# R語言簡介

王貿

國立臺灣大學行為與資料科學研究中心助理研究員
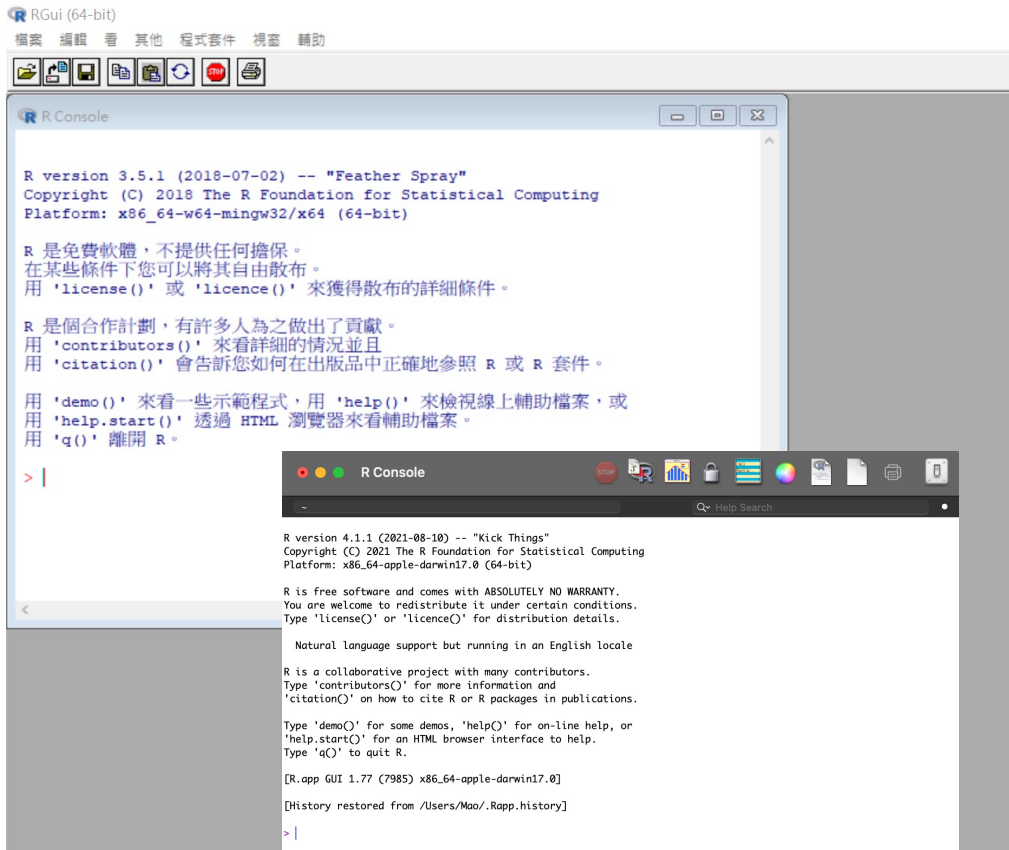
國立臺灣大學政治學系博士、兼任講師

maowang01@gmail.com
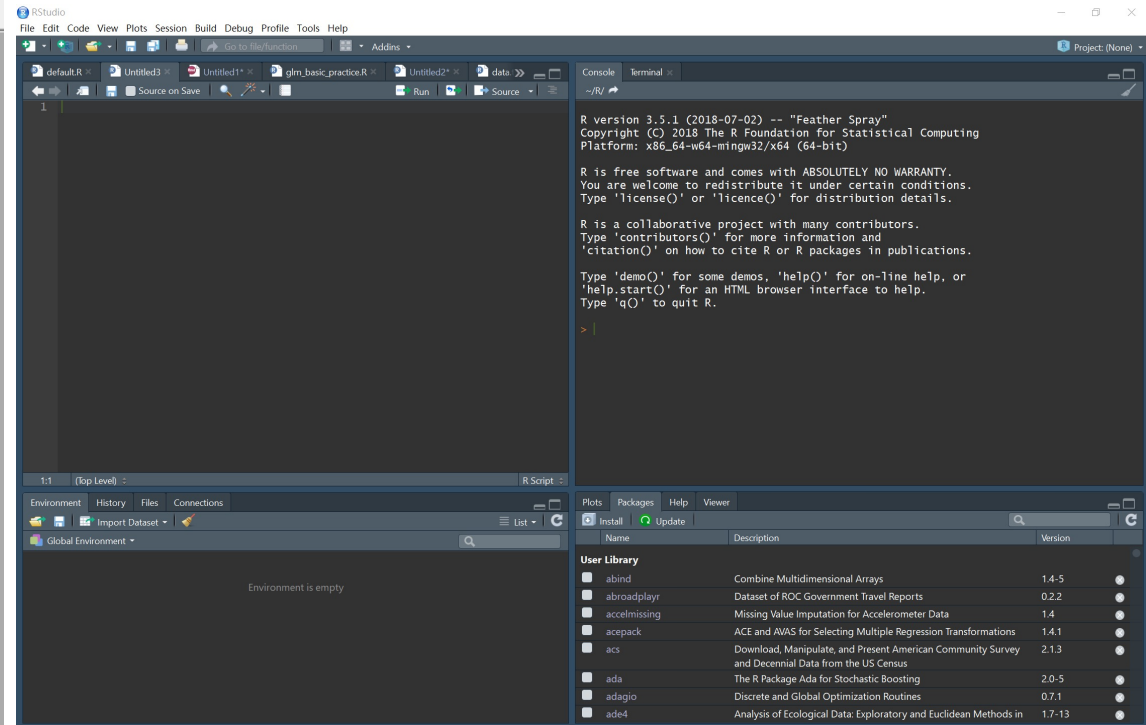
# 安裝 R 與 RStudio

- R
  R for Windows
  (https://cran.csie.ntu.edu.tw/bin/windows/base/)
  R for Mac
  (https://cran.csie.ntu.edu.tw/bin/macosx/)

- **RStudio（請先安裝R）**
  (https://www.rstudio.com/products/rstudio/download/#download)

- **https://bids.github.io/2019-01-17-bids/** (安裝介紹影片在最下方)

# R



# RStudio (IDE)

# 為什麼要學程式語言？

- 簡化繁瑣重複的工作
- make your life easier

# 為什麼程式語言要學 R？

- 免費！
- 功能強大（各種分析都可以辦到，繪圖尤其強大）。
- 易學（對人文社會科學背景者較容易）。
- 友善的學習社群。

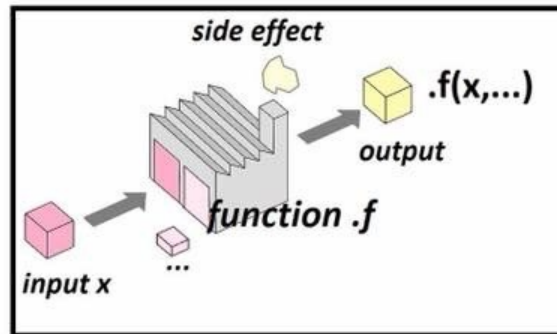# R 語言

- Ross Ihaka



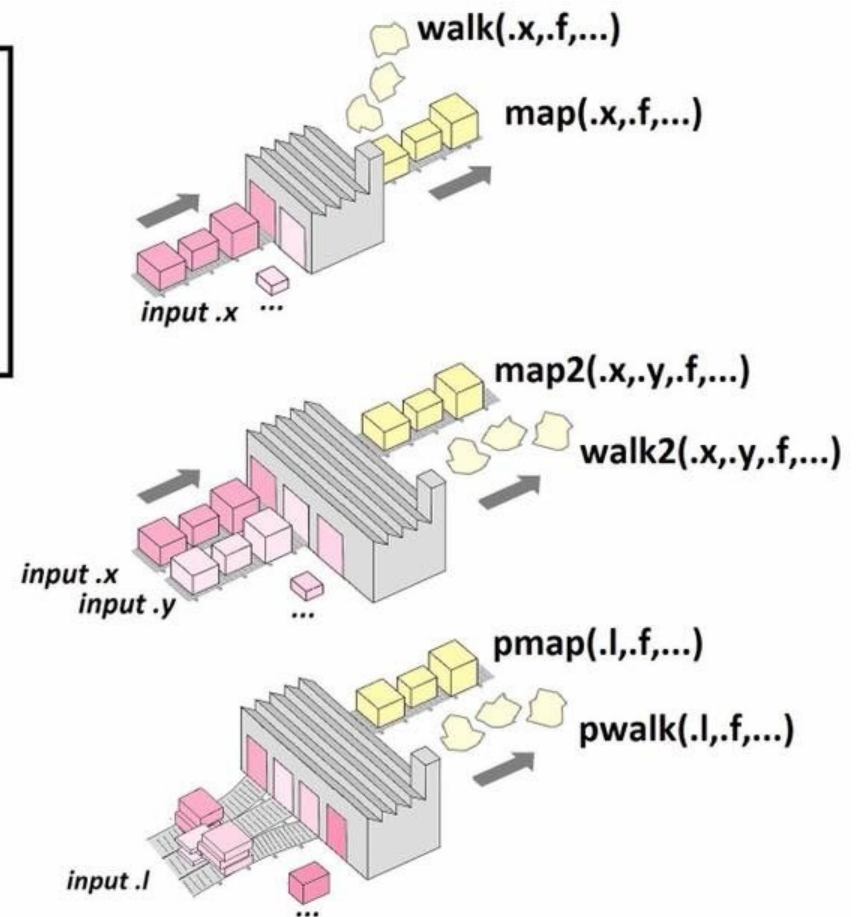- Robert C. Gentleman

# R 語言的基本要素

- 物件（Object）
- 函式（Function）



About computation in R

"To understand computations in R, two slogans are helpful:

- Everything that exists is an object.

- Everything that happens is a function call."
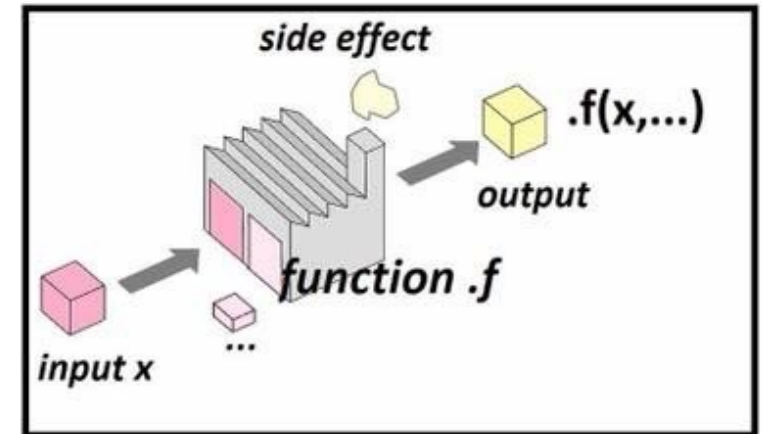
— John Chambers

# R 語言範例



- print(x,

  digits = getOption("digits"), ...)


- 參數（Arguments）

**object** an object for which a summary is desired.

**digits** minimal number of significant digits
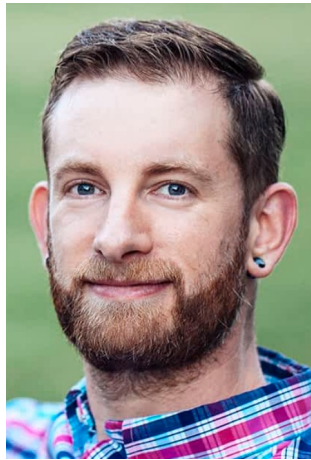
- R語言的程式碼是**有區分大小寫**（case sensitive）

# R 語言參數說明

- **必要參數**
- **x    a numeric vector, matrix or data frame**
- **預設參數**
- **y    NULL (default) or a vector, matrix or data frame with compatible dimensions to x. The default is equivalent to y = x (but more efficient).**

- 依參數位置（position）
- 依參數名稱

```
cor( x,
     y = NULL,
     use = "everything",
     method = c("pearson", "kendall", "spearman"))
```

# R 語言開發生態

- R Development Core Team (Base R)
- 其他套件（package）製作者

# R 能做什麼事？

- 統計分析
- 文字探勘
- 社會網絡分析
- 空間分析
- 網路爬蟲
- ……

## CRAN Task Views

CRAN task views aim to provide some guidance which packages on CRAN are relevant for tasks related to a certain topic. They give a brief overview of the included packages and can be automatically installed using the ctv package. The views are intended to have a sharp focus so that it is sufficiently clear which packages should be included (or excluded) - and they are *not* meant to endorse the "best" packages for a given task.

- To automatically install the views, the ctv package needs to be installed, e.g., via
  `install.packages("ctv")`
  and then the views can be installed via `install.views` or `update.views` (where the latter only installs those packages are not installed and up-to-date), e.g.,
  `ctv::install.views("Econometrics")`
  `ctv::update.views("Econometrics")`
- The task views are maintained by volunteers. You can help them by suggesting packages that should be included in their task views. The contact e-mail addresses are listed on the individual task view pages.
- For general concerns regarding task views contact the ctv package maintainer.

### Topics

| | |
|---|---|
| Bayesian | Bayesian Inference |
| ChemPhys | Chemometrics and Computational Physics |
| ClinicalTrials | Clinical Trial Design, Monitoring, and Analysis |
| Cluster | Cluster Analysis & Finite Mixture Models |
| Databases | Databases with R |
| DifferentialEquations | Differential Equations |
| Distributions | Probability Distributions |
| Econometrics | Econometrics |
| Environmetrics | Analysis of Ecological and Environmental Data |
| ExperimentalDesign | Design of Experiments (DoE) & Analysis of Experimental Data |
| ExtremeValue | Extreme Value Analysis |
| Finance | Empirical Finance |
| FunctionalData | Functional Data Analysis |
| Genetics | Statistical Genetics |
| Graphics | Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization |
| HighPerformanceComputing | High-Performance and Parallel Computing with R |
| Hydrology | Hydrological Data and Modeling |
| MachineLearning | Machine Learning & Statistical Learning |

https://cran.r-project.org/web/views/

# 為什麼有人覺得 R 不好學？

- **難處1：要記各種函式（function）**
- 解答：沒有人真的能全部記住，重點是該函式的說明文件（documentation）清楚嗎？

- **難處2：入門的門檻不低，要學的套件（package）太多！**
- 解答：理解基本的資料結構，可以降低後續學習的門檻；學習具有同樣設計邏輯的套件（如tidyverse）。

# 怎麼問問題？

- **Repr**oducible **ex**ample (**reprex**)
  https://github.com/tidyverse/reprex

- 讓別人能最小化的**重現**你的問題，才能
  夠幫你處理問題。
  - A **minimal dataset**, necessary to
    reproduce the error
  - The **minimal runnable code** necessary
    to reproduce the error, which can be run
    on the given dataset.

https://stackoverflow.com/questions/5963269/how-to-make-a-
great-r-reproducible-example

# Coding style

- [https://style.tidyverse.org/](https://style.tidyverse.org/)
- 讓你自己及合作者更容易讀你的code。

# 其他學習資源

- DataCamp

- Coursera

- edX

- Cheat sheet
  https://www.rstudio.com/resources/cheatsheets/

# R 語言基礎

- [R語言基礎：簡介](#)
- R語言基礎：資料篩選與整理
- R語言基礎：資料探索與分析
- R語言基礎：程式設計基礎

- [R語言進階：資料整理](#)

**The Social Science Data Lab at UC Berkeley**

# 實際操作

• R Markdown files

| | | | |
|---|---|---|---|
| 📁 .git | 2018/11/11 下午 08:... | 檔案資料夾 | |
| 📁 .Rproj.user | 2018/11/11 下午 08:... | 檔案資料夾 | |
| 📁 data | 2018/11/11 下午 08:... | 檔案資料夾 | |
| 📄 | 2018/11/11 下午 08:... | 文字文件 | 1 KB |
| 🖼 hfs | 2018/11/11 下午 08:... | PNG 檔案 | 131 KB |
| 📄 LICENSE | 2018/11/11 下午 08:... | 檔案 | 14 KB |
| Ⓡ Part 4 script | 2018/11/11 下午 08:... | R 檔案 | 2 KB |
| Ⓡ R Fundamentals Part 1 Introduction | 2018/11/11 下午 08:... | RMD 檔案 | 34 KB |
| Ⓡ R Fundamentals Part 2 Subsetting and reshaping | 2018/11/11 下午 08:... | RMD 檔案 | 20 KB |
| Ⓡ R Fundamentals Part 3 Data exploration and analysis | 2018/11/11 下午 08:... | RMD 檔案 | 26 KB |
| Ⓡ R Fundamentals Part 4 Project | 2018/11/11 下午 08:... | RMD 檔案 | 3 KB |
| Ⓡ R Fundamentals What is R markdown | 2018/11/11 下午 08:... | RMD 檔案 | 3 KB |
| 🦊 R_Fundamentals_Bonus_-_For-loops_and_functions | 2018/11/11 下午 08:... | Firefox HTML Docu... | 1,434 KB |
| 🦊 R_Fundamentals_Part_1_Introduction | 2018/11/11 下午 08:... | Firefox HTML Docu... | 905 KB |
| 🦊 R_Fundamentals_Part_2_Subsetting_and_reshaping | 2018/11/11 下午 08:... | Firefox HTML Docu... | 1,077 KB |
| 🦊 R_Fundamentals_Part_3_Data_exploration_and_analysis | 2018/11/11 下午 08:... | Firefox HTML Docu... | 1,756 KB |
| 🦊 R_Fundamentals_Part_4_Project | 2018/11/11 下午 08:... | Firefox HTML Docu... | 853 KB |
| 🦊 R_Fundamentals_What_is_R_markdown | 2018/11/11 下午 08:... | Firefox HTML Docu... | 1,167 KB |
| 📄 README.md | 2018/11/11 下午 08:... | MD 檔案 | 5 KB |
| Ⓡ R-Fundamentals | 2018/11/11 下午 08:... | R Project | 1 KB |
| Ⓡ solutions | 2018/11/11 下午 08:... | RMD 檔案 | 15 KB |

# 作業範例

- Introduction to R
https://www.datacamp.com/courses/free-introduction-to-r

- Intermediate R
https://www.datacamp.com/courses/intermediate-r

```
1  ---
2  title: "HW3"
3  author: "Your_name"
4  date: "2019/4/8"
5  output: html_document
6  ---
7
8  ```{r setup, include=FALSE}
9  knitr::opts_chunk$set(echo = TRUE)
10 ```
11
12 ```{r message=FALSE}
13 library(tidyverse)
14 ```
15
16 ## 1. 英文資料前處理與斷詞
17
18 請記得把分析的檔案下載到你的工作目錄。資料取自 [Martijn Schoonvelde Dataverse](https://dataverse.harvard.edu/file.xhtml?persistentId=doi:10.7910/DVN/2PNZNU/I0I7GM&version=2.0)，COOL課程平台也有放資料檔(comb.corpus.Rdata)。
19
20 ```{r load}
21 # 使用load讀Rdata檔，會在讀取的同時，自動創建原始檔案的object，所以不用重新assign
22 load("data/comb.corpus.Rdata")
23 glimpse(corpus)
24 corpus <- as_tibble(corpus)
25 ```
26
27 請使用stringr的函式，進行文字資料前處理。
28
29 1. 篩選country欄位中，含有"DK"的觀察值（請注意：原始資料中有"DK"與"DK"都要選取），並存成 `dk_corpus`。
30
31 ```{r dk_corpus}
32 unique(corpus$country)
33
34 ```
35
36 2. 請將所有存在於text欄位的數字都移除，同樣assign回 `dk_corpus`。
37
```

# 資料類型

- 數值（numeric）
- 字串（character, string）
- 邏輯判斷（logical）

# 儲存格式

|    | Homogeneous | Heterogeneous |
|----|-------------|---------------|
| 1d | Atomic vector | List |
| 2d | Matrix | Data frame |
| nd | Array | |

# 什麼是文字探勘？

- Github repository
- https://github.com/aleszu/textanalysis-shiny

- Shiny App實做
  https://storybench.shinyapps.io/textanalysis/