Show me the data!
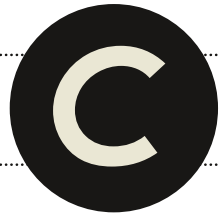
# Big Data & Social Analysis R

Instructors: Chung-pei Pien

ZU1942001/266868001/Z23937001/ZM1941001

International College of
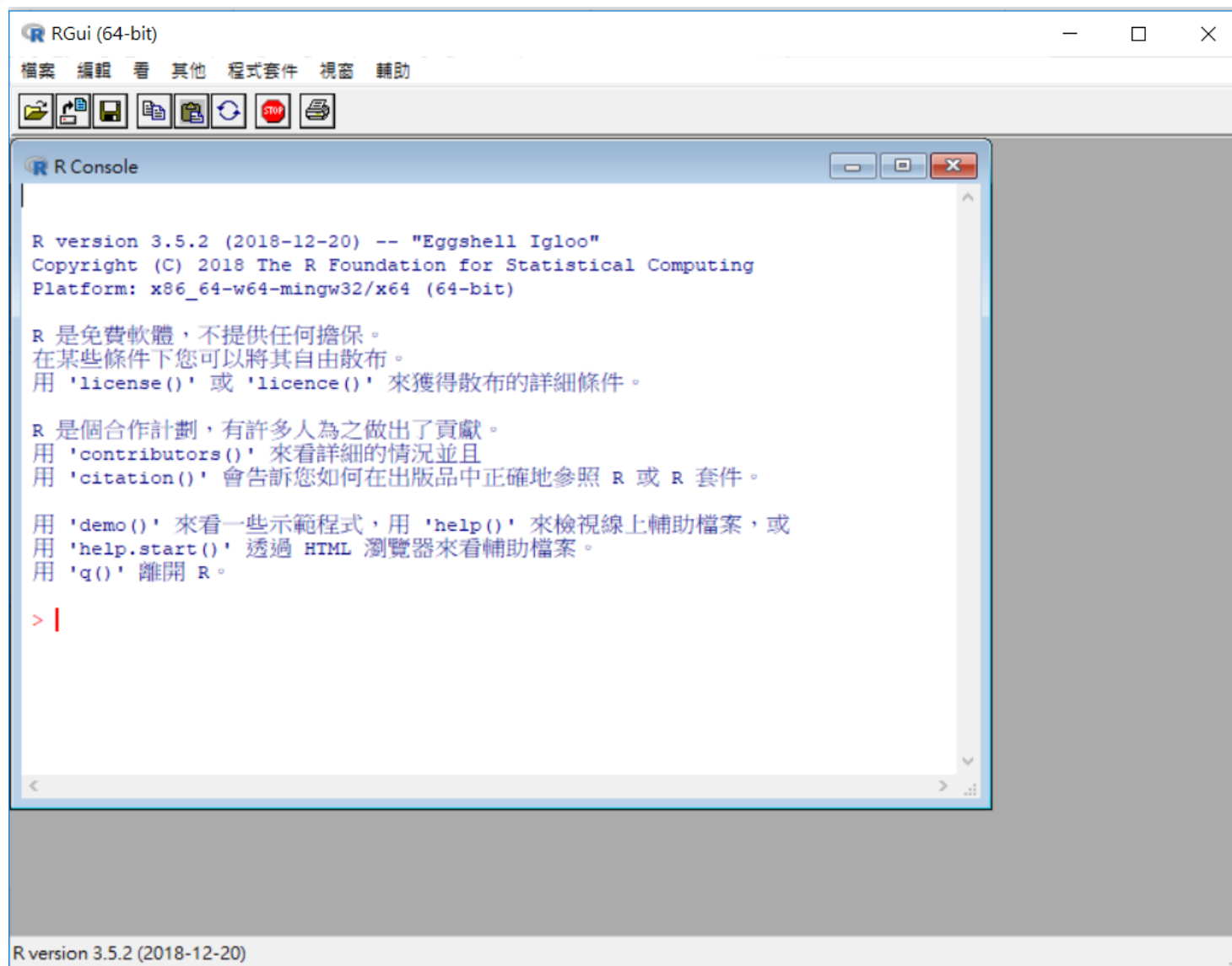INNOVATION
National Chengchi University
國立政治大學創新國際學院

# **C**ONTENTS

# R and R-Studio

# 01 R and R-Studio

# R and R-Studio

R： Base computer
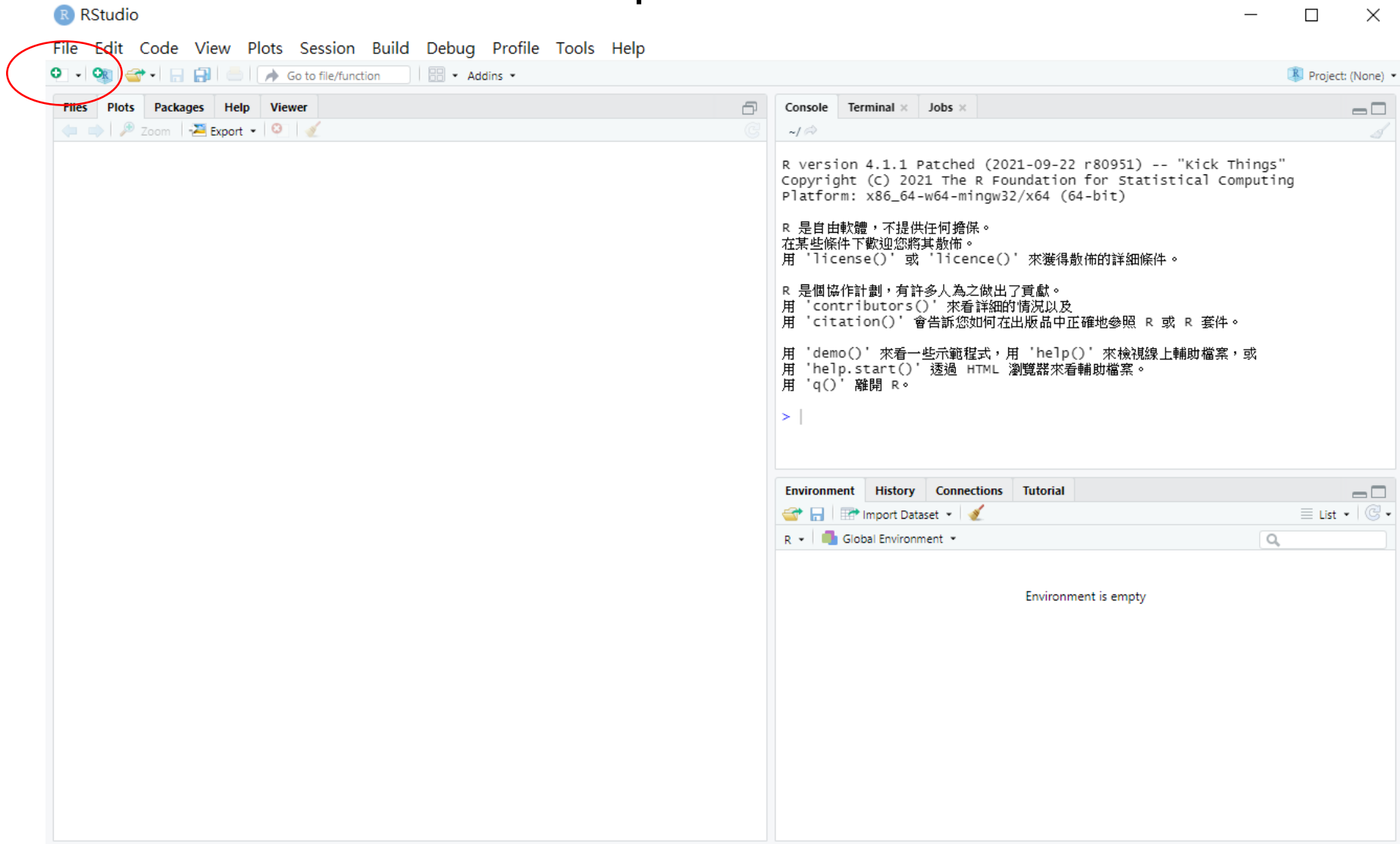R-Studio： add RAMs, two monitors

# 01 R and R-Studio

- Console：Input commands and show results

- Environment：the list of objects

- Source：R scripts and content of objects

- Plots：the plots

- Packages：the list of packages

- Help

# R and R-Studio
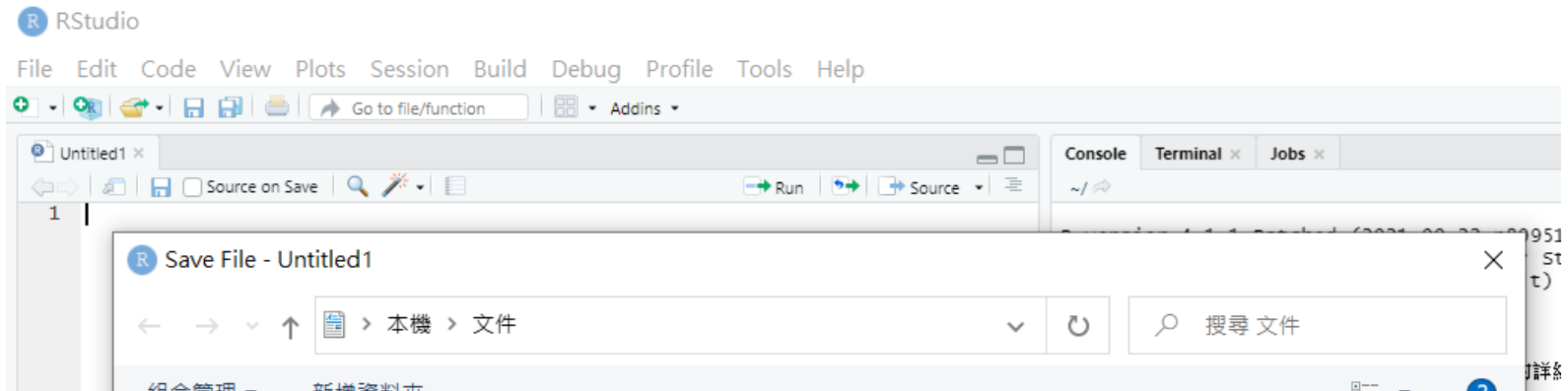
Create a new R script file

# R and R-Studio

1. Create a new project folder
2. Save the R script file into this folder

# R and R-Studio

R-Studio Cloud: A cloud-based solution that allows anyone to do, share, teach and learn data science online
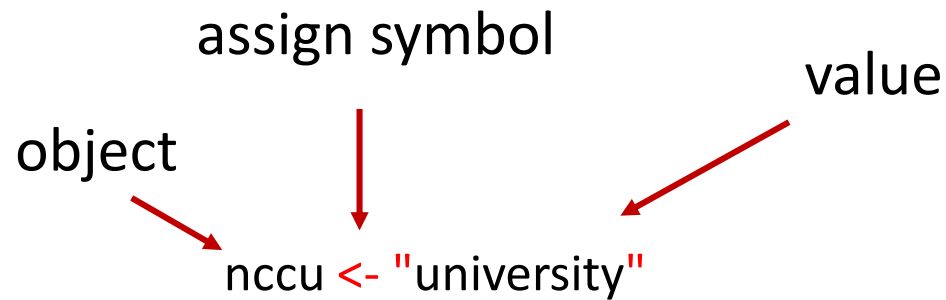
# Basic Skills of R

# The basic concept of R:

## Assign values to objects

The basic concept of coding in R is to assign a set of values to objects. Then you can transform, analyze, calculate, or represent the objects

# Assign values to an object

assign symbol

value

object

nccu <- "university"

# Assign values to an object

assign symbol

value

object

nccu <- "2022"

Objects is the key element in R：

1.  An object can involve one or million data.

2.  There are different kinds of objects. They offer different ways to store data.

The kinds of objects in R

1. Vector
2. List
3. Matrix
4. Table (dataframe)
5. Others:

The kinds of objects in R

1.  Vector
2.  List
3.  Matrix
4.  Table (dataframe)
5.  Others:

Vector: a set of numbers (numeric) or characters is created by c() function.

nccu <- c("taipei", "ici", "zoo", "beautiful")

z <- c(1, 5, 10, 15, 27)

# Basic Skills of R

If a vector includes numerics and or characters?

nccu <- c("taipei", "ici", "zoo", "beautiful", 1, 6, 12)

R recognizes it as Character

# Basic Skills of R

To create/identify an object as numerics or characters is extremely important in R!!!

lawmaker_id <- c("0123", "1123", "0894", "1305")

lawmaker_id2 <- c(0123, 1123, 0894, 1305)

What is coding?

Coding is to pick up specific <span style="color:red">elements</span> in your objects and calculate/change their value.

```
nccu <- c("taipei", "ici", "zoo", "beautiful")

z <- c(1, 5, 10, 15, 27)
```

What is coding?

Coding is to pick up specific elements in your objects and calculate/change their value.

```
nccu <- c("taipei", "ici", "zoo", "beautiful")

nccu[1]
nccu[3]
nccu[5]
```

What is coding?

Coding is to pick up specific elements in your objects and calculate/change their value.

```
nccu <- c("taipei", "ici", "zoo", "beautiful")

nccu[1]
nccu[3]
nccu[5]

nccu[1] <- c("home")

nccu

nccu[5] <- c("international")
```

A Function in R is to do a specific task, such as creating objects, picking up elememnts, making plots, creating models, and so on.



R Reference Card

# Basic Skills of R

There is a vector student to show a class's students' gender information.

student <- c("m", "f", "f", "m", "m", "m", "f", "f", "f", "f")

1. How many students are in this class?

2. How many male students are in this class?

3. How many female students are in this class?

There is a vector student to show a class's students' gender information.

student <- c("m", "f", "f", "m", "m", "m", "f", "f", "f", "f")

1. How many students are in this class?

   length(student)

2. How many male students are in this class?

   student[student == "m"]
   length(student[student == "m"])

There is a vector student to show a class's students' gender information.

student <- c("m", "f", "f", "m", "m", "m", "f", "f", "f", "f")

3. How many female students are in this class?

Answer this question in Moodle (Practice 1)

# Basic Skills of R

There is a vector age to show a company's workers' age.

age <- c(45, 60, 22, 61, 34, 59, 64, 54, 29, 31)

1. How many workers are in this company?

2. How many workers' age are larger than 60?

3. How many workers' age are smaller than 30?

There is a vector age to show a company's workers' age.

age <- c(45, 60, 22, 61, 34, 59, 64, 54, 29, 31)

## 1. How many workers are in this company?

    length(age)

## 2. How many workers' age are larger than 60?

    age[age >= 60]
    length(age[age >= 60])

There is a vector age to show a company's workers' age.

age <- c(45, 60, 22, 61, 34, 59, 64, 54, 29, 31)

3. How many workers' age are smaller than 30?

Answer this question in Moodle (Practice 2)

# Basic Skills of R

There is a vector student to show a class's students' gender information.

student <- c("m", "f", "f", "m", "m", "m", "f", "f", "f", "f")

4. The ratio of male students in this class

length(student[student == "m"]) / length(student)

# Basic Skills of R

There is a vector id to show a NCCU class's undergraduate students' university id.

id <- c("110111222", "109222111", "109222333", "108333444", "110555666")

1. How many freshman students are in this class?

There is a vector id to show a NCCU class's undergraduate students' university id.

id <- c("110111222", "109222111", "109222333", "108333444", "110555666")

1.   How many freshman students are in this class?

```
nchar(id)
year <- substr(id, 1, 3)
length(year[year == "110"])
```

There is a vector uid to show a NCCU class's students' university id.

uid <- c("110101222", "109252111", "109202333", "108351444", "110151666")

1. How many undergraduate students are in this class?
2. How many undergraduate and freshman students are in this class?

   Answer these questions in Moodle (Practice 3)

# Basic Skills of R

There is a vector birth to show a NCCU class's students' birth year.

birth <- c("2002", "1999", "2001", "1998", "2000")

1.  What is the average of the students' age

# Basic Skills of R

There is a vector birth to show a NCCU class's students' birth year.

birth <- c("2002", "1999", "2001", "1998", "2000")

1. What is the average of the students' age

birth <- as.numeric(birth)
age_b <- 2022 - birth
mean(age_b)

median(age_b)
max(age_b)

There is a vector birthdate to show a baseball team players' birth date.

birthdate <- c("1983-11-01", "1995-01-19", "2001-06-23", "1987-12-09", "1999-10-21", "1999-03-31")

1. What is the average of the players' age

   Answer this question in Moodle (Practice 4)

# 115<sup>th</sup> US Congress Data

# 115<sup>th</sup> US Congress Data

The term of office of 115<sup>th</sup> US Congress was from January 3, 2017, to January 3, 2019.

It was reelected in November 2018, two years after Trump's win in 2016.

The best form to use this data is a table. We will do it in the next week. Today, we still use vectors to handle this data.

# 115<sup>th</sup> US Congress Data

Gender, year, party are three vectors to represent the 115th US House lawmakers' characteristics.

Please answer the following question in Moodle (Practice 5):

1. The female ratio of the lawmakers
2. The mean of lawmakers' age
3. Which party dominate the House?

Before you conduct this project, please think about why lawmakers' gender ratio, mean of age, and party domination are important?

# Assignment

This week's homework will answer the questions about the changes in gender inequality, age, party's control from 115$^{th}$ to 117$^{th}$ US congress.

# Assignment

```
1   #Class: Week 02
2   #Course: Big Data and Social Analysis
3   #Semester: Spring 2021
4   #Lesson: R, Vector, and Object
5   #Instructor: Chung-pei Pien
6   #Organization: ICI, NCCU
7
8 ▾ ### Student Information --------
9
10  #Chinese Name:
11  #English First Name:
12  #UID:
13  #E-mail:
14
15 ▾ ### Questions --------
16
17  #Run the following vector codes. You will have 9 vector objects.
18
19  #115th House Gender
20
21  gender_115 <- c("M", "F", "M", "M", "M", "M", "M", "M", "M", "M", "M", "M", "M",
22
23  #115th House Birth Date
24
25  birth_115 <- c("1954-09-16", "1946-05-27", "1965-07-22", "1979-06-19", "1951-11-(
26              "1946-05-12", "1952-12-23", "1947-12-21", "1955-08-11", "1963-12-22",
27
28  #115th House Party
29
```

# Assignment

115th -117th US Congress Data

```
59
60  #Please answer the following questions. Remember, No Comments, No Points!!!!!!!!
61
62  #Question 1: (3 points)
63  #From 115-117 terms of US house, which term has the largest number of lawmakers?
64
65  #Question 2: (9 points)
66  #From 115-117 terms of US house, which term has worse gender inequality performance?
67
68  #Question 3: (9 points)
69  #From 115-117 terms of US house, which term's age is oldest?
70
```

43

# Assignment

115th -117th US Congress Data