

# INTRODUCTION TO DATA SCIENCE

**Pierre-Alexandre Balland**

# Key issues

- A data-driven revolution
- The tech giants and the rest of the world (a story of positive feedback loops)
- Key applications of data science (business, policy and science)
- From data to stories
- A lot of data > really precise data
- Algorithmic bias, ethics and echo chambers
- Network Science, Big Data, AI, Machine Learning and Deep Learning

# Computer lab: R & RStudio

- In this course we will perform structural network analysis with packages implemented in the R statistical software
- R is the software – but we will use Rstudio as an interface
- R is an open-source project lifted by a virtual community of thousands of developers and million of users worldwide

# Why R?

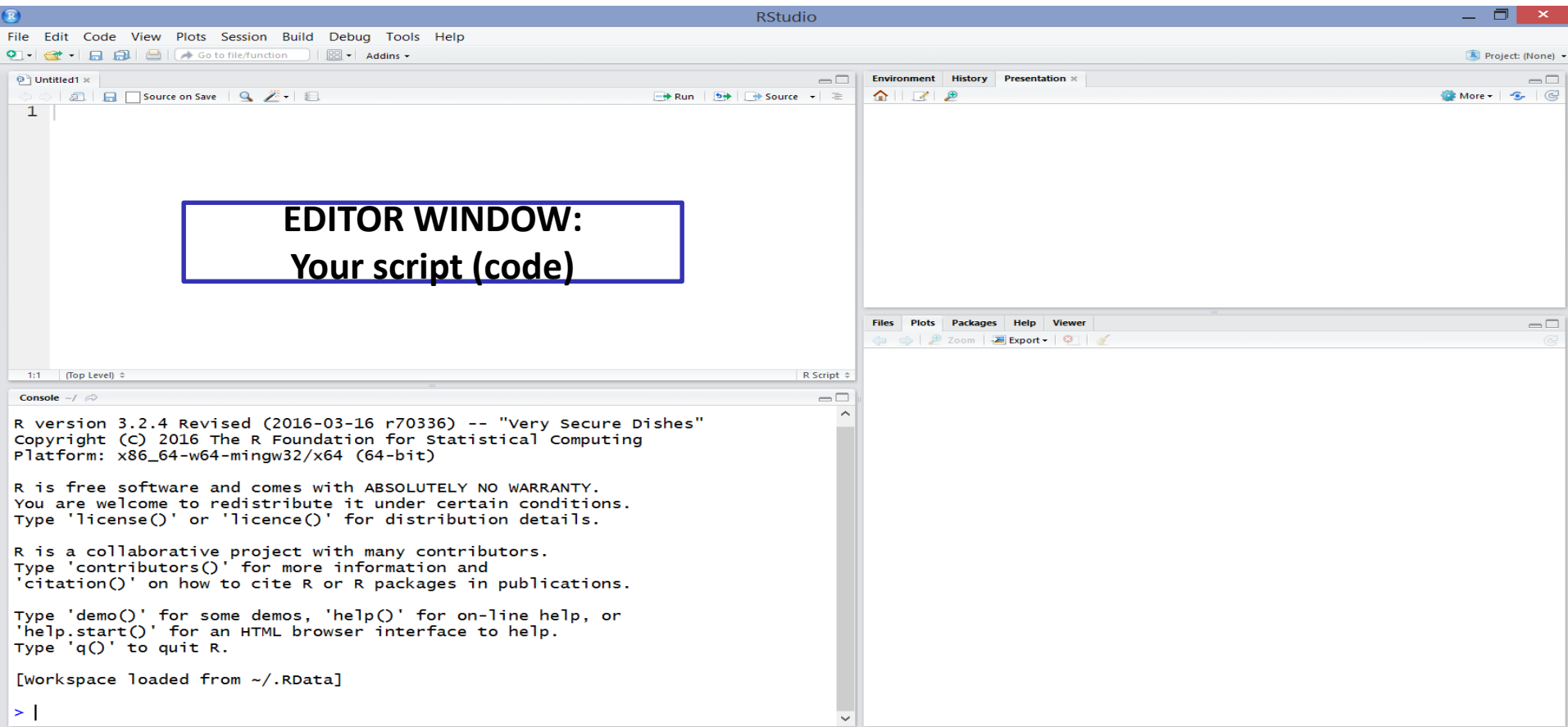
- Reproducibility – R scripts
- Today R offers the most elegant and comprehensive language for the structural and dynamic analysis of networks
- It's free and contains state-of-the-art statistical and graphical routines not yet available in other software
- You can do all your analysis in R, but also data scrapping, create a webpage, or write your research paper

# Getting started with R

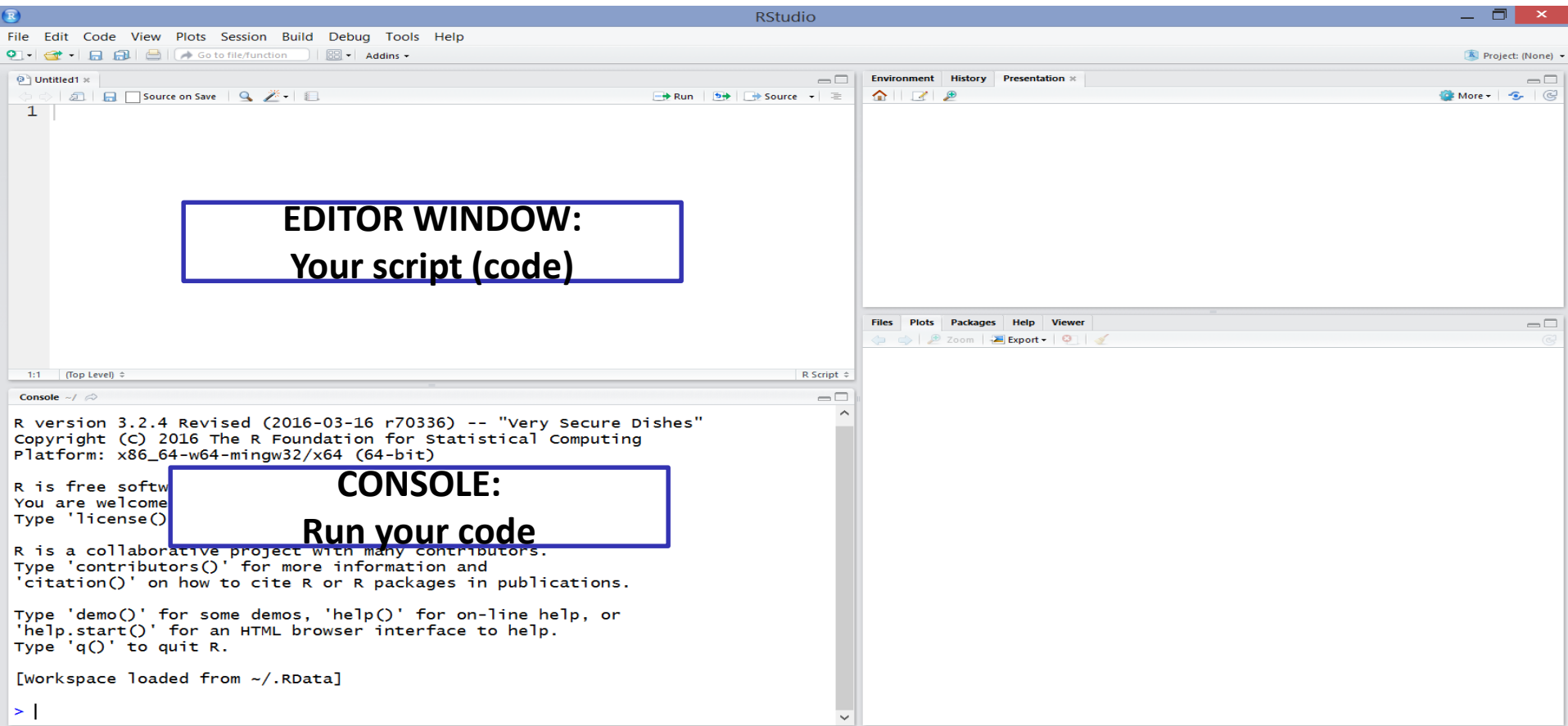
- Using R is easier than it looks like. And once you master it, you save a ridiculous amount of time
- Afraid of R? It is just a big calculator (a very smart one)
- R is case sensitive
- The `#` character at the beginning of a line signifies a comment, it is ignored by R



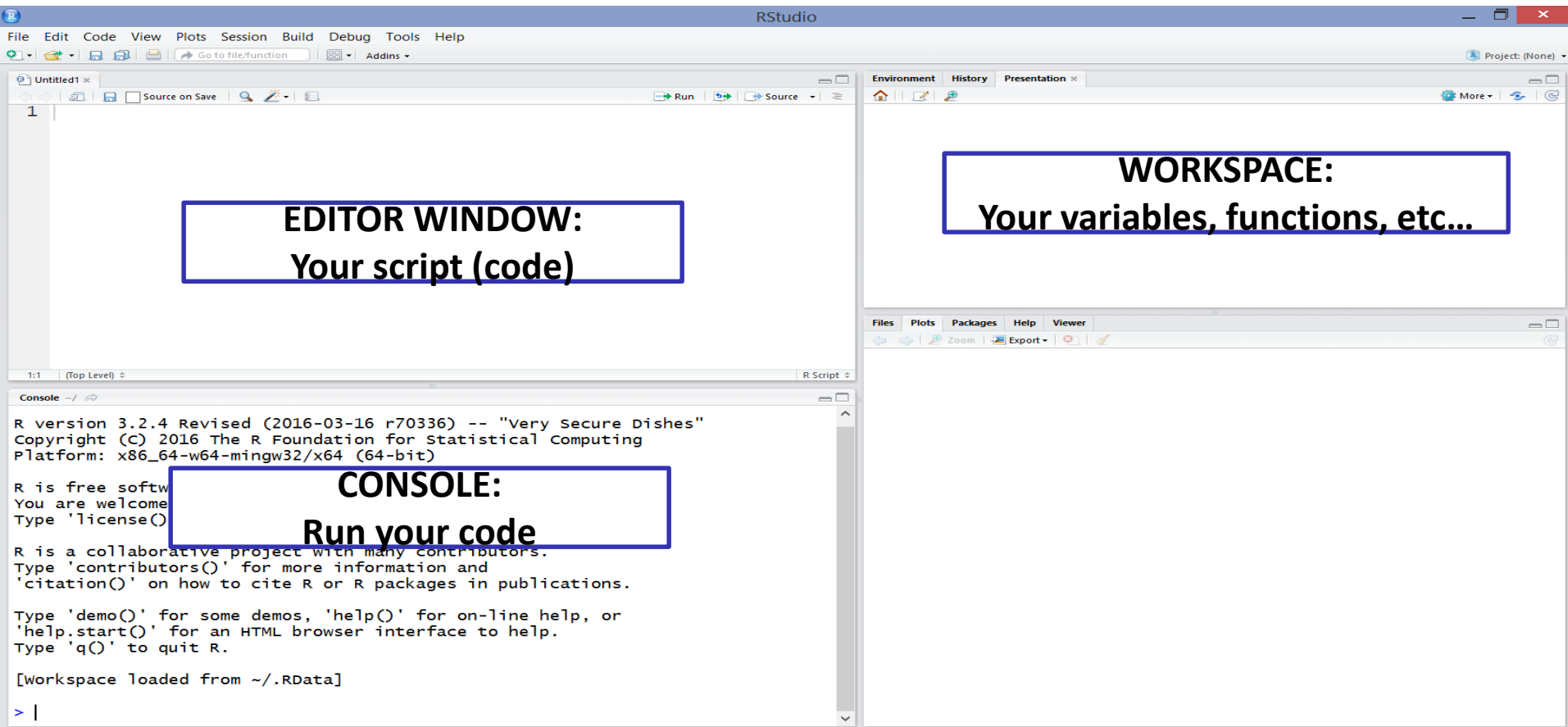
# RStudio



# RStudio

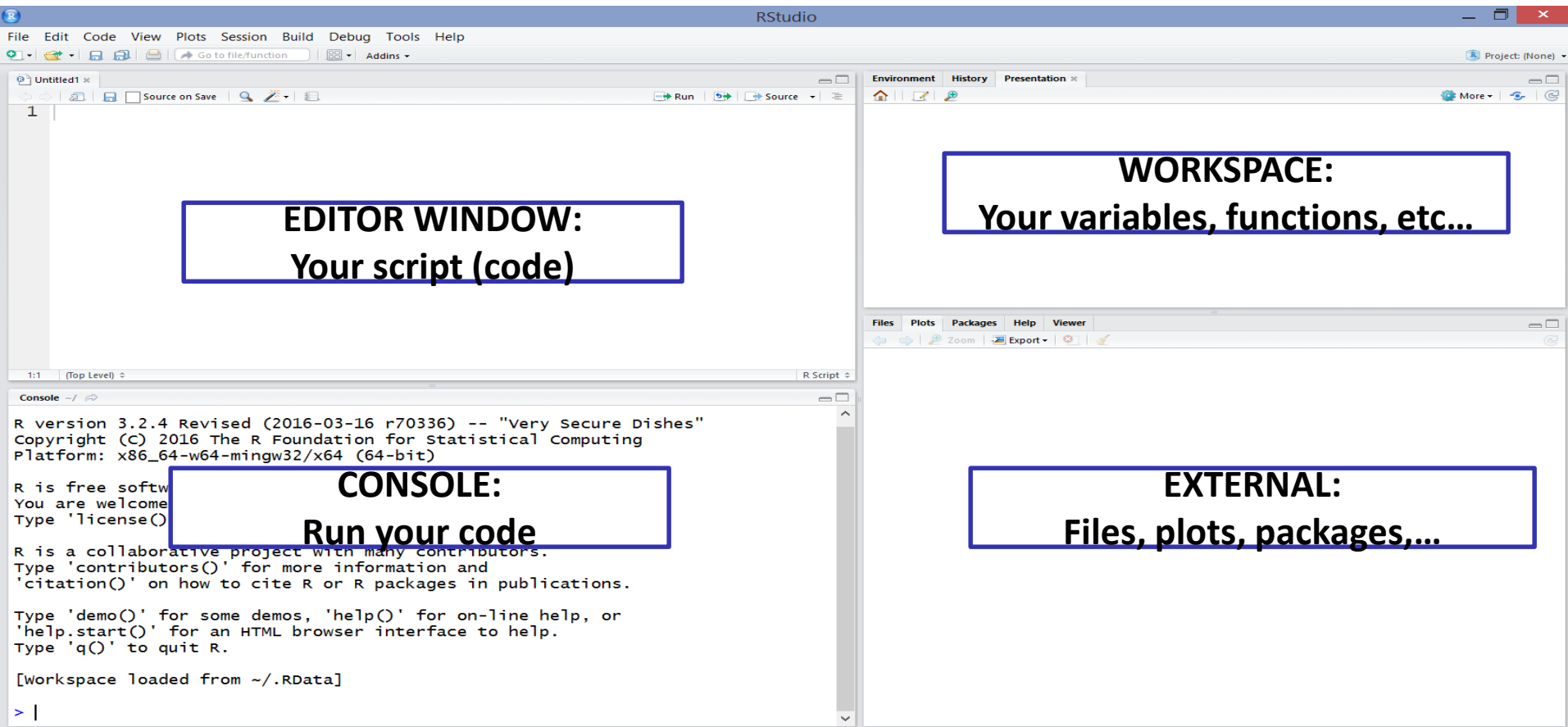


# RStudio





# RStudio



# Let's get started

The screenshot displays the RStudio interface. The top-left pane shows a script editor with the code `1 10+5`. A blue arrow points from this code to the 'Run' button in the top toolbar. The bottom-left pane shows the console output, which includes the R startup message and the result of the calculation `10+5`.

Environment History Presentation

Files Plots Packages Help Viewer

1:5 (Top Level) R Script

Console

```
Copyright (C) 2016 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

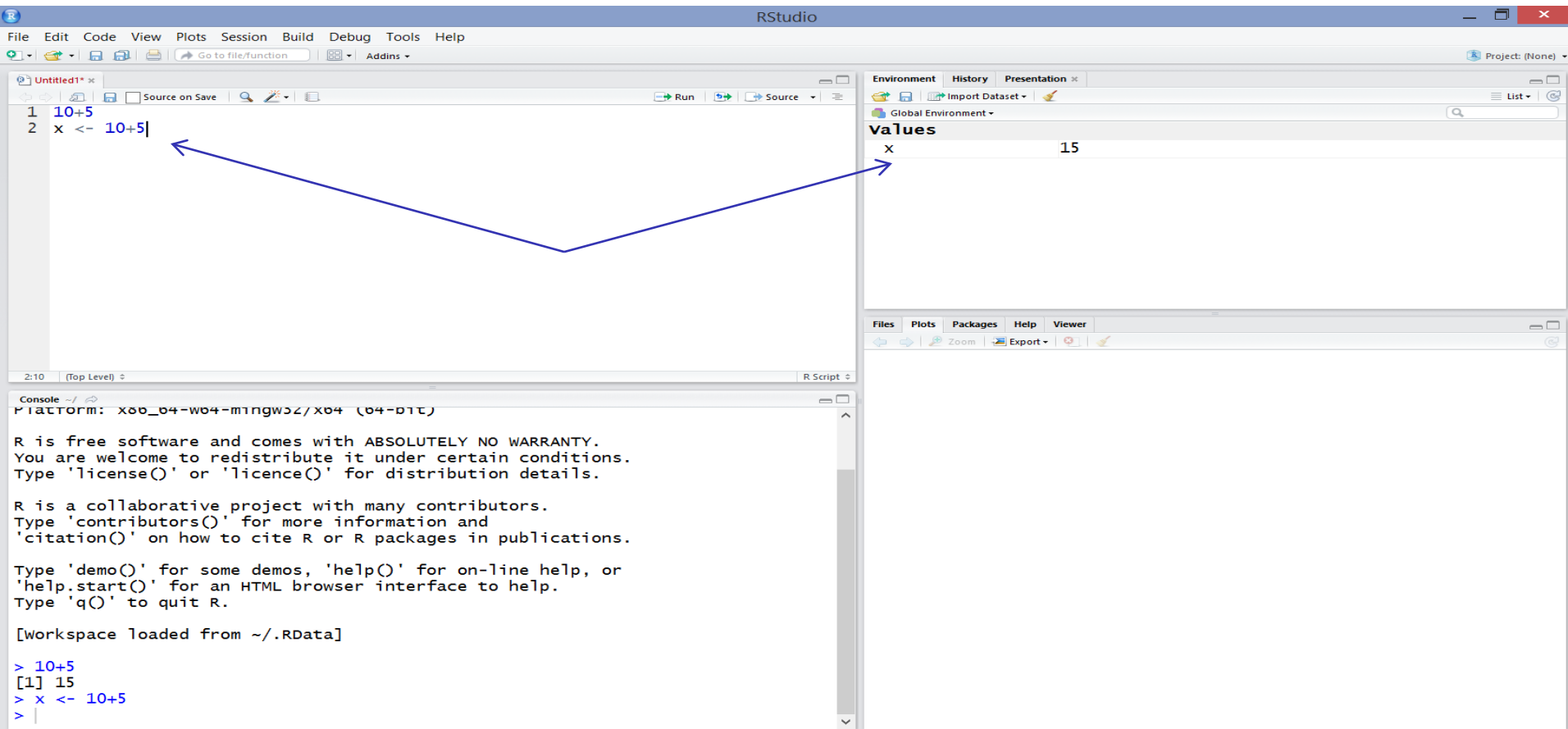
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> 10+5
[1] 15
>
```

# Create a variable "x"



The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains two lines of R code:

```
1 10+5  
2 x <- 10+5|
```

A blue arrow points from the expression `10+5` in line 2 to the `values` pane.
- Environment Pane:** Shows the **Global Environment** with a table of values:

values	
x	15
- Console:** Shows the R startup message and the execution of the code:

```
Platform: x86_64-w64-mingw32/x64 (64-bit)  
  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.  
  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
  
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.  
  
[Workspace loaded from ~/.RData]  
  
> 10+5  
[1] 15  
> x <- 10+5  
> |
```

# Create a variable y

The screenshot displays the RStudio environment with three main panels: the Source editor, the Environment pane, and the Console.

**Source Editor:** Contains the following R code:

```
1 10+5
2 x <- 10+5
3 y <- "geography"
```

**Environment Pane:** Shows the Global Environment with the following values:

values	
x	15
y	"geography"

**Console:** Shows the output of the R session, including the R startup message and the results of the code execution:

```
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> 10+5
[1] 15
> x <- 10+5
> y <- "geography"
>
```

# Check the value of x and y

The screenshot displays the RStudio environment with three main panels:

- Source Editor (Left):** Contains an R script with the following code:

```
1 10+5
2 x <- 10+5
3 y <- "geography"
4 # check the value of x and y
5 x
6 y
7 |
```
- Environment Panel (Top Right):** Shows the 'Global Environment' with the following values:

values	
x	15
y	"geography"
- Console (Bottom):** Shows the output of the script execution:

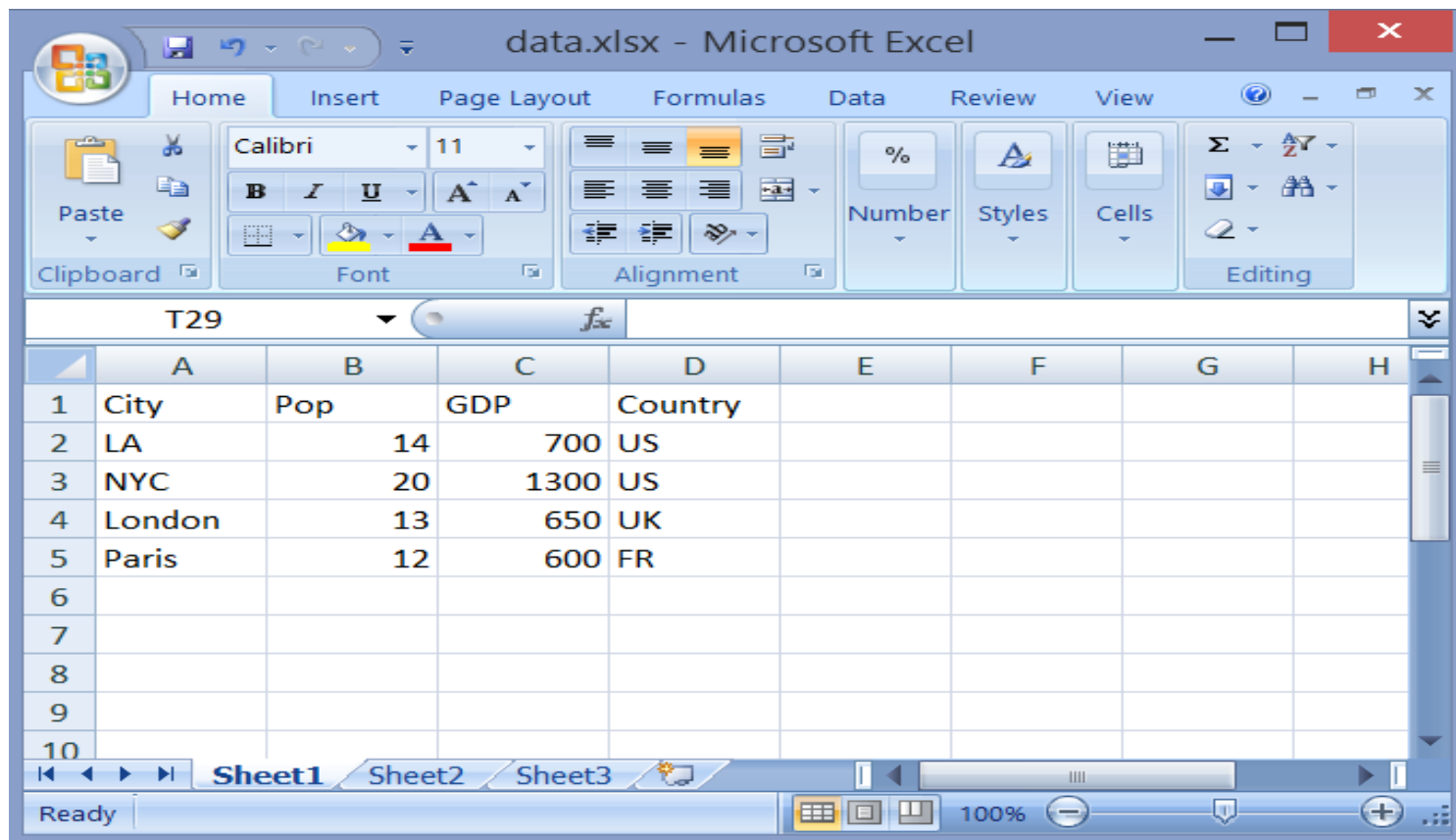
```
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

> 10+5
[1] 15
> x <- 10+5
> y <- "geography"
> # check the value of x and y
> x
[1] 15
> y
[1] "geography"
> |
```

# Let's create a toy dataset

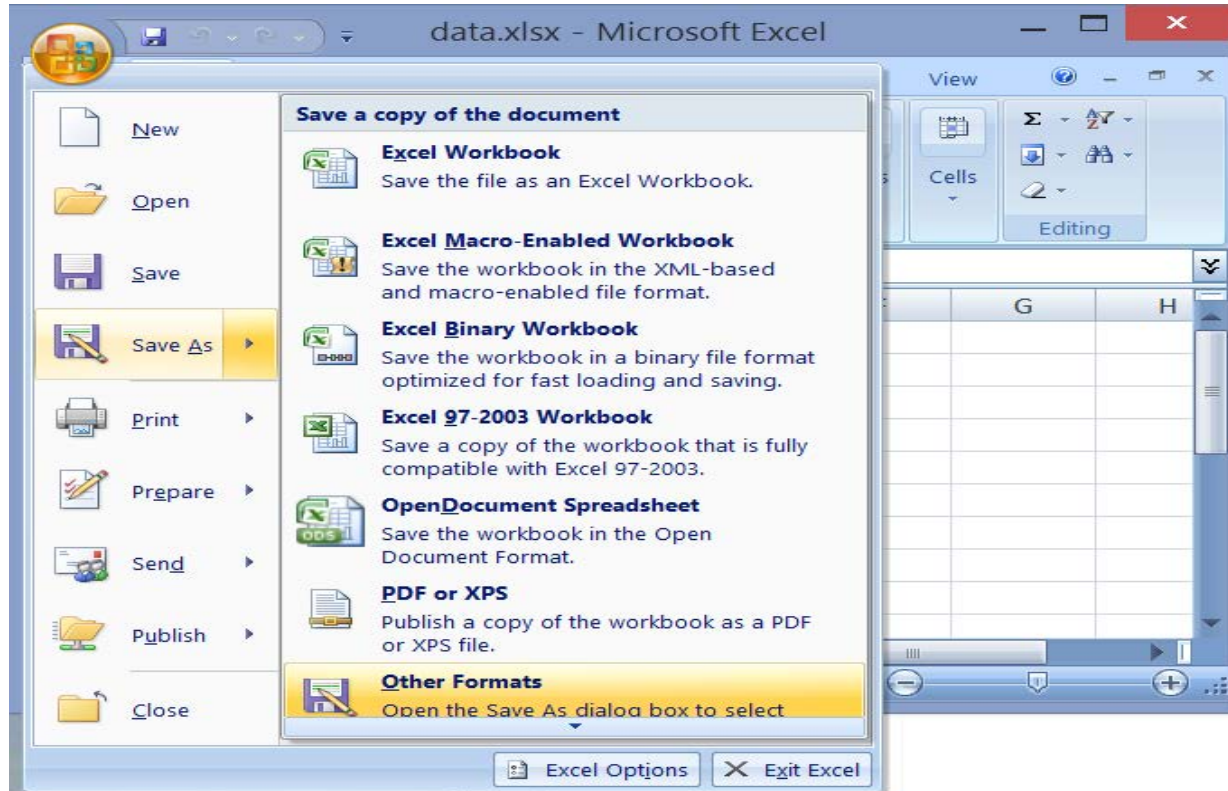


The screenshot shows the Microsoft Excel interface with the file name "data.xlsx". The ribbon is set to "Home", and the "Font" group is active. The spreadsheet contains a table with 4 columns: City, Pop, GDP, and Country. The data is as follows:

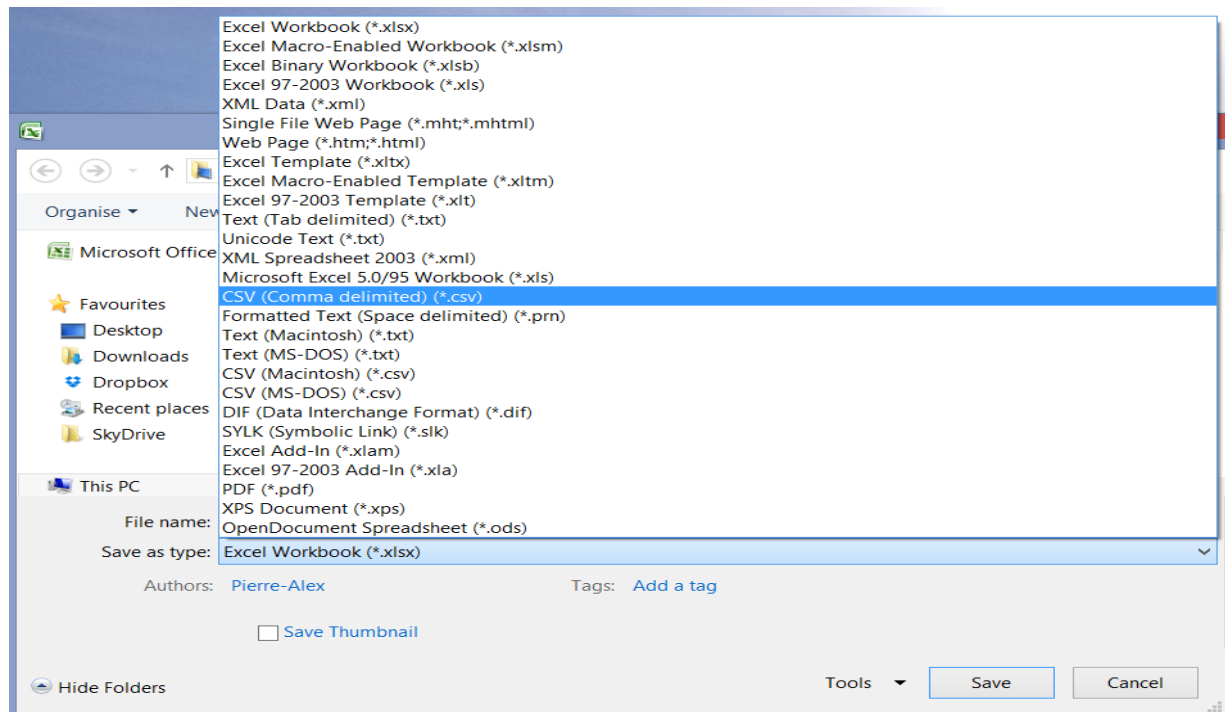
	A	B	C	D	E	F	G	H
1	City	Pop	GDP	Country				
2	LA	14	700	US				
3	NYC	20	1300	US				
4	London	13	650	UK				
5	Paris	12	600	FR				
6								
7								
8								
9								
10								

The status bar at the bottom indicates "Ready" and "100%".

# Save as a .csv file

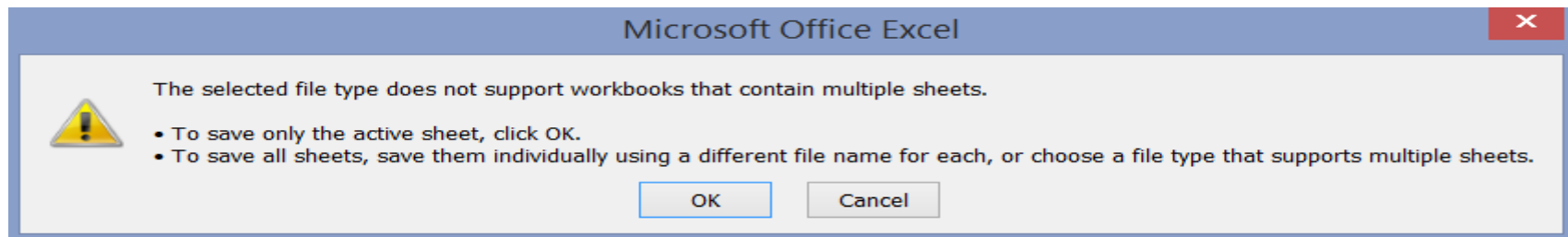


# Save as a .csv file

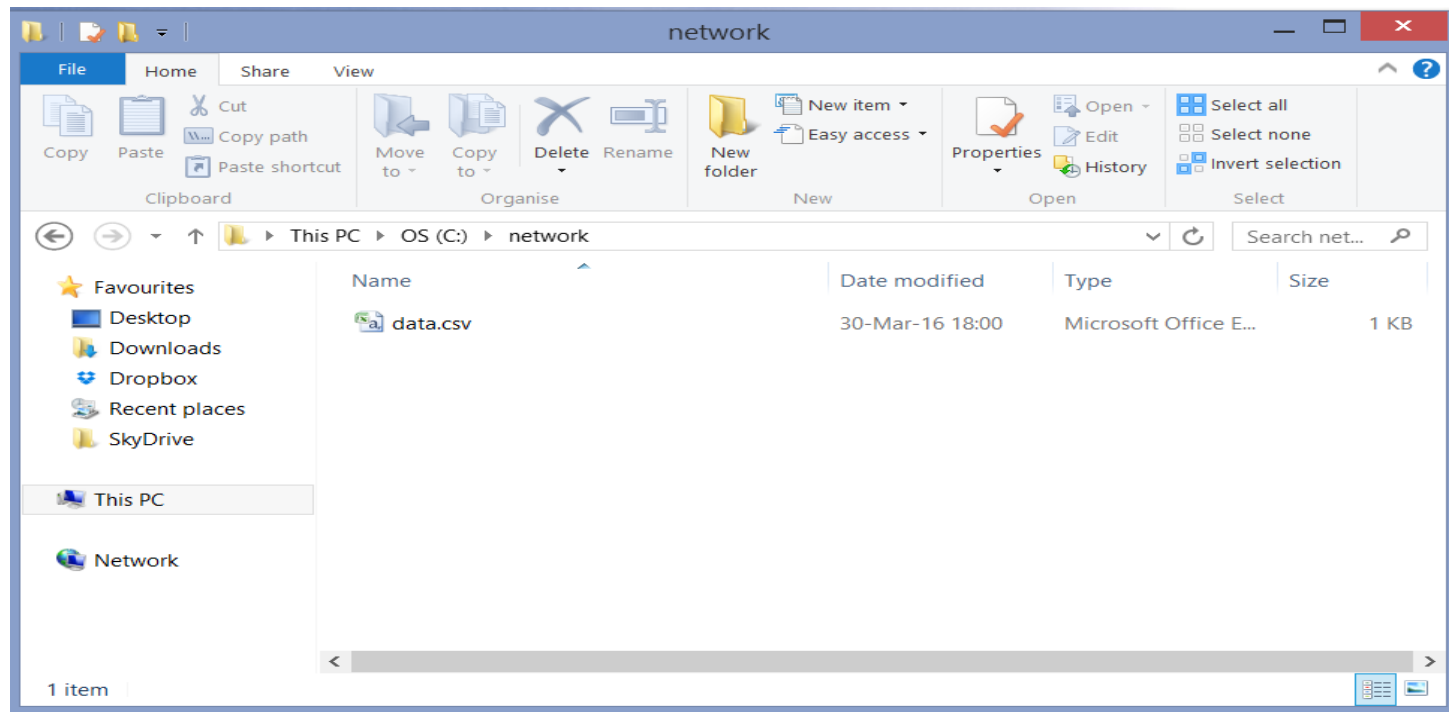




# Two warnings: ok



# Create a new folder and move the .csv



# This is your file path

