory**, but can be tricky if you actually want to use the variables as strings.
- You can specify the option `stringsAsFactors = FALSE` to make sure you prevent R from turning strings into factors.



## Exercise 4: Test file paths

1. Save your code.
2. Start a new R session: go to `Session > New Session`. This session should be completely blank.
3. Open the code you just saved.
4. Add a line opening the data set in your code

```
# Load data set
whr <- read.csv(file.path(finalData, "whr_panel.csv"),
                 header = T)
```

1. Run the whole script. If it worked, your environment should now include the `whr` dataset and the path locals.

# Exploring a dataset

---



# Exploring a dataset



Some useful functions:

- **View()**: open the data set
- **class()**: reports object type or type of data stored
- **dim()**: reports the size of each one of an object's dimension
- **names()**: returns the variable names of a data set
- **str()**: general information on an R object
- **summary()**: summary information about the variables in a data frame
- **head()**: shows the first few observations in the dataset
- **tail()**: shows the last few observations in the dataset



## Exercise 5: Explore a dataset

Use some of the functions listed above to explore the `whr` data set.

# Exploring a dataset



```
# View the data set (same as clicking on it in the Environment pane)  
View(whr)
```

# Exploring a dataset



```
class(whr)
```

```
## [1] "data.frame"
```

```
dim(whr)
```

```
## [1] 470 8
```

# Exploring a dataset



```
str(whr)
```

```
## 'data.frame':    470 obs. of  8 variables:
## $ country      : chr  "Switzerland" "Iceland" "Denmark" "Norway" ...
## $ region       : chr  "Western Europe" "Western Europe" "Western Europe" "Western Europe" ...
## $ year         : int   2015 2015 2015 2015 2015 2015 2015 2015 2015 2015 2015 ...
## $ happiness_rank : int   1 2 3 4 5 6 7 8 9 10 ...
## $ happiness_score : num  7.59 7.56 7.53 7.52 7.43 ...
## $ economy_gdp_per_capita: num  1.4 1.3 1.33 1.46 1.33 ...
## $ health_life_expectancy: num  0.941 0.948 0.875 0.885 0.906 ...
## $ freedom       : num  0.666 0.629 0.649 0.67 0.633 ...
```

# Exploring a dataset



```
summary(whr)
```

```
##      country          region          year  happiness_rank
## Length:470      Length:470      Min.    :2015      Min.    :  1.00
## Class :character Class :character 1st Qu.:2015      1st Qu.: 40.00
## Mode  :character Mode  :character Median :2016      Median : 79.00
##                                     Mean  :2016      Mean   : 78.83
##                                     3rd Qu.:2017      3rd Qu.:118.00
##                                     Max.   :2017      Max.   :158.00
## happiness_score economy_gdp_per_capita health_life_expectancy  freedom
## Min.    :2.693      Min.    :0.0000      Min.    :0.0000      Min.    :0.0000
## 1st Qu.:4.509      1st Qu.:0.6053      1st Qu.:0.4023      1st Qu.:0.2976
## Median :5.282      Median :0.9954      Median :0.6301      Median :0.4183
## Mean    :5.371      Mean    :0.9278      Mean    :0.5800      Mean    :0.4028
## 3rd Qu.:6.234      3rd Qu.:1.2524      3rd Qu.:0.7683      3rd Qu.:0.5169
## Max.    :7.587      Max.    :1.8708      Max.    :1.0252      Max.    :0.6697
```

# Exploring a dataset



```
head(whr)
```

```
##      country      region year happiness_rank happiness_score
## 1 Switzerland Western Europe 2015           1           7.587
## 2   Iceland Western Europe 2015           2           7.561
## 3   Denmark Western Europe 2015           3           7.527
## 4    Norway Western Europe 2015           4           7.522
## 5    Canada North America 2015           5           7.427
## 6    Finland Western Europe 2015           6           7.406
## economy_gdp_per_capita health_life_expectancy freedom
## 1           1.39651           0.94143 0.66557
## 2           1.30232           0.94784 0.62877
## 3           1.32548           0.87464 0.64938
## 4           1.45900           0.88521 0.66973
## 5           1.32629           0.90563 0.63297
## 6           1.29025           0.88911 0.64169
```

# Exploring a data set



Didn't get all of those?

Don't worry, you'll see them again soon.



# Commenting

---

# Commenting



- To comment a line, write `#` as its first character
- You can also add `#` half way through a line to comment whatever comes after it
- In Stata, you can use `/*` and `*/` to comment part of a line's code. That is not possible in R: whatever comes after `#` will be a comment
- To comment a selection of lines, press `Ctrl + Shift + C`



## Exercise 6: Commenting

1. Go the `Lab 2 - Intro II.R` script. Select the lines under `PART 2: Set folder paths`.
2. Use the keyboard shortcut to comment these lines.
3. Use the keyboard shortcut to comment these lines again. What happened?

# Creating a document outline in RStudio

---

# Creating a document outline in RStudio



- RStudio also allows you to **create an interactive index** for your scripts
- To add a section to your code, create a commented line with the title of your section and add at least 4 trailing dashes, pound signs or equal signs after it



## Exercise 7: Headers

1. Open the script index and make `PART 1` a section header. Do the same for parts 2 and 3.
2. Note that once you create a section header, an arrow appears right next to it. Click on the arrows of parts 2 and 3 to see what happens.

# Creating a document outline in RStudio



- The outline can be accessed by clicking on the button on the top right corner of the script window. You can use it to jump from one section to another
- You can also use the keyboard shortcuts `Alt + L` (`Cmd + Option + L` on Mac) and `Alt + Shift + L` to collapse and expand sections

# Using packages

---



# Using packages



- Since there is a lot of people developing for R, it can have many different functionalities.
- To make it simpler, these functionalities are bundled into packages.
- A package is just **a unit of shareable code**.

# Using packages



- It may contain new functions, but also more complex functionalities, such as a Graphic User Interface (GUI) or settings for parallel processing (similar to Stata MP).
- They can be shared through R's official repository - CRAN (13,000+ packages reviewed and tested).
- There are many other online sources such as GitHub, but it's important to be careful, as these probably haven't gone through a review process as rigorous as those in CRAN.

# Using packages



- To install and use packages you can either do it with the user interface or by the command prompt.

```
# Installing a package  
install.packages("stargazer",  
                 dependencies = T)  
# the dependencies argument also installs all other packages  
# that it may depend upon to run
```

- You only have to install a package once, but you have to **load it every new session**.



## Exercise 8

1. Now load the package we just installed. Use the `library()` function to do it.

# Using packages



```
library(stargazer)
```

Notice that we used quotes around the name to install the package, but we don't need them anymore to load it.



## Warnings vs errors

What if this happens?

```
> library(stargazer)
package stargazer was built under R version 3.4.4
Please cite as:

Hlavac, Marek (2018). stargazer: well-Formatted Regression and Summary Statistics Tables.
R package version 5.2.2. https://CRAN.R-project.org/package=stargazer
```



## Warnings vs errors

R has two types of error messages, `warnings` and actual `errors`:

- `Errors` - break your code, usually preventing it from running.
- `Warnings` - usually mean that nothing went wrong yet, but you should be careful.

RStudio's default is to print warning messages, but not stop the code at the lines where they occur. You can configure R to stop at warnings if you want.

# Functions inception

---



# Functions inception



- In R, you can **write one function inside another**
- In fact, you have already done this a few times in this course
- Here's an example:

# Functions inception



```
# Doing it the long way -----  
# Create a vector with the log of the happiness score  
log_score ← log(whr$happiness_score)  
  
# Get descriptive statistics for the log vector  
summary(log_score)
```

```
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.  
## 0.9907  1.5061  1.6644  1.6576  1.8300  2.0264
```

```
# Shortcut to get to the same place -----  
summary(log(whr$happiness_score))
```

```
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.  
## 0.9907  1.5061  1.6644  1.6576  1.8300  2.0264
```

# Functions inception



- This is a simple example of **metaprogramming** (that's the real name of this technique) and may seem trivial, but it's not
- For starters, you can't do it in Stata!

# Functions inception



```

      _____ (R)
     /  /  /  /  /
    /  /  /  /  /
   /  /  /  /  /
  /  /  /  /  /
 /  /  /  /  /
/  /  /  /  /
Statistics/Data Analysis

MP - Parallel Edition

15.1 Copyright 1985-2017 StataCorp LLC
      StataCorp
      4905 Lakeway Drive
      College Station, Texas 77845 USA
      800-STATA-PC      http://www.stata.com
      979-696-4600      stata@stata.com
      979-696-4601 (fax)

681-user 4-core Stata network perpetual license:
      Serial number: 501506002486
      Licensed to: WBG User
                  World Bank Group

Notes:
1. Unicode is supported; see help unicode\_advice.
2. More than 2 billion observations are allowed; see help obs\_advice.
3. Maximum number of variables is set to 120000; see help set\_maxvar.
4. New update available; type -update all-

running C:\Program Files (x86)\Stata15\sysprofile.do ...

. sysuse auto
(1978 Automobile Data)

. summarize log(make)
variable log not found
r(111);

.
```

# Functions inception



- This is a **very powerful technique**, as you will soon see
- It's **also a common source of error**, as you can only use one function inside the other if the output of the inner function is the same as the input of the outer function
- It can also get quite tricky to follow what a line of code with multiple functions inceptions is doing
- Which is why we sometimes use pipes: `%>%`

# Functions inception



```
# 1: Doing it the long way -----  
# Create a vector with the log of the happiness score  
log_score ← log(whr$happy_score)  
  
# Get descriptive statistics for the log vector  
mean(log_score)  
  
# 2: Shortcut to get to the same place -----  
mean(log(whr$happy_score))  
  
# 3: Now with pipes -----  
whr$happy_score %>%  
  log() %>%  
  mean()
```

In a few words, `x %>% f()` is the same as `f(x)`

# Functions inception



Now that you know piping exists in R, you should know that it can **drastically improve code readability**. And from now on you can also laugh if you see this in some tidyverse nerd laptop sticker or t-shirt:



# Looping

---



# Looping



- One thing that usually gives people away as Stata users writing R code are loops
- In Stata, we use for loops quite a lot
- The equivalent to that in R would be to write a for loop like this

```
# A for loop in R  
for (number in 1:5) {  
  print(number)  
}
```

# Looping



```
# A for loop in R  
for (number in 1:5) {  
  print(number)  
}
```

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



- R, however, has a whole function family that allows users to loop through an object *in a more efficient way*\_, without using explicit loops
- They're called `apply` and there are many of them, with different use cases
- If you look for the `apply` help file, you can see all of them
- For the purpose of this training, we will only use two of them, `sapply` and `apply`



`sapply(X, FUN, ... )`: applies a function to all elements of a vector or list and returns the result in a vector. Its arguments are:

- **X**: a data frame (or matrix) the function will be applied to
- **FUN**: the function you want to apply
- **...**: possible function options

# Looping



```
# A for loop in R
for (number in c(1.2,2.5)) {
  print(round(number))
}
```

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

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

```
# A much more elegant loop in R
sapply(c(1.2,2.5), round)
```

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



A more general version is the `apply` function.

`apply(X, MARGIN, FUN, ... )`: applies a function to all columns or rows of matrix. Its arguments are

- **X**: a data frame (or matrix) the function will be applied to
- **MARGIN**: 1 to apply the function to all rows or 2 to apply the function to all columns
- **FUN**: the function you want to apply
- **...**: possible function options

# Looping



```
# Create a matrix
matrix ← matrix(c(1, 24, 9, 6, 9, 4, 2, 74, 2),
                 nrow = 3)
```

```
# Look at the matrix
matrix
```

```
##      [,1] [,2] [,3]
## [1,]    1    6    2
## [2,]   24    9   74
## [3,]    9    4    2
```

# Looping



```
# Row means  
apply(matrix, 1, mean)
```

```
## [1]  3.00000 35.66667  5.00000
```

```
# Column means  
apply(matrix, 2, mean)
```

```
## [1] 11.333333  6.333333 26.000000
```



# Custom functions

---

# Custom functions



- As we have said several times, **R is super flexible**
- One example of that is that it's **super easy and quick to create custom functions**
- Here's how:

# Custom functions



```
square ← function(x) {  
  y ← x ^ 2  
  return(y)  
}  
square(2)
```

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



## Exercise 9: Create a function

1. Create a function that calculates the z-score of a vector.

# Custom functions



```
zscore <- function(x) {  
  mean <- mean(x, na.rm = T)  
  sd   <- sd(x, na.rm = T)  
  z    <- (x - mean)/sd  
  
  return(z)  
}  
  
summary(zscore(whr$happiness_score))
```

```
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.  
## -2.3551 -0.7579 -0.0776  0.0000  0.7590  1.9492
```

# Indentation

---

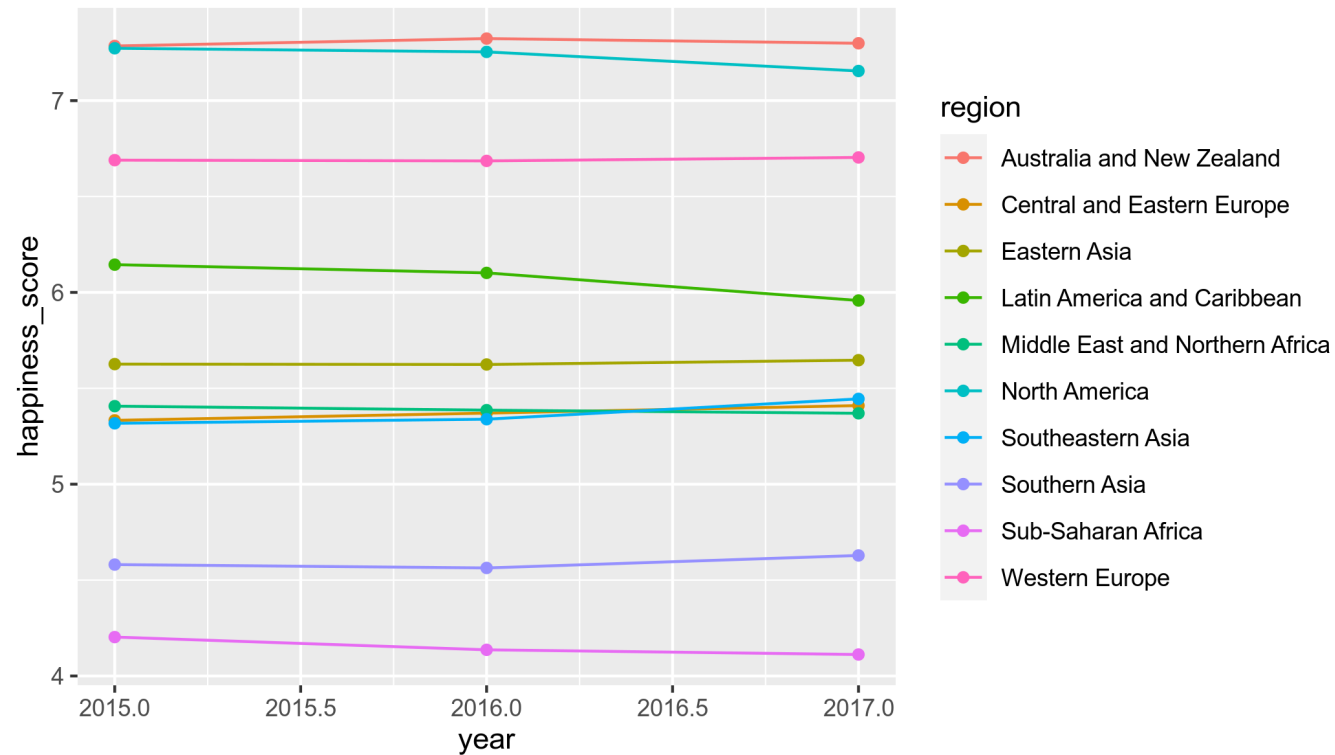


## Why indent?

```
# Here's some code
annualHappy_reg ← aggregate(happy_score ~ year + region, data = whr, FUN = mean)
plot ← ggplot(annualHappy_reg,aes(y = happy_score,x = year, color = region,
group = region)) + geom_line() + geom_point()
print(plot)
```



## Why indent?







## Why indent?

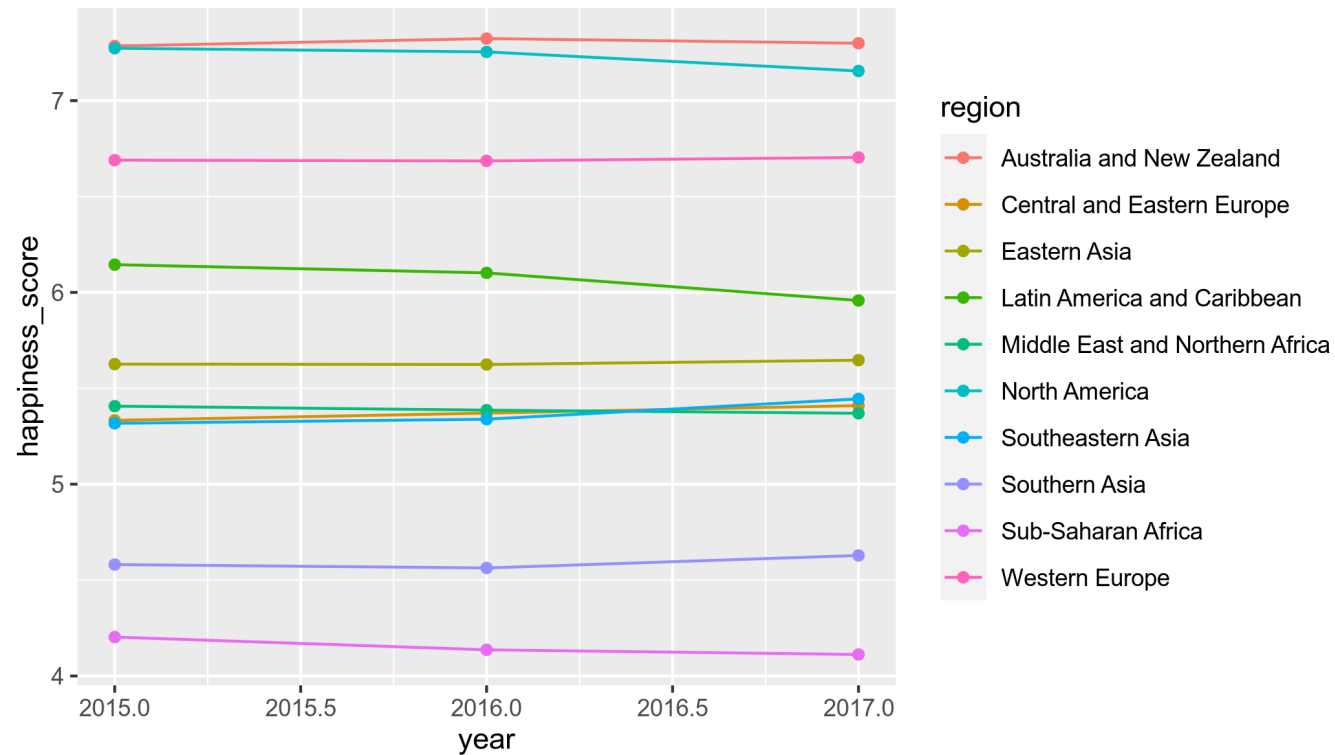
```
# Here's the same code
annualHappy_reg ←
  aggregate(happiness_score ~ year + region,
            data = whr,
            FUN = mean)

plot ←
  ggplot(annualHappy_reg,
        aes(y = happiness_score,
            x = year,
            color = region,
            group = region)) +
  geom_line() +
  geom_point()

print(plot)
```



## Why indent?





## Why indent?

- Even though R understands what unindented code says, it can be **quite difficult for a human being to read it**
- On the other hand, **white space does not have a special meaning for R**, so it will understand code that is more readable for a human being



- Indentation in R looks different than in Stata:
  - To indent a whole line, you can select that line and press `Tab`
  - To unindent a whole line, you can select that line and press `Shift + Tab`
  - However, this will not always work for different parts of a code in the same line
- In R, we typically don't introduce white space manually
- It's rather introduced by RStudio for us



## Exercise 8: Indentation in R

To see an example of how indenting works in RStudio, go back to our first example with `sapply`:

```
# A much more elegant loop in R  
sapply(c(1.2, 2.5), round)
```

1. Add a line between the two arguments of the function (the vector of numbers and the `round` function)
2. Now add a line between the numbers in the vector.

# Indentation



Note that RStudio formats the different arguments of the function differently:

```
# A much more elegant loop in R  
sapply(c(1.2,  
         2.5),  
       round)
```

Thank you!

# Appendix





## Assignment 1

Create a function that

1. Takes as argument a vector of packages names
2. Loops through the packages listed in the input vector
3. Install the packages
4. Loads the packages



## If statements

- Installing packages can be time-consuming, especially as the number of packages you're using grows, and each package only needs to be installed once
- We often use locals in Stata to create section switches to install packages
- In R, the equivalent to that would be to create a new object as a section switch

## Exercise 9: Creating an if statement

Create a dummy scalar object called PACKAGES.

- TIP: Section switches can also be Boolean objects.

# Appendix - If statements



- Now we need to create an if statement using this switch
- If statements in R look like this:

```
# Turn switch on
PACKAGES ← 1

# Install packages
if (PACKAGES == 1) {
  install.packages(packages,
                   dependencies = TRUE)
}
```



## If statements

Possible variations would include

```
# Turn switch on
PACKAGES ← TRUE

# Using a Boolean object
if (PACKAGES == TRUE) {
  install.packages(packages, dep = T)
}

# Which is the same as
if (PACKAGES) {
  install.packages(packages, dep = T)
}
```

# Appendix - Assignment 2



Create a function that

1. Takes as argument a vector of packages names
  2. Loops through the packages listed in the input vector
  3. Tests if a package is already installed
  4. Only installs packages that are not yet installed
  5. Loads the packages
- TIP: to test if a package is already installed, use the following code:

```
# Test if object x is contained in  
# the vector of installed packages  
x %in% installed.packages()
```



## File paths best practices

- We at DIME Analytics recommend always using **explicit** and **dynamic** file paths
- **Explicit** means you're explicitly stating where the file will be saved -- instead of setting the working directory, for example
- **Dynamic** means that you don't need to adjust every file path in the script when you change from one machine to another -- they're updated based on a single line of code to be changed



## File paths

- Explicit and dynamic file path:

```
# Define dynamic file path
finalData <- "C:/Users/luiza/Documents/GitHub/
             dime-r-training/
             DataWork/DataSets/Final"

# Load data set
whr <- read.csv(file.path(finalData, "whr_panel.csv"),
                header = T)
```



## Using packages

Once a package is loaded, you can use its features and functions. Here's a list of some useful and cool packages:

- `Rcmdr` - easy to use GUI
- `swirl` - an interactive learning environment for R and statistics.
- `ggplot2` - beautiful and versatile graphics (the syntax is a pain, though)
- `stargazer` - awesome latex regression and summary statistics tables
- `foreign` - reads `.dta` and other formats from inferior statistical software
- `zoo` - time series and panel data manipulation useful functions
- `data.table` - some functions to deal with huge data sets
- `sp` and `rgeos` - spatial analysis
- `multiwayvcov` and `sandwich` - clustered and robust standard errors
- `RODBC`, `RMySQL`, `RPostgreSQL`, `RSQLite` - For relational databases and using SQL in R.





## Resources

- A discussion of folder structure and data management can be found here:  
[https://dimewiki.worldbank.org/wiki/DataWork\\_Folder](https://dimewiki.worldbank.org/wiki/DataWork_Folder)
- For a broader discussion of data management, go to [https://dimewiki.worldbank.org/wiki/Data\\_Management](https://dimewiki.worldbank.org/wiki/Data_Management)



## Git

Git is a version-control system for tracking changes in code and other text files. It is a great resource to include in your work flow.

We didn't cover it here because of time constraints, but below are some useful links, and DIME Analytics provides trainings on Git and GitHub, so keep an eye out for them.

- **DIME Analytics git page:** <https://worldbank.github.io/dimeanalytics/git/>
- **A Quick Introduction to Version Control with Git and GitHub:** <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004668>



## R projects

If you have used R before, you may have heard of RStudio Projects. It's RStudio's suggested tool for workflow management. DIME Analytics has found that it is not the best fit for our needs, because

1. In DIME, we mainly use Stata, and we prefer to keep a similar structure in R (Stata 15 also has a projects feature, but it is not yet widely adopted)
2. We need to keep our code and data in separate folders, as we store code in GitHub and data in DropBox

However, if you want to learn more about it, we recommend starting here: <https://r4ds.had.co.nz/workflow-projects.html>