

BIG Data Application

(Business Management Seminar)

Ming-Chang Lee, Ph.D.

Department and Graduate Institute of Information Management

Yu Da University of Science and Technology

November 19, 2018

WORKING TYPES
SIVELY GICAL REDS CREATED

DESKTOP CURRENTLY FC WORLD'S TENS CAPACITY FORMS PRESENTATIONS NEW PRACTITIONERS

COMPLEX DIFFICULTY TARGET ABILITY SENSOR ARCHIVES

TERABYTES MP RESEARCH DATABASES SETS EXAMPLES APPLIED DESCRIPTING AMOUNT

ONE SINCE USE SOFTWARE BUSINESS MOVING UBQUITOUS

RADIO-FREQUENCY SOLID COMPLEXITY WIRELESS

TOLERABLE SAN PARALLEL SIZE NEEDED QUALITIES PETABYTES

PROCESSING LOGS

INTERNET TECHNOLOGIES USED DISTRIBUTED CAPTURE MANAGE GROW

MAY MANAGEMENT SOCIAL LARGE EVERY LARGER

DEFINING CASE STORE

CURRENT ELAPSED THOUGHT

<http://rwepa.blogspot.com/>



The screenshot shows the homepage of the RWEPA blog. At the top left is the RWEPA logo with the text "Since 2013". To its right is a block of R code for creating a map of Taiwan. Below the code are two maps: a scatter plot and a choropleth map of Taiwan. To the right of the maps is a welcome message in Chinese and English, followed by more R code. The main content area features a heading "主題式地圖(Thematic map) - 政府開放資料為例" with an image of a thematic map of Taiwan. On the left, there's a "G+" button. On the right, there's a search bar and a sidebar with a "快速連結" section containing links to various R-related resources. A red arrow points from the bottom right towards the "關於作者" link in the sidebar.

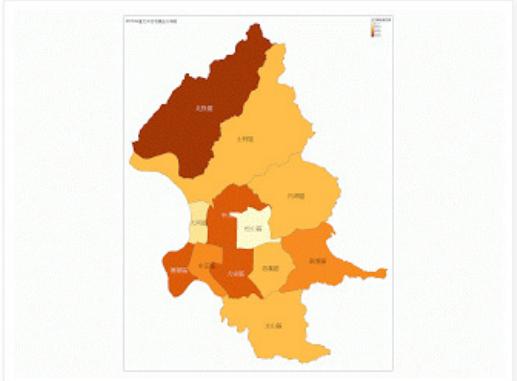
歡迎光臨 RWEPA 部落格
Welcome to RWEPA blog

歡迎來到 RWEPA blog, 成立宗旨為提供免費R軟體的相關資訊。R包括大量套件可應用於不同領域, 例如: 2D/3D互動式繪圖, 資料視覺化, 資料探勘, 線性與非線性最佳化問題, 時間序列, 空間資料, 財務分析, 多變量分析, 問卷調查, 實驗設計, 統計製程管制, 存活分析, 臨床實驗分析, 社會網絡分析, 生物資訊, 醫學統計等。

G+

2018年10月27日 星期六

主題式地圖(Thematic map) - 政府開放資料為例



主題式地圖
Thematic map
開放式資料
open data
地圖資料與社會經濟資料合併
rgdal 套件
tmap 套件

搜尋此網誌

快速連結

- ★★★R教學-基礎篇
- ★★★R軟體教學影片
- ★R-3.5.1-Wndows下載
- ★RStudio-1.1.463下載
- ★RStudio Daily
- R使用者論壇
- R
- RStudio - IDE整合工具
- R-bloggers
- 中華R軟體學會(CARS)
- 臺灣資料科學與商業應用協會(DSBA)
- R-Taichi 太極(USA)
- 育達科技大學資訊管理系
- WEPA
- 關於作者

About Me



- Ph. D., Department of Industrial and Systems Engineering, Chung Yuan Christian University.



- Shipping analyst
- Marketing analyst
- Invited Talks: +1800 hours



- Data mining
- Deep learning
- Visualization
- Optimization



- R/Python
- Office / VBA
- Database/SQL

Content

1. Big Data
2. R / RStudio
3. Big Data Application

1. Big Data

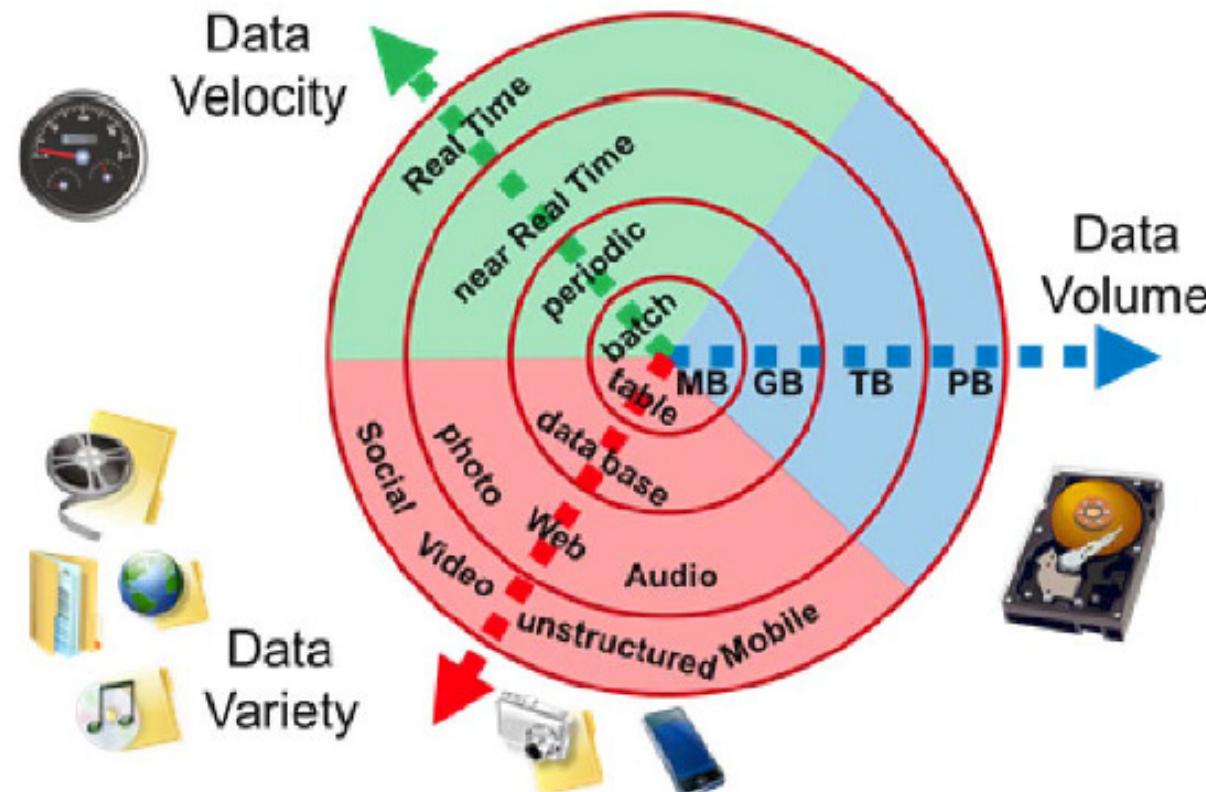
Big Data?

- Data sets that are too **large** or **complex** for traditional data-processing application software to **adequately** deal with.

Reference: https://en.wikipedia.org/wiki/Big_data

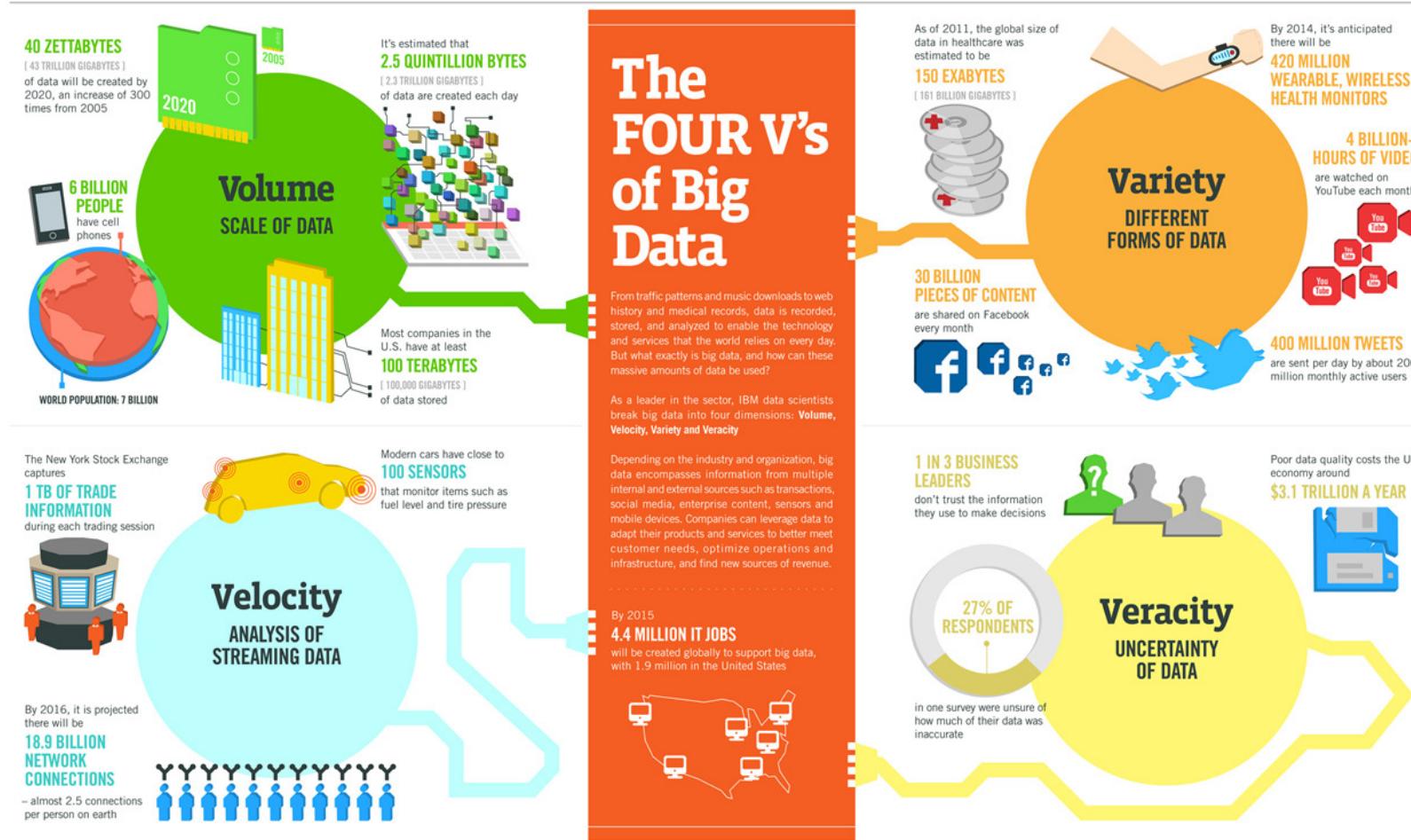
Big Data - 3V

- Doug Laney, 2001:



Reference: Lancy, D., 3D Data Management: Controlling Data Volume, Velocity, and Variety, Gartner, 2001.

IBM Big Data - 4V



Reference: <http://www.ibmbigdatahub.com/infographic/four-vs-big-data>

Veracity

- Veracity refers to the quality or trustworthiness of the data.
- A common complication is that the data is saturated with both useful signals and lots of noise (data that can't be trusted).

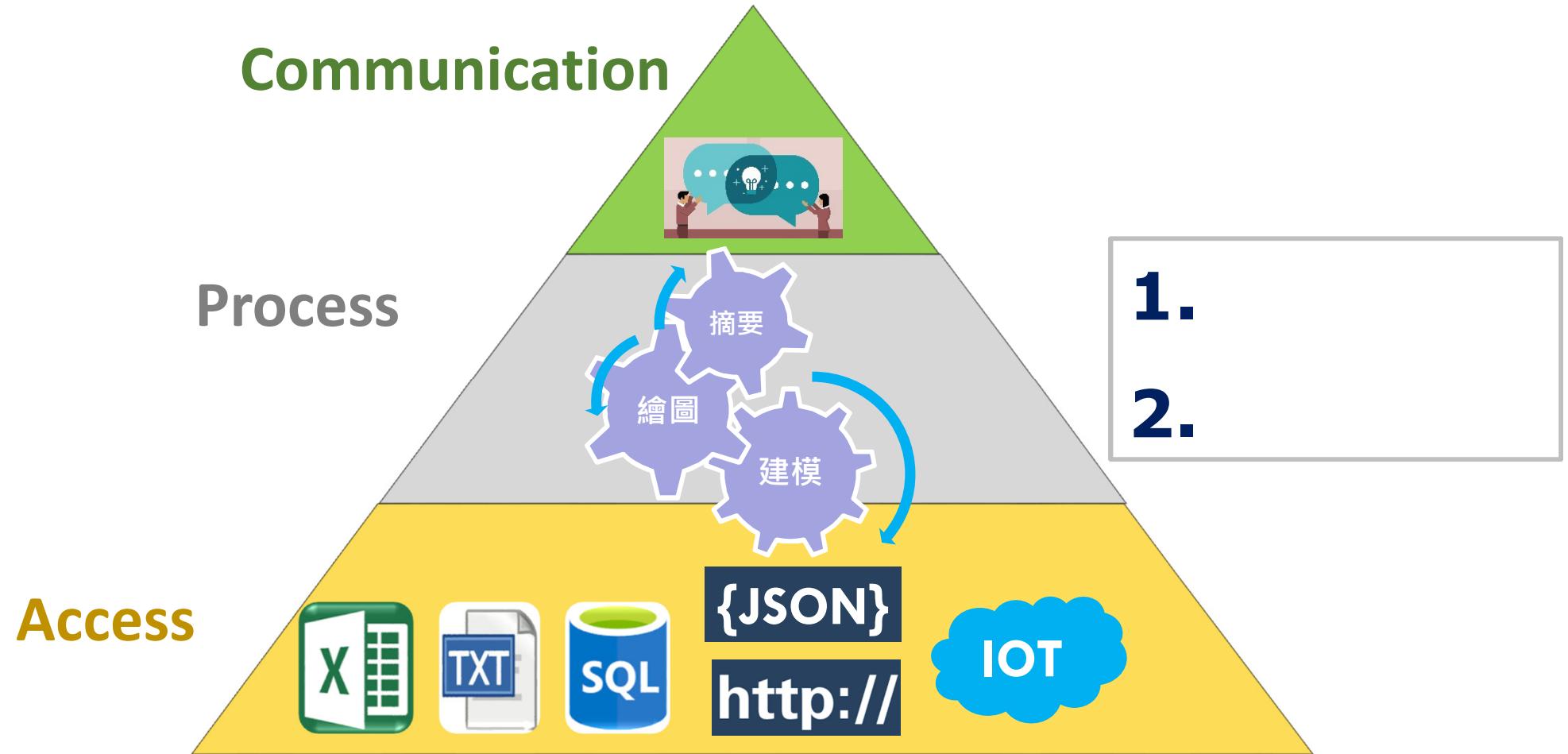
Why Big Data now?

- More data are being collected and stored
- Open sources
- Commodity hardware
- Cloud
- IOT, ...

Big Data Challenges

- Capturing data
- Data storage
- Data analysis, querying
- Search
- Sharing
- Transfer
- Visualization
- Communication, ...

APC approach



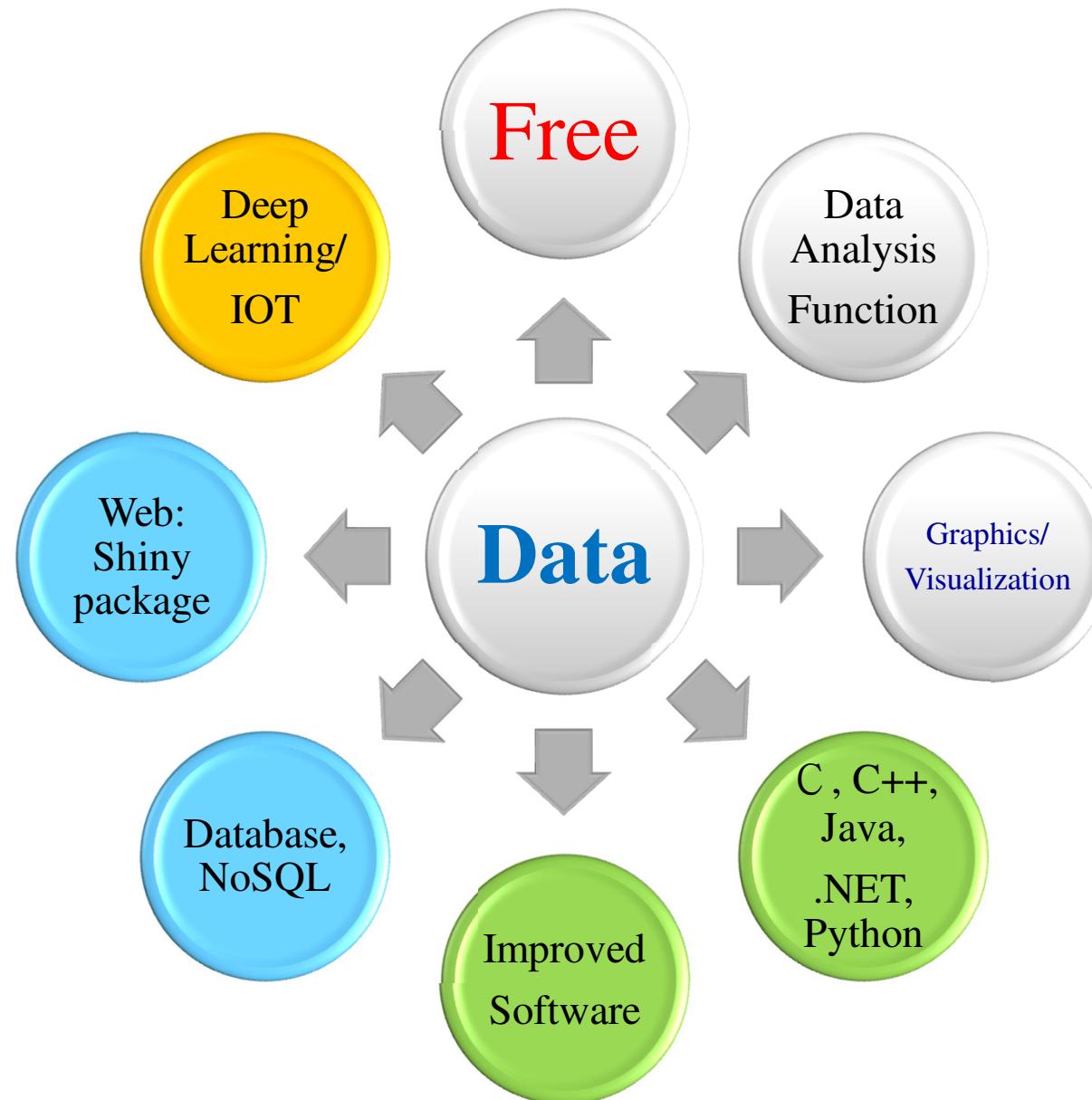
2. R / RStudio

R!

- 1976 - John Chambers, Rick Becker, and Allan Wilks developed S language at Bell Labs.
- 1993 - Ross Ihaka and Robert Gentleman developed R at University of Auckland in New Zealand.
 - R is a programming language with the features:
 - Statistical analysis
 - Graph
 - Visualization
- 1997年 – R development core team
 - February 2000 – R 1.0.0
 - March 2013 – R 2.15.3
 - July 2018 – R 3.5.1



R - Features



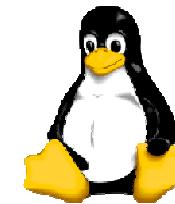
R- Download

- Web : <http://www.r-project.org/>
- Download \ CRAN
- Taiwan CRAN

<https://ftp.yzu.edu.tw/CRAN/>
<http://ftp.yzu.edu.tw/CRAN/>
<http://cran.csie.ntu.edu.tw/>

- Download R for Windows

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)



R- Download

- base:

[Download R 3.5.1 for Windows](#) (62 megabytes, 32/64 bit)

- R install guide, 2006
 - http://web.ydu.edu.tw/~alan9956/docu/refer/R01_install.pdf
- Installed path
 - C:\R\R-3.5.1

<http://www.r-project.org/>

The R Project for Statistical Computing

Download



[Home]

Download

CRAN

R Project

About R

Logo

Contributors

What's New?

Reporting Bugs

Development Site

Conferences

Search

Manuals

R Foundation

Foundation

Board

Members

Donors

Donate

Help With R

Getting Help

Documentation

Manuals

FAQs

The R Journal

Books

Certification

Other

Links

Bioconductor

Related Projects

Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To [download R](#), please choose your preferred [CRAN mirror](#).

Free

If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

News

- You can now support the R Foundation with a renewable subscription as a [supporting member](#)
- **R version 3.5.1 (Feather Spray)** has been released on 2018-07-02.
- The R Foundation has been awarded the Personality/Organization of the year 2018 award by the professional association of German market and social researchers.

News via Twitter

 The R Foundation

@_R_Foundation



If you would like to support the R Foundation financially as a supporting member you can now choose an annual subscription that renews automatically. Or you can choose a one-off donation as before. For more details see [r-project.org/foundation/don...](#)



Sep 13, 2018

 The R Foundation Retweeted



 Peter Dalgaard

@pdalgaard

Sep 13, 2018

Genetic analysis

 useR! 2019

@UseR2019_Conf

Jul 28, 2018

R - Manuals

The R Manuals

edited by the R Development Core Team.

The following manuals for R were created on Debian Linux and may differ from the manuals for Mac or Windows on platform-specific pages, but most parts version of the manuals for each platform are part of the respective R installations. The manuals change with R, hence we provide versions for the most recent version for the patched release version (R-patched) and finally a version for the forthcoming R version that is still in development (R-devel).

Here they can be downloaded as PDF files, EPUB files, or directly browsed as HTML:

Manual	R-release	R-patched
An Introduction to R is based on the former "Notes on R", gives an introduction to the language and how to use R for doing statistical analysis and graphics.	HTML PDF EPUB	HTML PDF EPUB
R Data Import/Export describes the import and export facilities available either in R itself or via packages which are available from CRAN.	HTML PDF EPUB	HTML PDF EPUB
R Installation and Administration	HTML PDF EPUB	HTML PDF EPUB
Writing R Extensions covers how to create your own packages, write R help files, and the foreign language (C, C++, Fortran, ...) interfaces.	HTML PDF EPUB	HTML PDF EPUB
A draft of The R language definition documents the language <i>per se</i> . That is, the objects that it works on, and the details of the expression evaluation process, which are useful to know when programming R functions.		
R Internals : a guide to the internal structures of R and coding standards for the core team working on R itself.		
The R Reference Index : contains all help files of the R standard and recommended packages in printable form. (9MB, approx. 3500 pages)		

**contributed documentation
(Free!)**

Translations of manuals into other languages than English are available from the [contributed documentation](#) section (only a few translations are available).

R - Manuals

Contributed Documentation

Note: The CRAN area for contributed documentation is frozen and no longer actively maintained.

English --- Other Languages

Manuals, tutorials, etc. provided by users of R. The R core team does not take any responsibility for contents, but we appreciate the effort very much and encourage everybody to contribute to this list! To submit, follow the submission instructions on the [CRAN main page](#). All material below is available directly from CRAN, you may also want to look at the list of [other R documentation](#) available on the Internet.

Note: Please use the [directory listing](#) to sort by name, size or date (e.g., to see which documents have been updated lately).

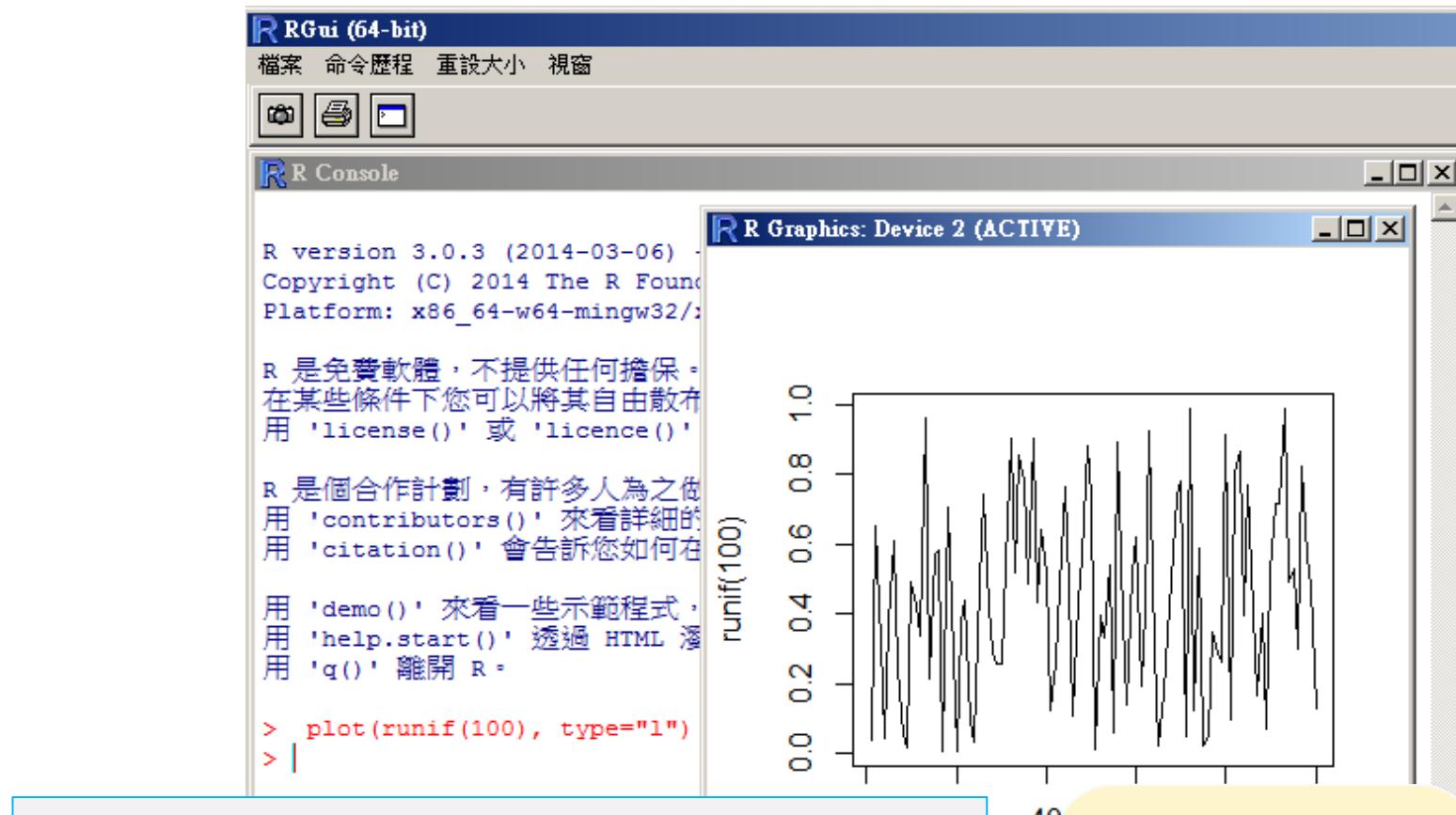
English Documents

Documents with more than 100 pages:



- “**Visual Statistics. Use R!**” by Alexey Shipunov ([PDF](#), 2017-05-15, 388 pages) materials are accessible from [Alexey Shipunov's English R page](#).
- “**Using R for Data Analysis and Graphics - Introduction, Examples and Commentary**” by John Maindonald ([PDF](#), data sets and scripts are available at [JM's homepage](#)).
- “**Practical Regression and Anova using R**” by Julian Faraway ([PDF](#), data sets and scripts are available at the [book homepage](#)).

R - Practice



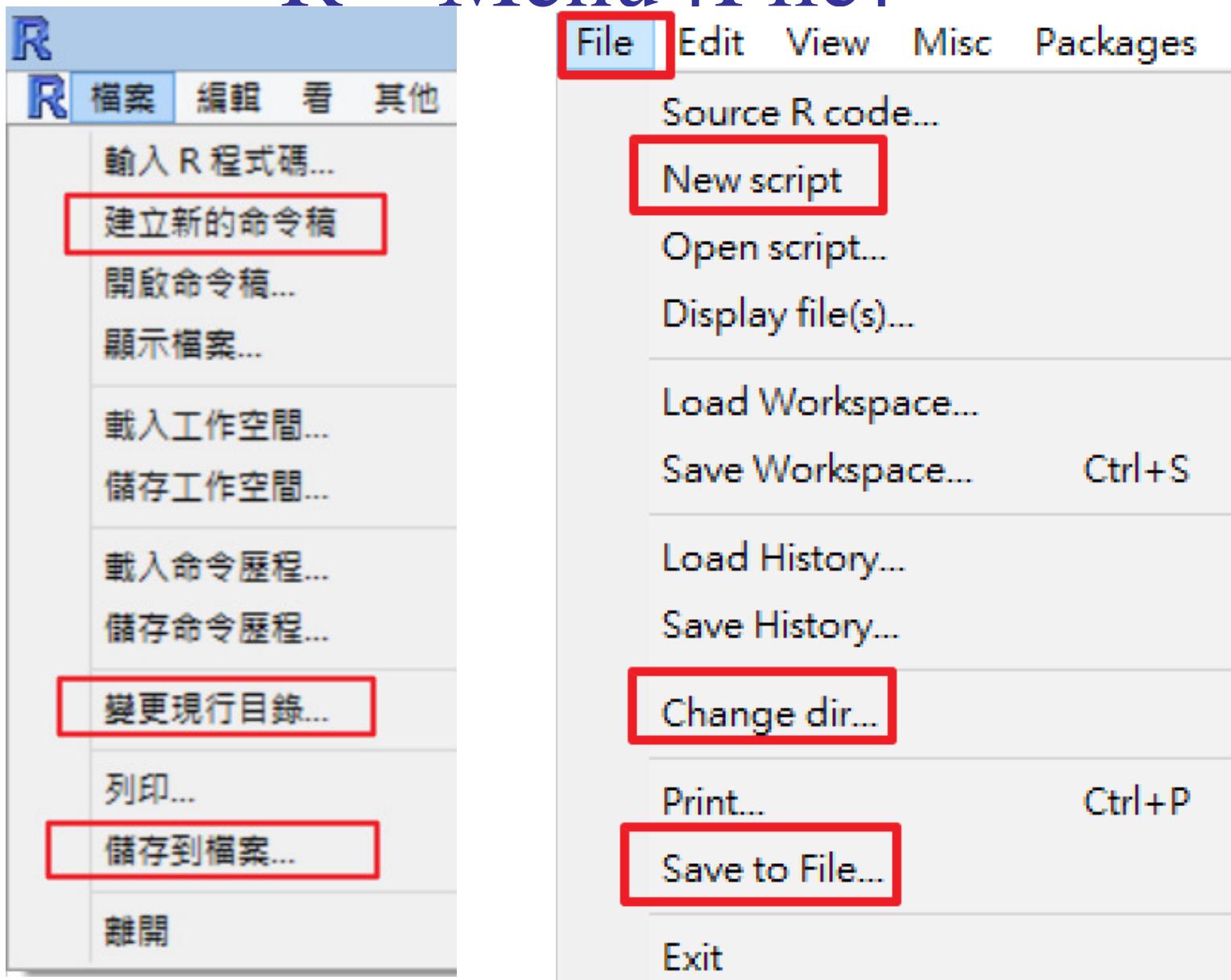
```
plot(runif(100), type="l")
```

```
demo(graphics)
```

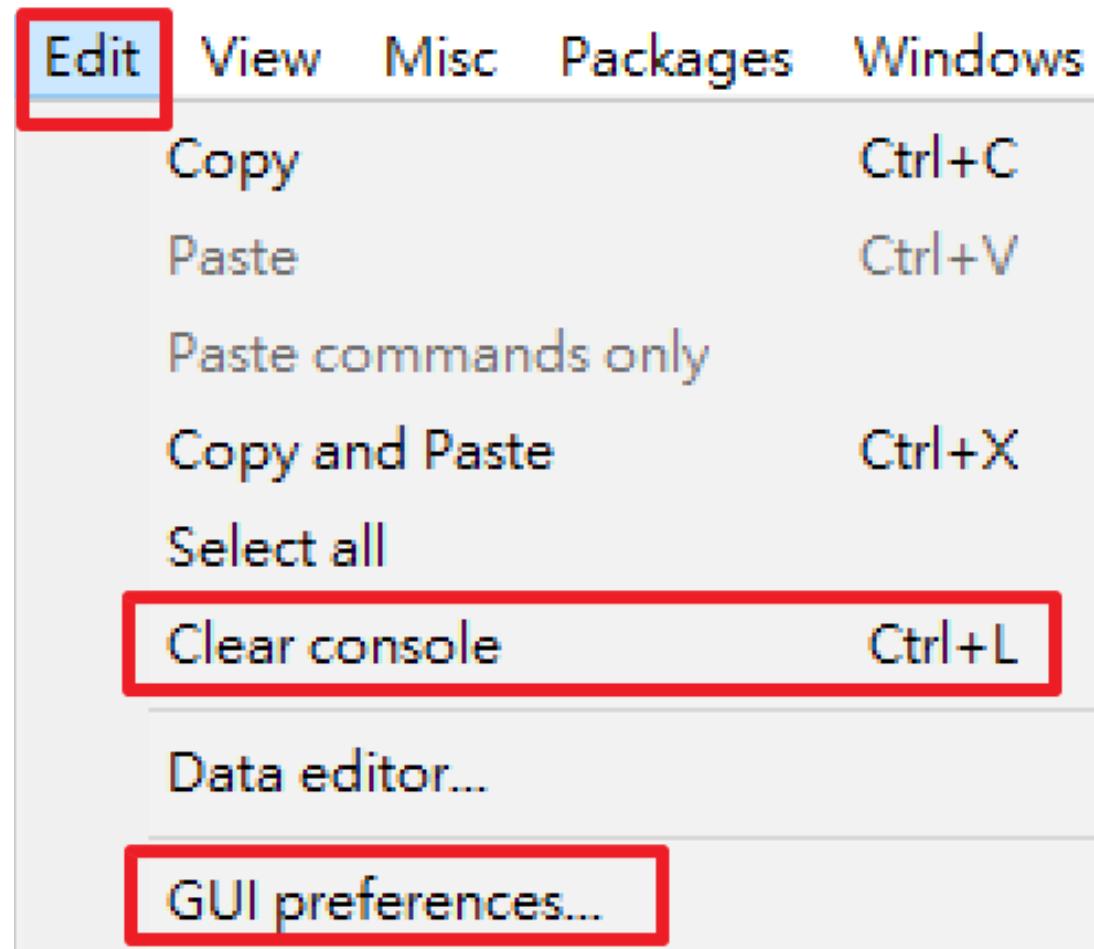
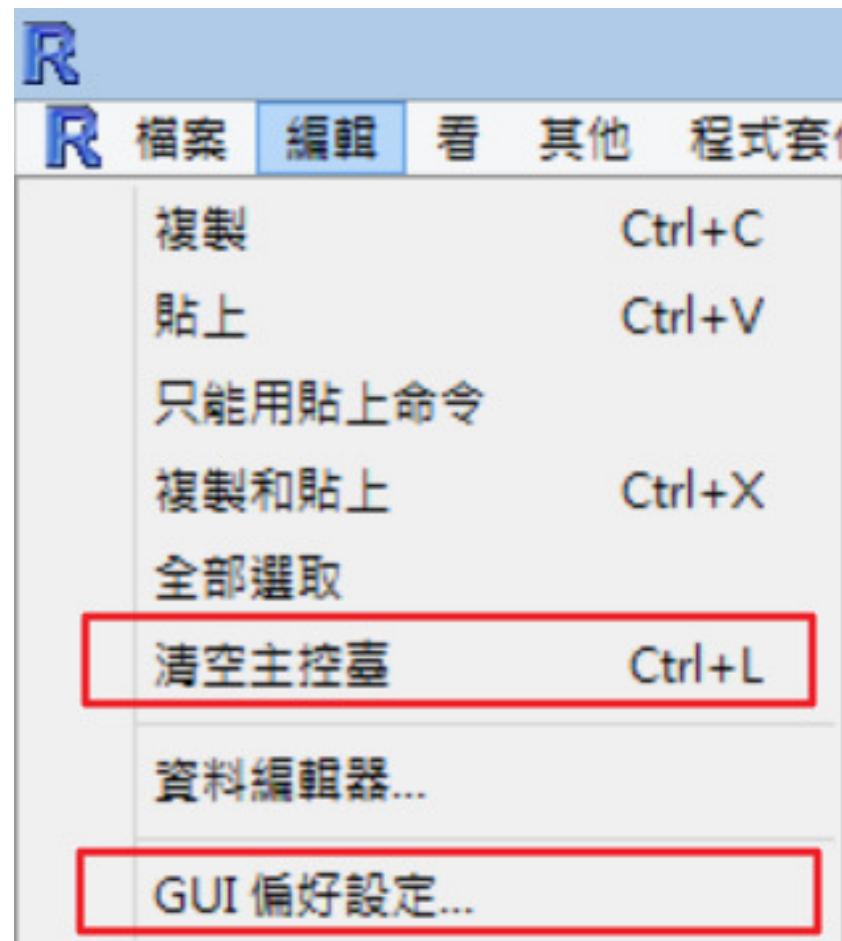
```
@RV demo(persp)
```

Notice:
Upper and lower case

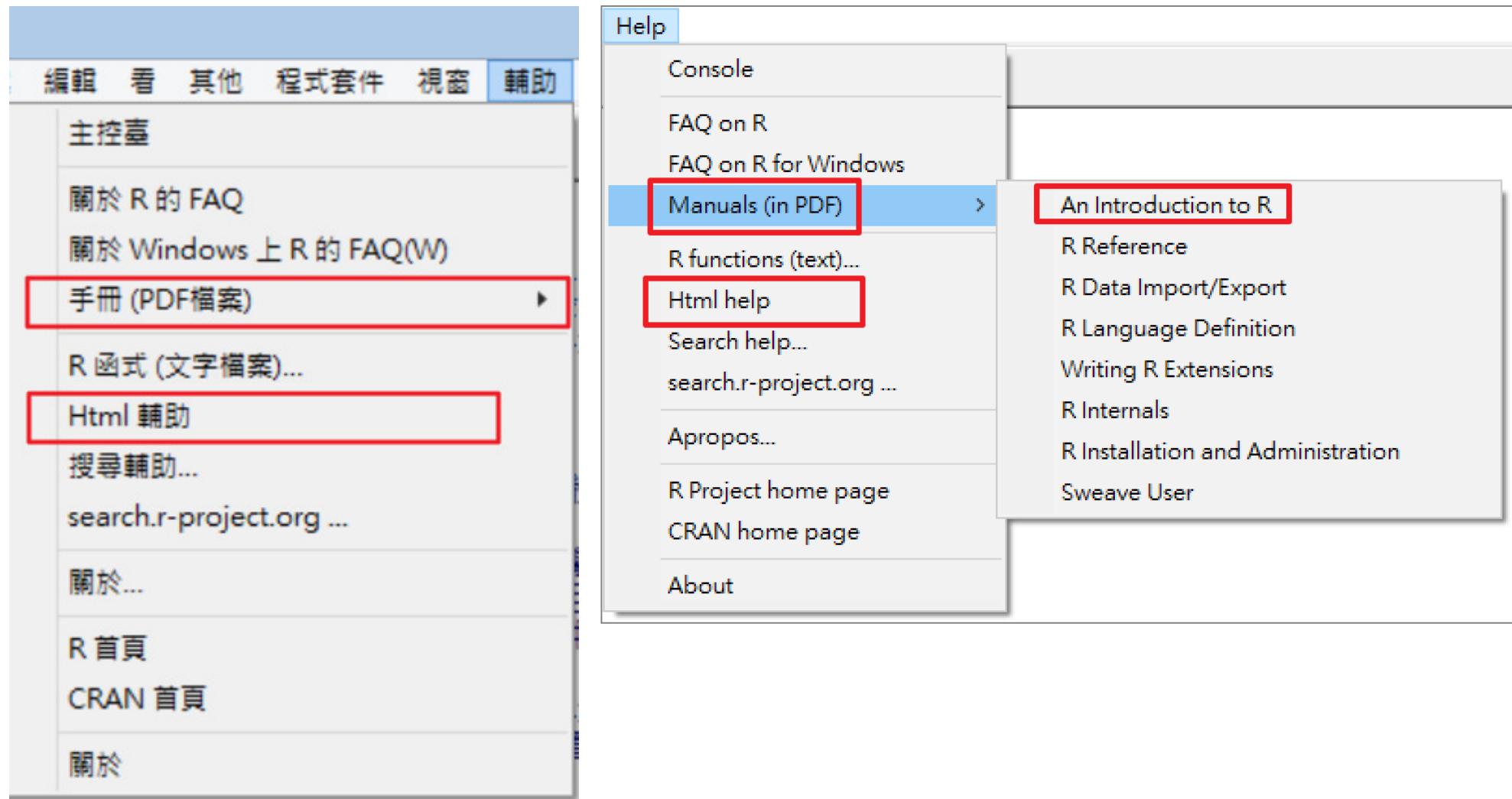
R – Menu [File]



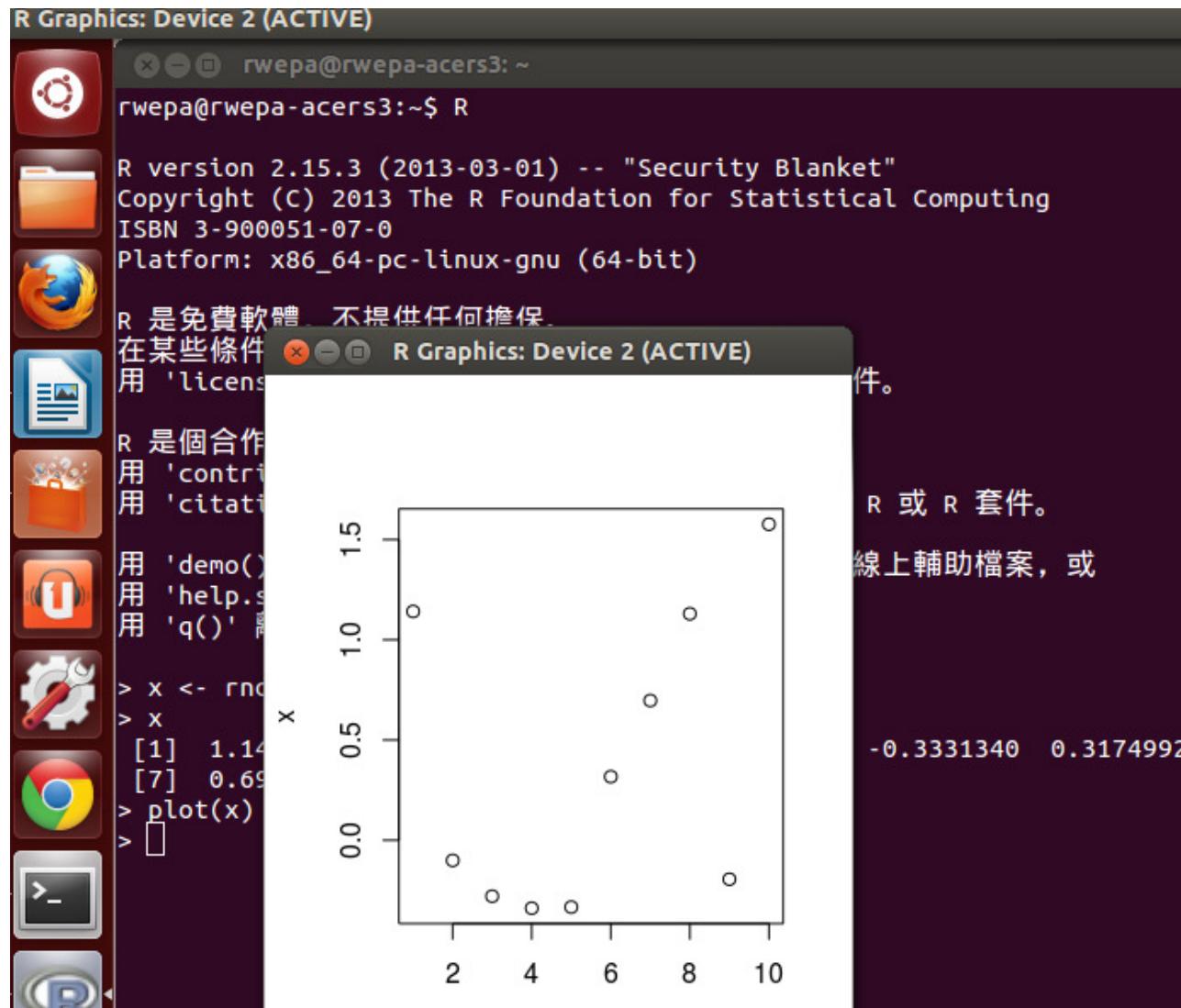
R – Menu [Edit]



R – Menu [Help]



Ubuntu - R

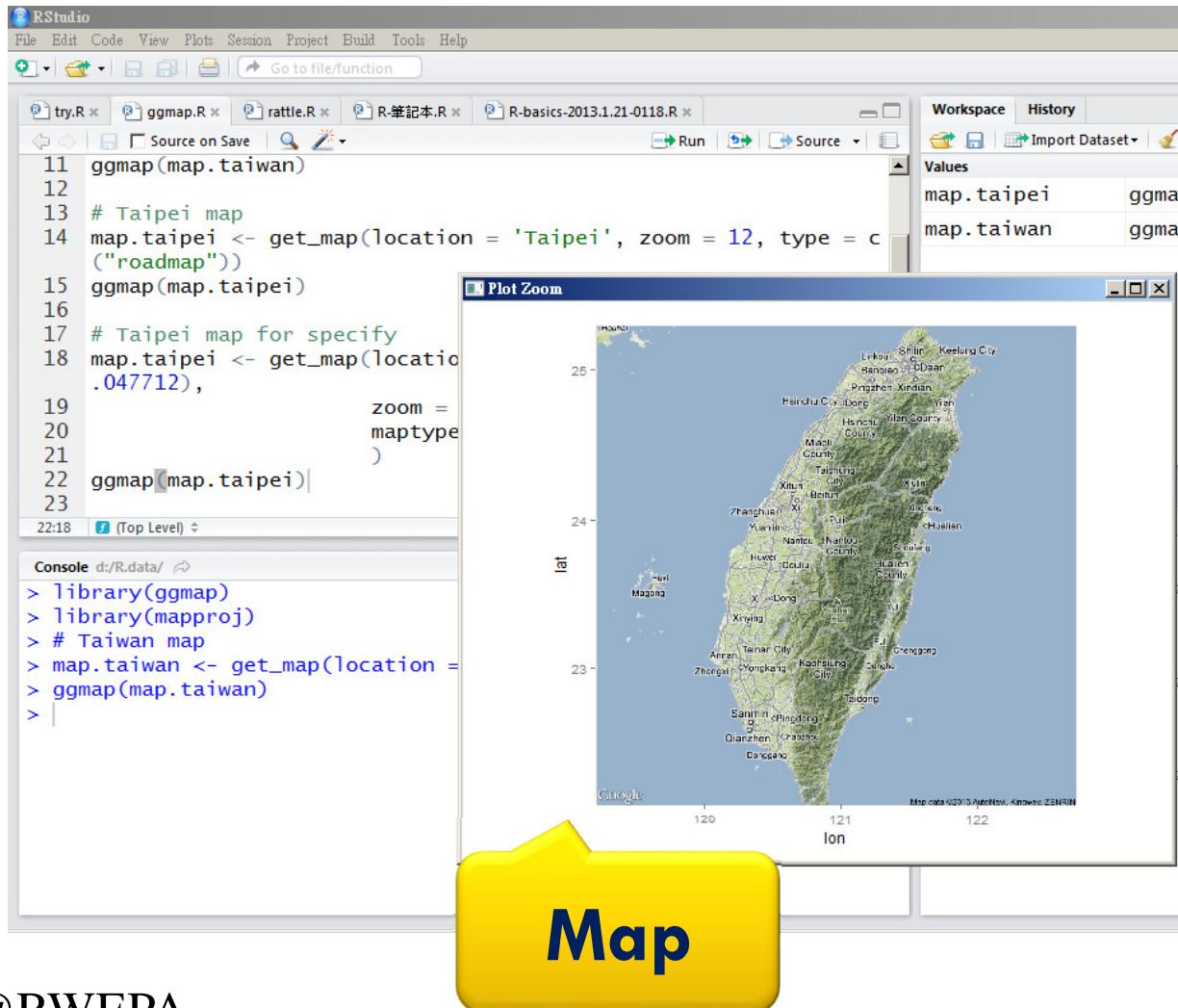


<http://rwepa.blogspot.tw/2013/05/ubuntu-r.html>

RStudio

Integrated Development Environment- RStudio

■ <http://www.rstudio.com/>



RStudio interface showing R code for mapping Taiwan:

```

11 ggmap(map.taiwan)
12
13 # Taipei map
14 map.taipei <- get_map(location = 'Taipei', zoom = 12, type = c("roadmap"))
15 ggmap(map.taipei)
16
17 # Taipei map for specify
18 map.taipei <- get_map(location = "Taipei", zoom = 12, maptype = "satellite")
19
20
21
22 ggmap(map.taipei)
23

```

Console output:

```

> library(ggmap)
> library(mapproj)
> # Taiwan map
> map.taiwan <- get_map(location = "Taiwan", zoom = 12, maptype = "satellite")
> ggmap(map.taiwan)

```

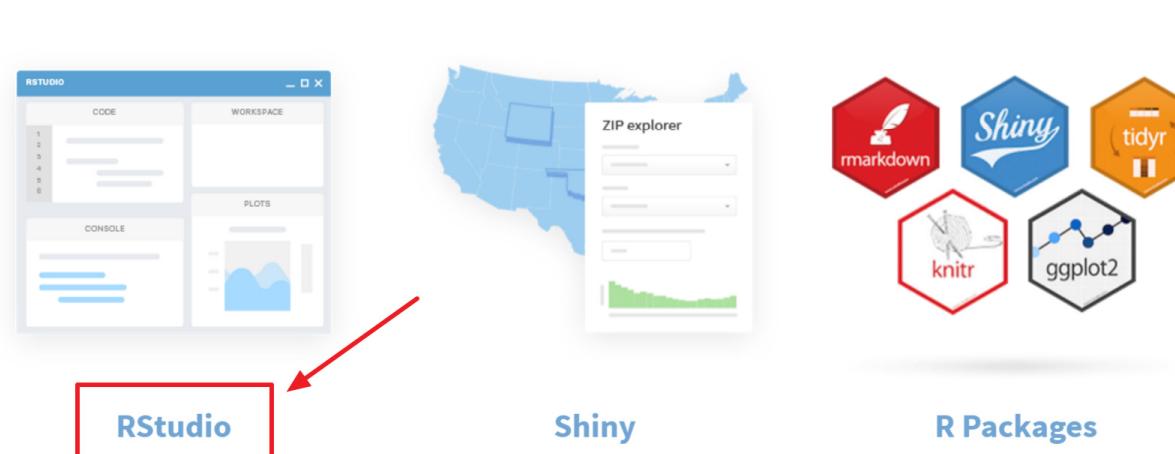
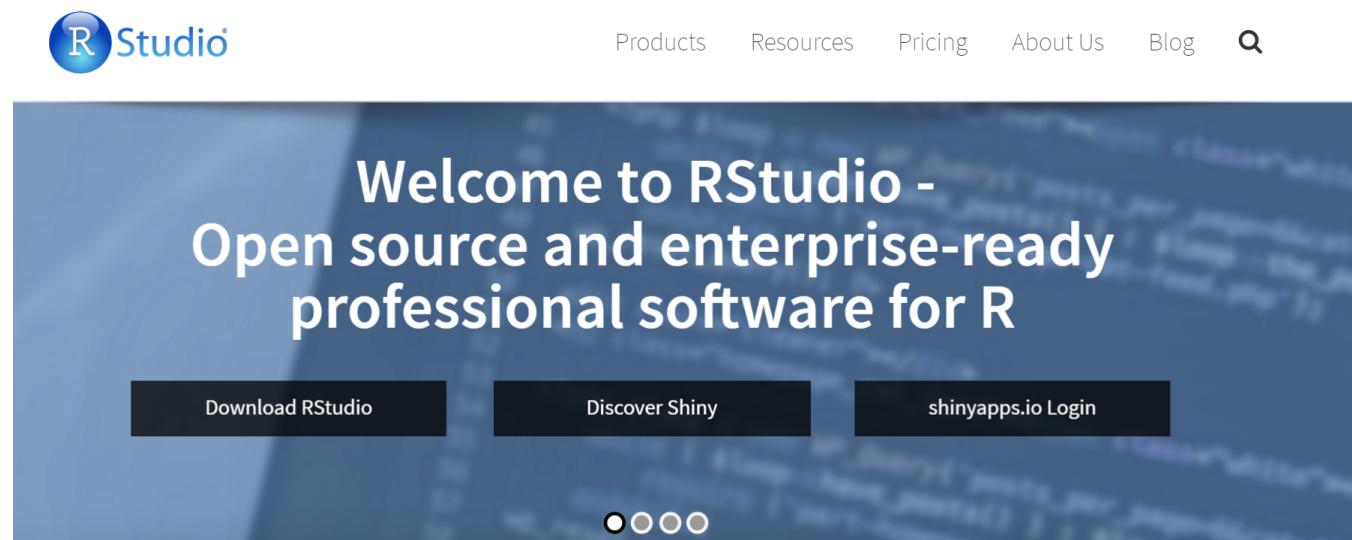
A yellow callout bubble labeled "Map" points to the map of Taiwan.



RStudio - Download

- <http://www.rstudio.com/>

Try: **RStudio daily**



RStudio - Download



Choose Your Version of RStudio

	Desktop	Server	
RStudio Desktop (Free License)		RStudio Server (Free License)	RStudio Server Pro (Commercial License)
Integrated Development Environment for R	✓	✓	✓
Priority support	✓		✓
Access via Web Browser		✓	✓
Enterprise Security and Access Controls			✓
Project Sharing			✓
Access to Multiple Versions of R			✓
Multiple Concurrent Sessions			✓
Administrative Dashboard			✓
Load Balancing and Resource Management			✓
License	AGPL	Commercial	AGPL
	DOWNLOAD	DOWNLOAD	DOWNLOAD



RStudio - Download

RStudio Desktop 1.1.463 — Release Notes

RStudio requires R 3.0.1+. If you don't already have R, download it [here](#).

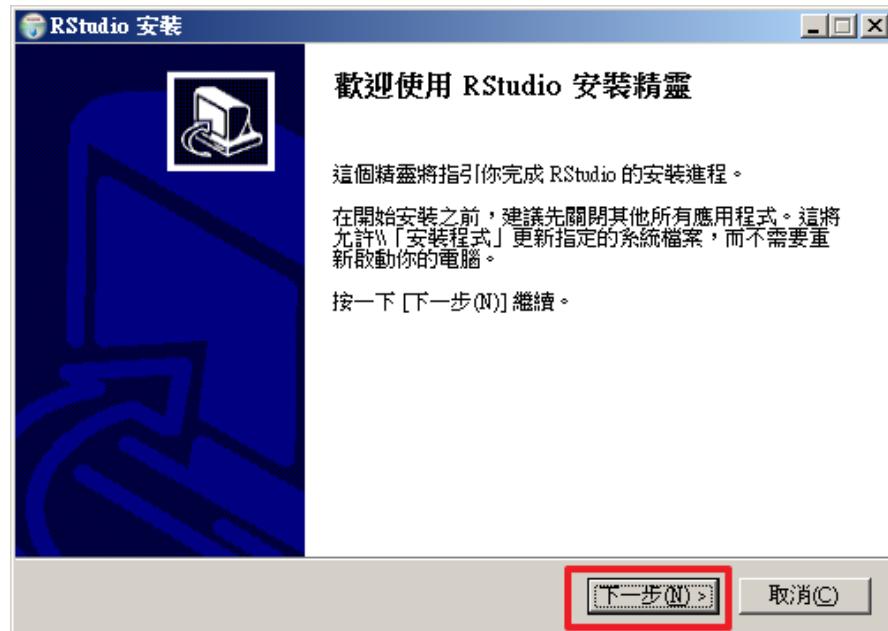
Linux users may need to [import RStudio's public code-signing key](#) prior to installation, policy.

Installers for Supported Platforms

Installers	Size	Date
RStudio 1.1.463 - Windows Vista/7/8/10	85.8 MB	2018-10-29
RStudio 1.1.463 - Mac OS X 10.6+ (64-bit)	74.5 MB	2018-10-29
RStudio 1.1.463 - Ubuntu 12.04-15.10/Debian 8 (32-bit)	89.3 MB	2018-10-29

RStudio - Installation

1



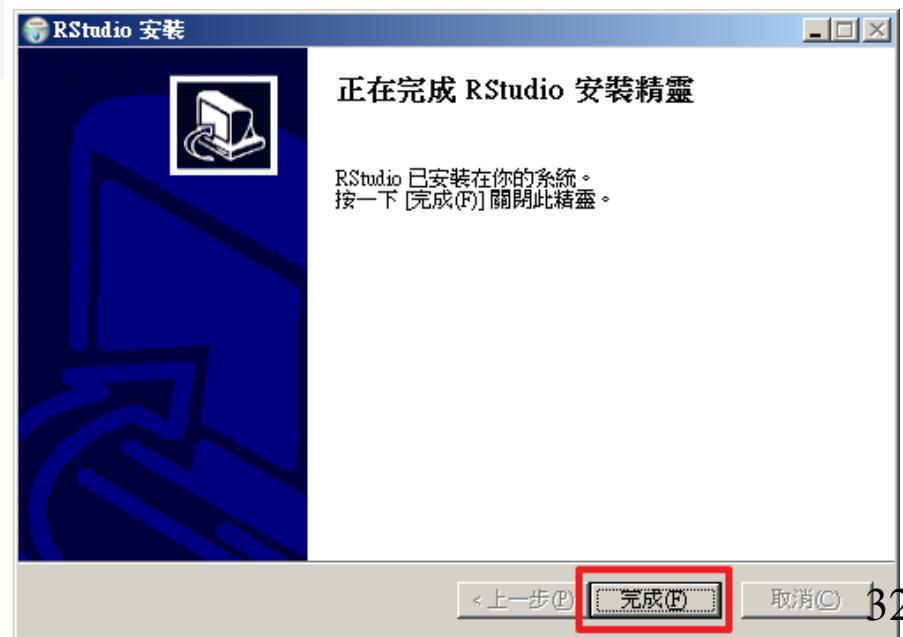
2



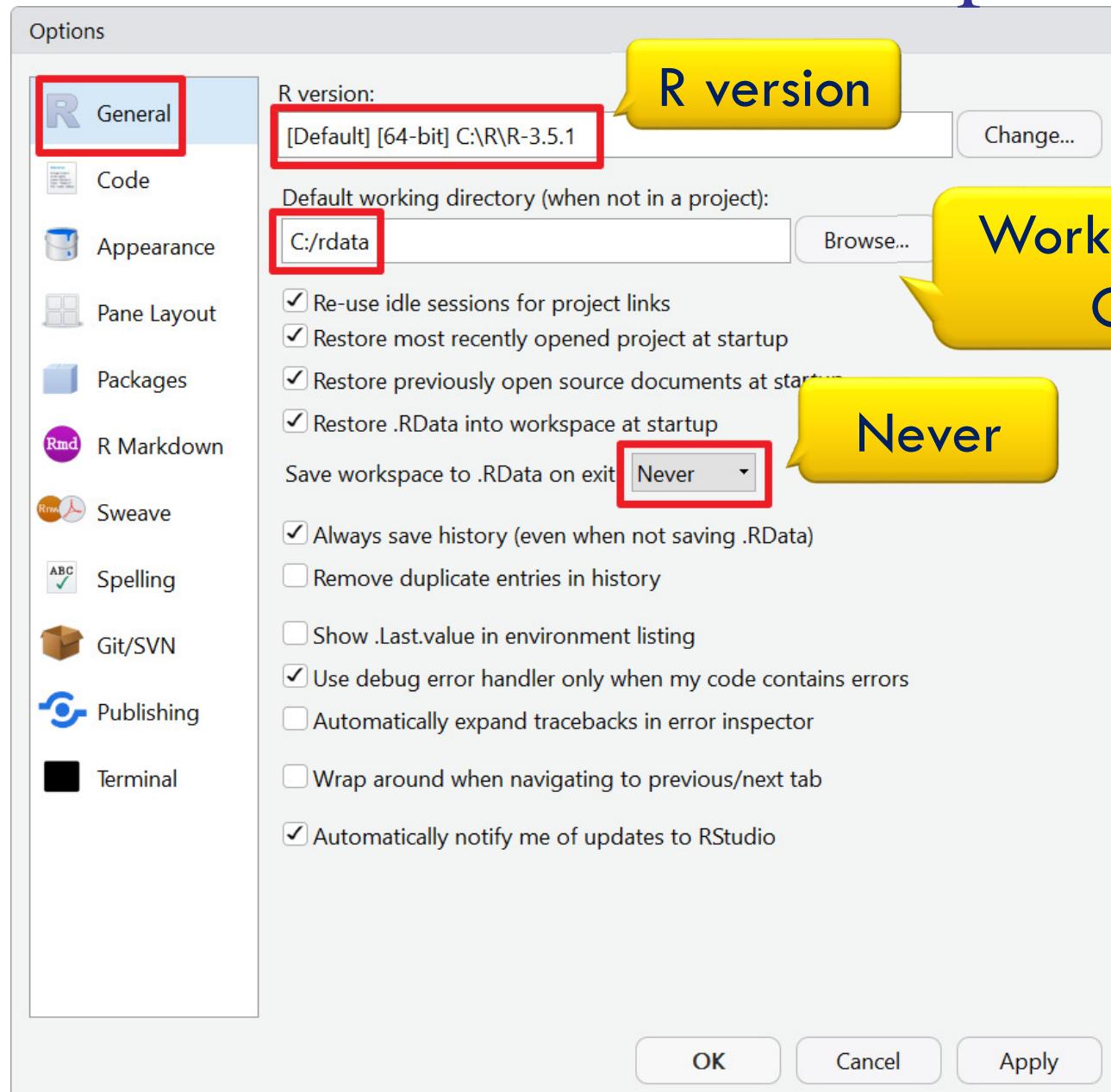
3



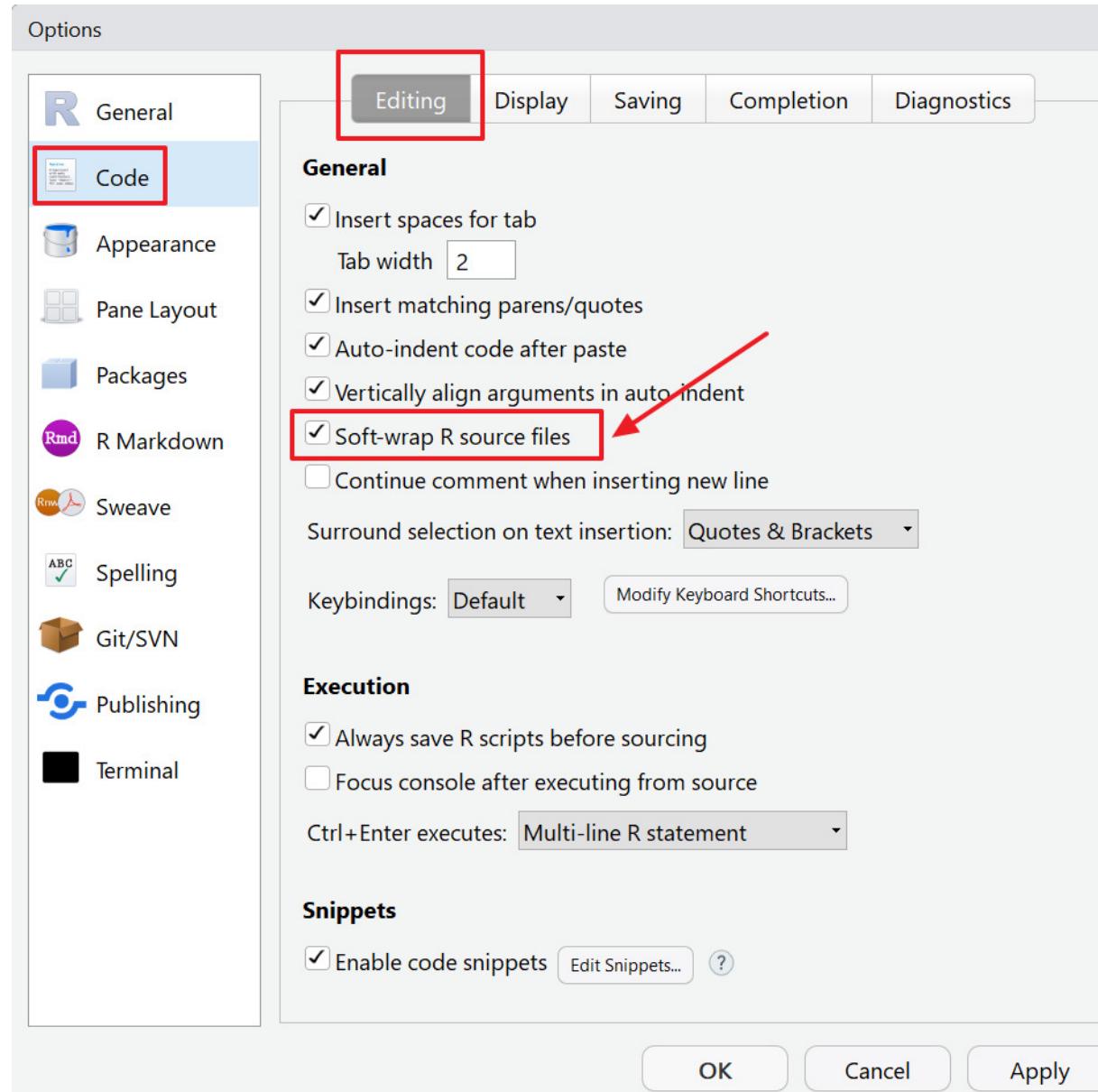
4



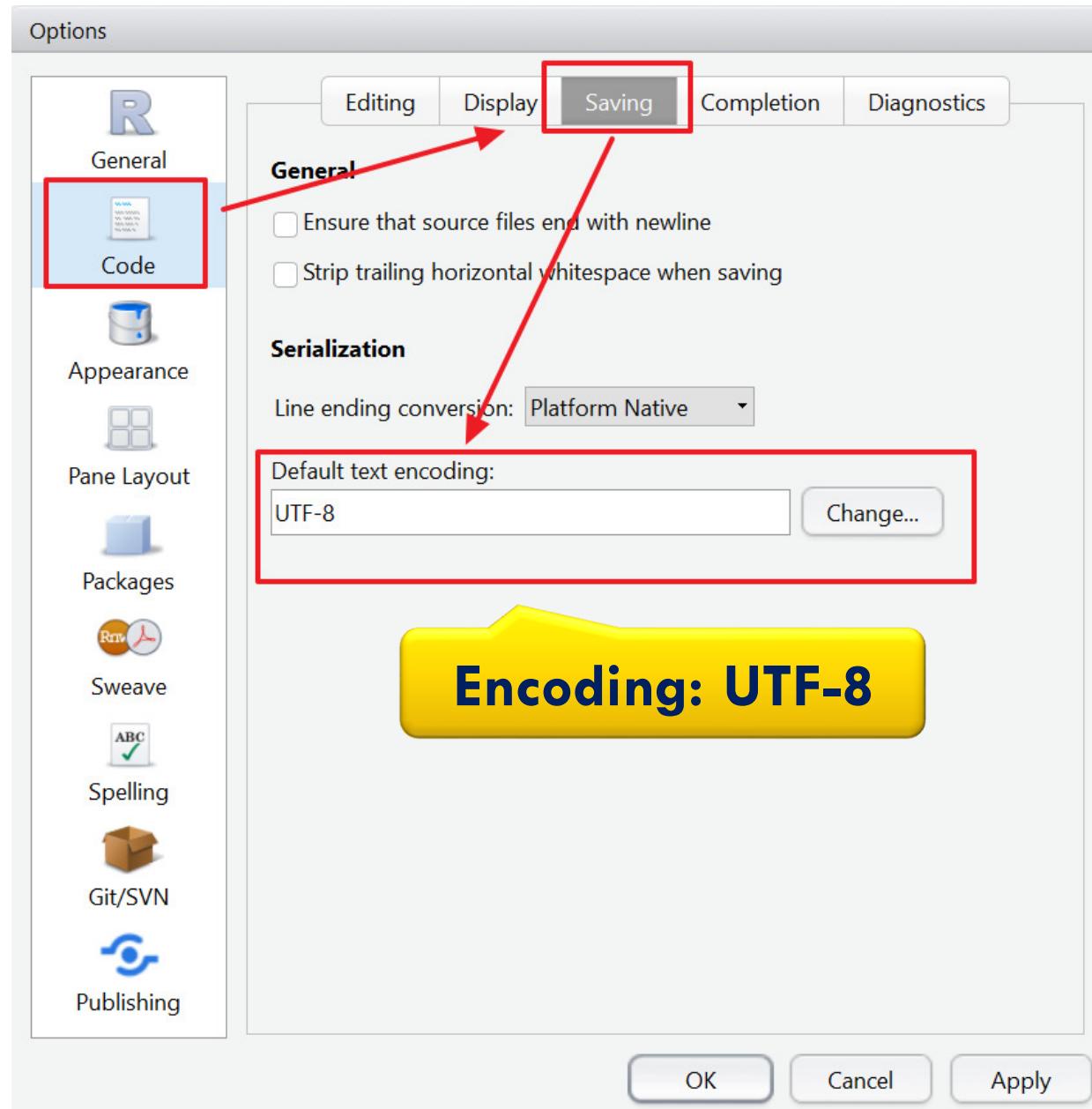
RStudio - Tools \ Global Options



RStudio - Options

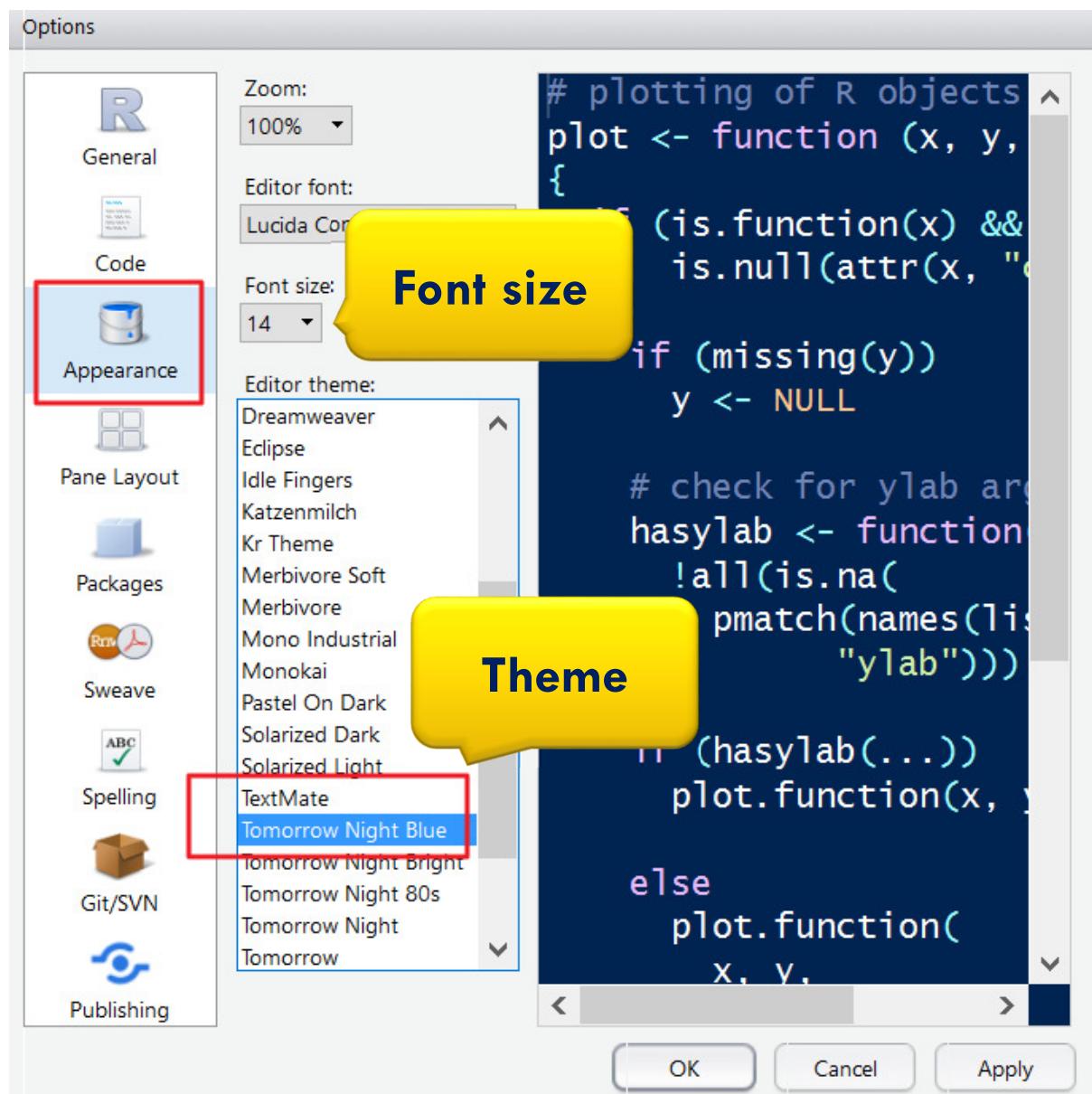


RStudio - Options

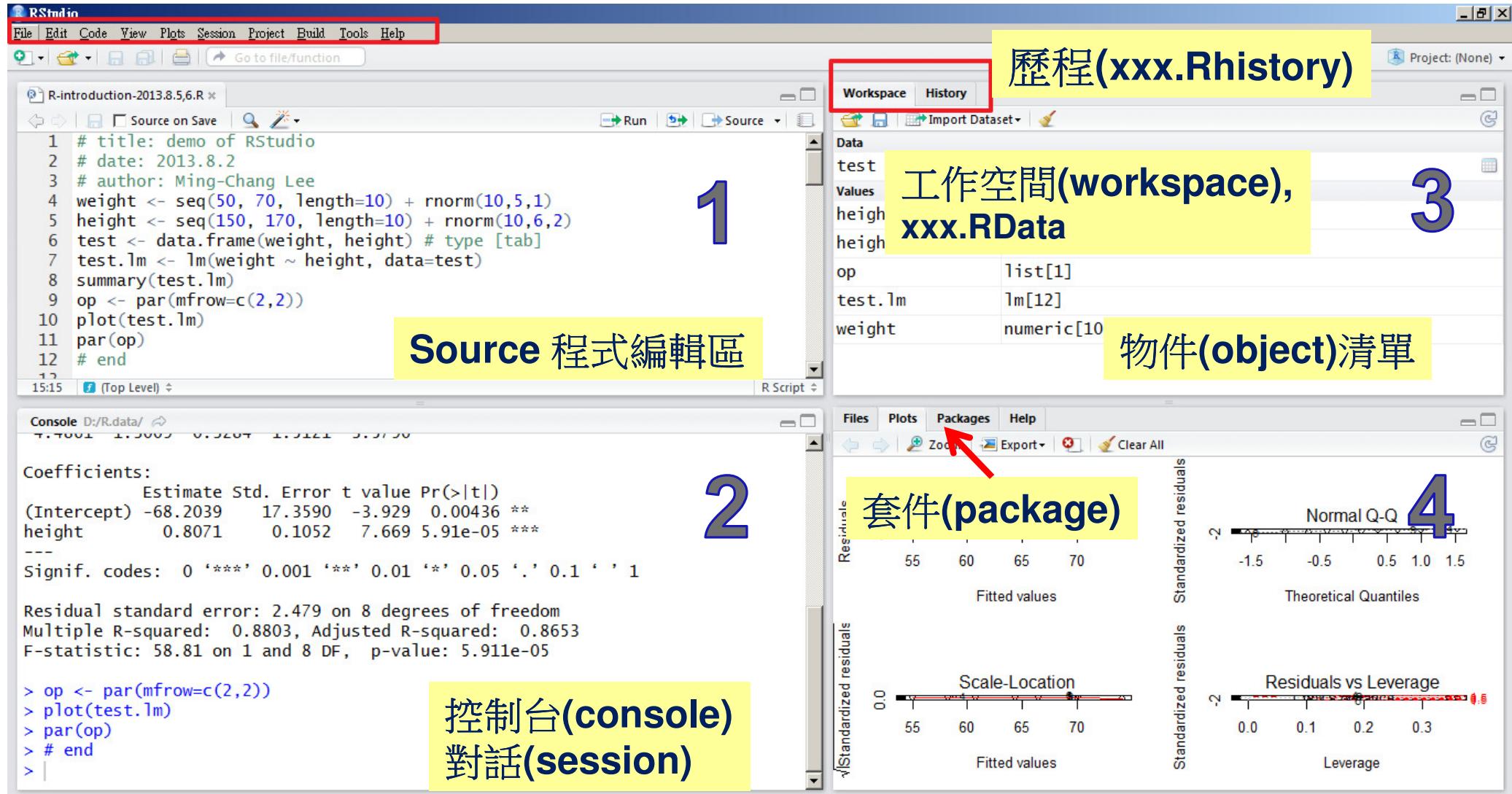


RStudio - Options

- Appearance \ Editor theme
- Default:
TextMate



RStudio - Environment



The screenshot illustrates the RStudio environment with four numbered sections:

- Source 程式編輯區**: The top-left pane shows an R script named "R-introduction-2013.8.5.R" containing code for a linear regression analysis.
- 控制台(console)
對話(session)**: The bottom-left pane displays the R console output, including the regression results and the command to plot the results.
- 歷程(xxx.Rhistory)**: The top-right pane shows the history of objects created in the workspace, with a yellow box highlighting the tab bar.
- 工作空間(workspace),
xxx.RData**: The middle-right pane displays the current objects in the workspace, with a yellow box highlighting the workspace tab.
- 物件(object)清單**: The bottom-right pane shows diagnostic plots for the regression model: Normal Q-Q plot, Scale-Location plot, and Residuals vs Leverage plot.

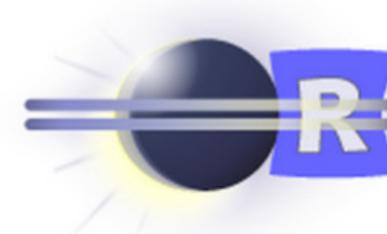
CTRL + SHIFT + F10: 重新啟動R

RStudio - Keyboard Shortcuts

Keyboard Shortcuts	Feature
Ctrl + Shift + N	Create new script
Ctrl + S	Save
Ctrl + Shift + R	Insert section
Alt + -	Insert assignment operator
Ctrl + Shift + C	Comments
Ctrl + Enter	Run selected lines
Ctrl + Shift + F10	Restart R
Alt + Shift + K	Shortcuts reference

R + Editor

- R – Naïve environment
- RStudio – IDE
- Eclipse
 - StatET 3.4.1:
An Eclipse based IDE (integrated development environment) plug-in for R.
 - <http://www.walware.de/goto/statet>
- R Tools for Visual Studio
 - <https://docs.microsoft.com/zh-tw/visualstudio/rts/installing-r-tools-for-visual-studio>



Packages 套件

Packages

- Shared R codes
- Components:
 - R codes
 - R functions
 - Datasets
 - Help document
- Basic packages: 30
- Installed directory: \R-3.5.1\library

Packages

- Step 1: `install.packages("PackageName")`
- Step 2: `library(PackageName)`
- Install and load “e1071” (machine learning)

```
> install.packages("e1071")
trying URL 'http://cran.cs.pu.edu.tw/bin/windows/contrib/3.0/e1071_1.6-1.zip'
Content type 'application/zip' length 514468 bytes (502 Kb)
opened URL
downloaded 502 Kb

package 'e1071' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:\Users\Administrator\AppData\Local\Temp\RtmpoHSOAk\downloaded_packages
> library(e1071)
Loading required package: class
>
```

loaded packages
search()

Packages - 38 views

Contributed Packages (2018.11.18)

Available Packages

Currently, the CRAN package repository features 13398 available packages.

[Table of available packages, sorted by date of publication](#)

[Table of available packages, sorted by name](#)

sorted by name

Installation of Packages

Please type `help("INSTALL")` or `help("install.packages")` in R for information on how to install packages from this repository. The manual [R Installation and Administration](#) (also contained in the R base sources) explains the process in detail.

[CRAN Task Views](#) allow you to browse packages by topic and provide tools to automatically install all packages for special areas of interest. Currently, 38 views are available.

38 views

Packages - 38 views

2013年10月8日 星期二

RWEPA → task

Task Views - R套件區分成38個類別

更新日期: 2018.11.18

CRAN Task View : <https://cran.r-project.org/web/views/>

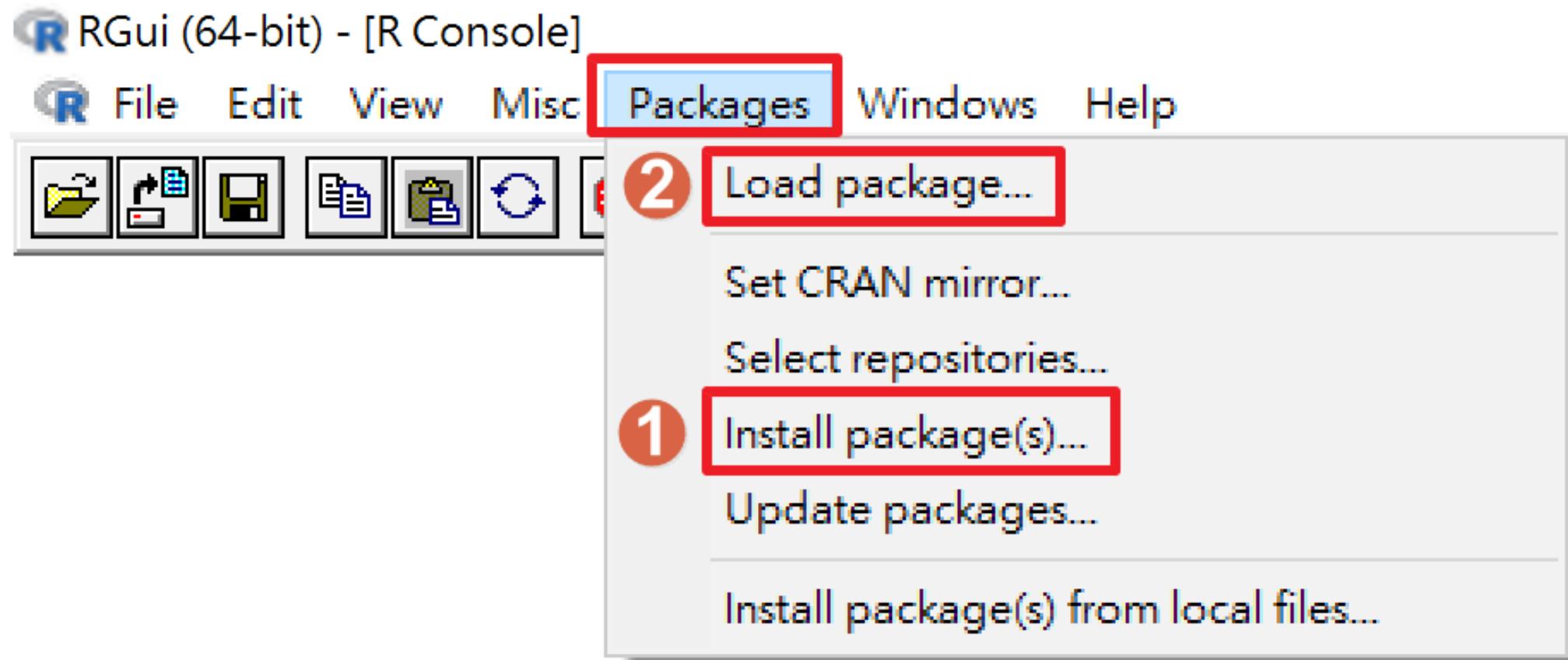
套件區分成 38 個類別, 中文說明如下:

-----	-----	-----	-----
編號	主題	英文說明	中文說明
-----	-----	-----	-----

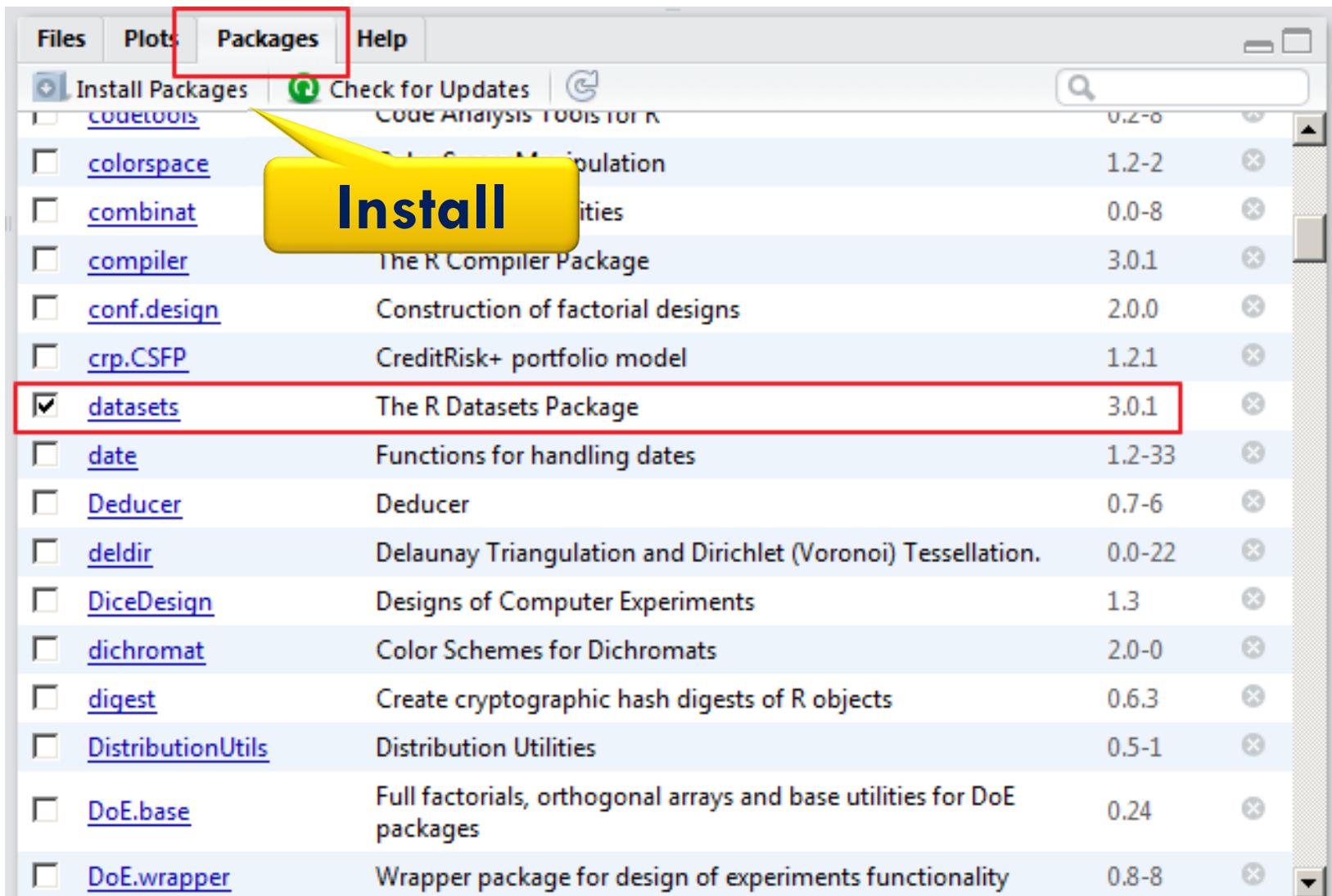
- 01, **Bayesian**, Bayesian Inference, 貝氏統計
- 02, **ChemPhys**, Chemometrics and Computational Physics, 計量化學學, 計算物理
- 03, **ClinicalTrials**, Clinical Trial Design, Monitoring, and Analysis, 臨床試驗設計, 監測和分析
- 04, **Cluster**, Cluster Analysis & Finite Mixture Models, 群集分析, 有限混合模型
- 05, **Databases**, Databases with R, R與資料庫連接
- 06, **DifferentialEquations**, Differential Equations, 微分方程
- 07, **Distributions**, Probability Distributions, 機率分配
- 08, **Econometrics**, Computational Econometrics, 計量經濟
- 09, **Environmetrics**, Analysis of Ecological and Environmental Data, 生態, 環境資料分析

```
# Installed list  
x <- installed.packages()  
x
```

R - Packages



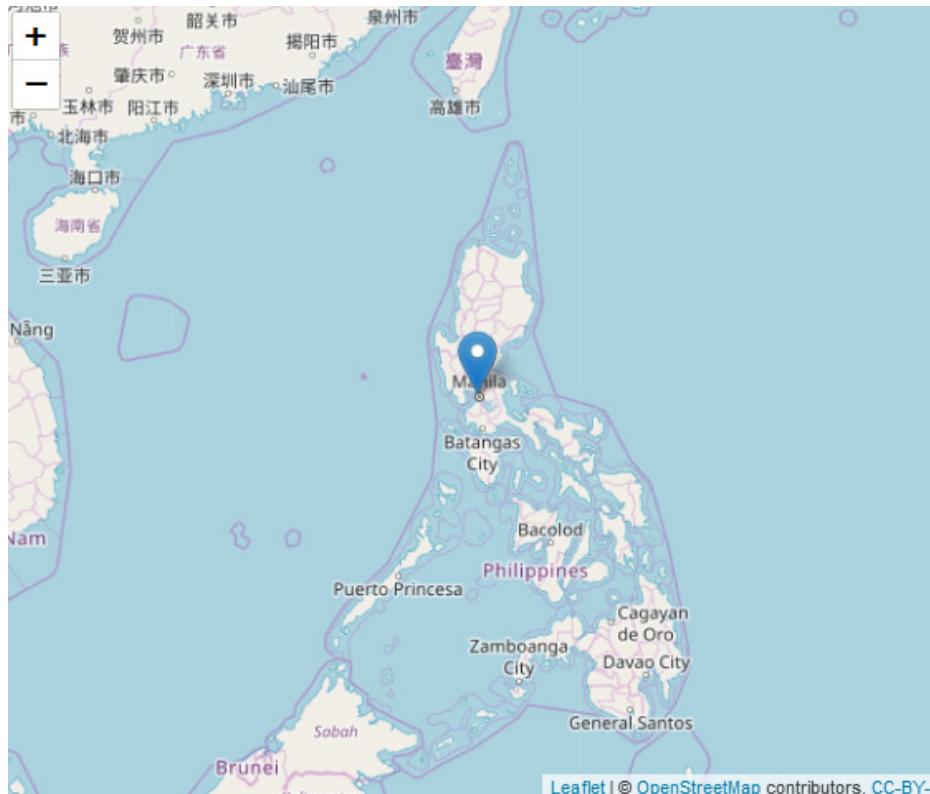
RStudio - Packages



leaflet package

<https://rstudio.github.io/leaflet/>

```
library(leaflet)
m <- leaflet() %>%
  addTiles() %>%
  addMarkers(lng=120.973551, lat=14.593105, popup="Home") %>%
  setView(lng = 120.973551, lat = 14.593105, zoom = 5)
m # Print the map
```



Data type

- Integer
- Numeric
- Character
- Logical

Date type - integer

```
> # 整數
> x1 <- c(1:10)
> x1
[1] 1 2 3 4 5 6 7 8 9 10
> typeof(x1)
[1] "integer"
> is.integer(x1)
[1] TRUE
> is.numeric(x1)
[1] TRUE
> class(x1)
[1] "integer"
```

Data type - numeric

```
> # 數值
> x2 <- c(1.1, 1.3, 1.5, 1.7, 2)
> x2
[1] 1.1 1.3 1.5 1.7 2.0
> typeof(x2)
[1] "double"
> is.integer(x2)
[1] FALSE
> is.numeric(x2)
[1] TRUE
> class(x2)
[1] "numeric"
```

Data type - character

```
> # 字串資料
> x3 <- c("台北市", "新北市", "台中市", "台南市", "高雄市")
> x3
[1] "台北市" "新北市" "台中市" "台南市" "高雄市"
> typeof(x3)
[1] "character"
> is.character(x3)
[1] TRUE
> class(x3)
[1] "character"
> x4 <- c(1, 2.3, "巨量資料")
> x4
[1] "1"          "2.3"        "巨量資料"
> class(x4)
[1] "character"
```

force transfer to
character

Data type - logical

- 邏輯值包括: TRUE 或 FALSE
- TRUE → 1
- FALSE → 0

```
> # 邏輯值
> TRUE*2
[1] 2
> FALSE*3
[1] 0
```

Data type - logical

```
> x5 <- x2 >= 1.4
> x5
[1] FALSE FALSE TRUE TRUE TRUE
> typeof(x5)
[1] "logical"
> is.integer(x5)
[1] FALSE
> is.numeric(x5)
[1] FALSE
> is.character(x5)
[1] FALSE
> is.logical(x5)
[1] TRUE
> class(x5)
[1] "logical"
```

Data object

Data object

向量 vector

北部	中部	南部
----	----	----

矩阵 matrix

1	3	5
2	4	6

阵列 array

1.1	4.4	7.7
2.2	5.5	8.8
3.3	6.6	9.9

资料框 data.frame

1	男	62
2	女	50
3	女	54
4	男	72

串列 list

北部	中部	南部
1	3	5
2	4	6

1	男	62
2	女	50
3	女	54
4	男	72

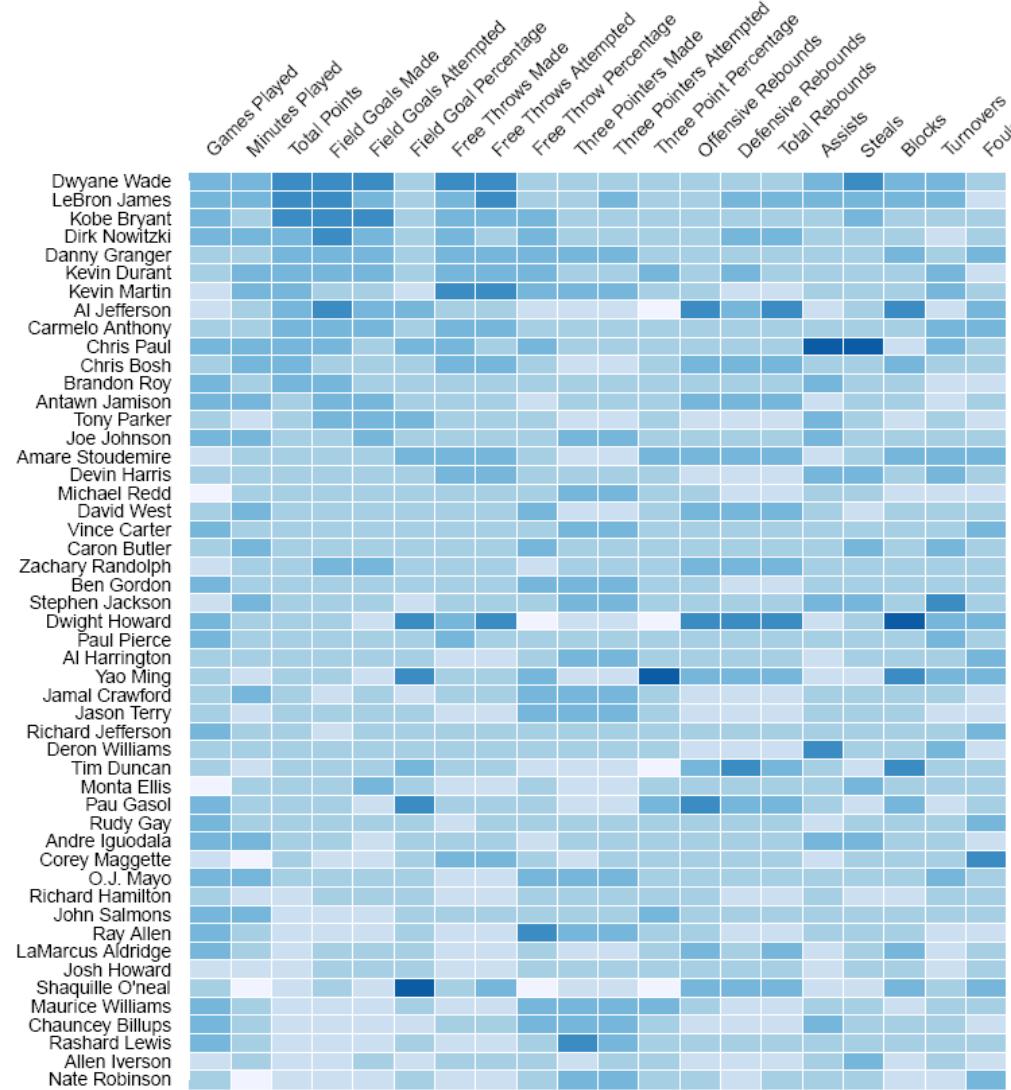
vector

matrix

data.frame

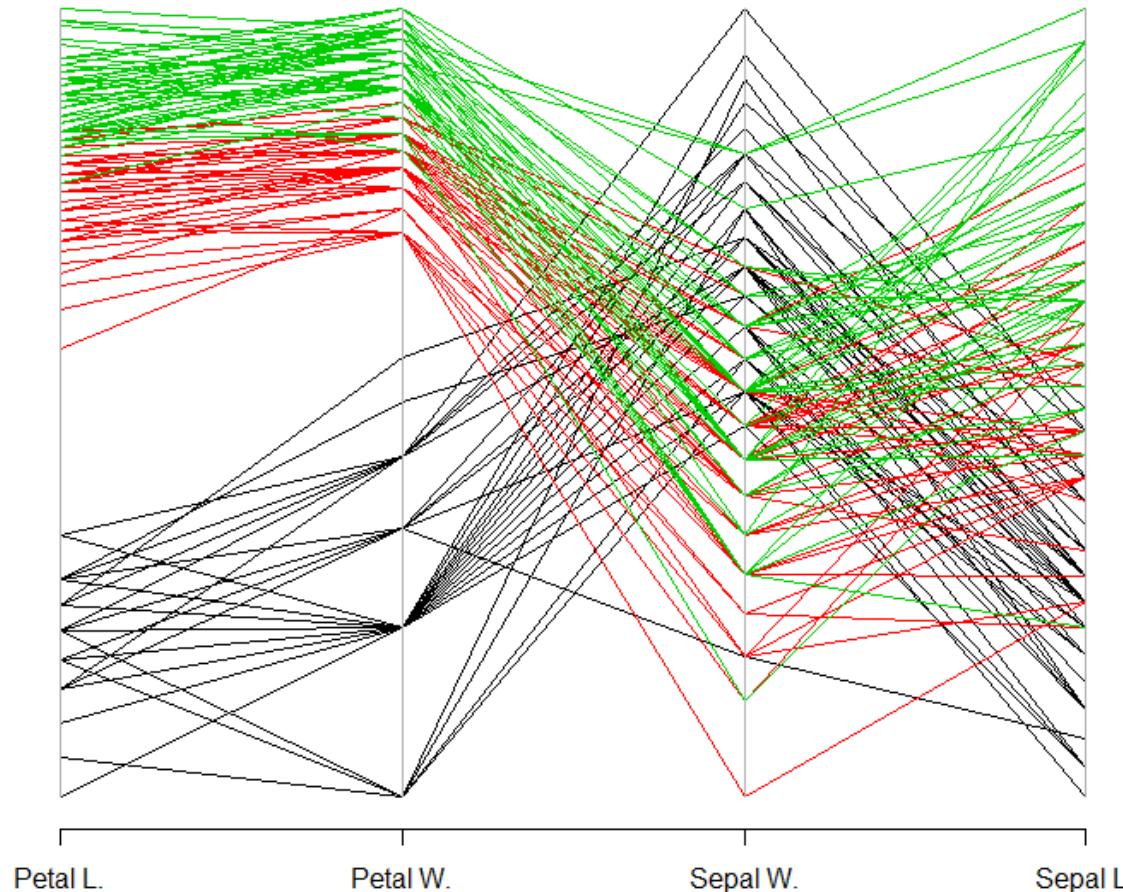
image {graphics}

NBA per game performance of top 50 scorers



parcoord {MASS}

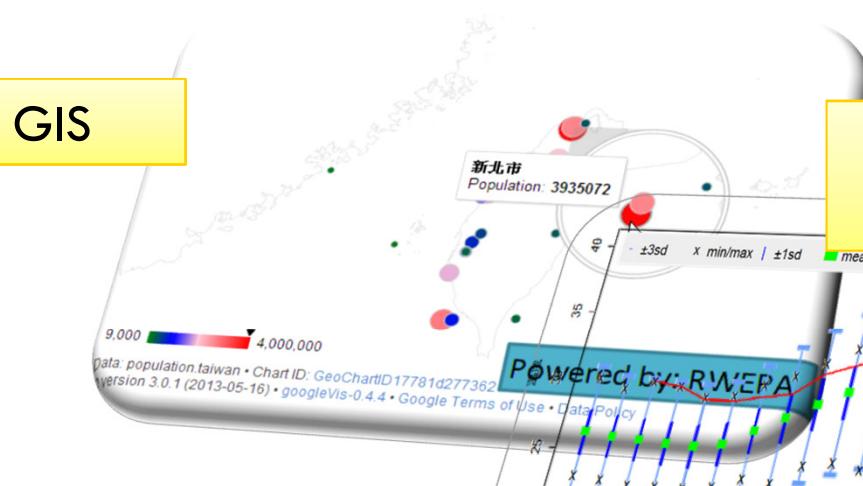
Parallel coordinates



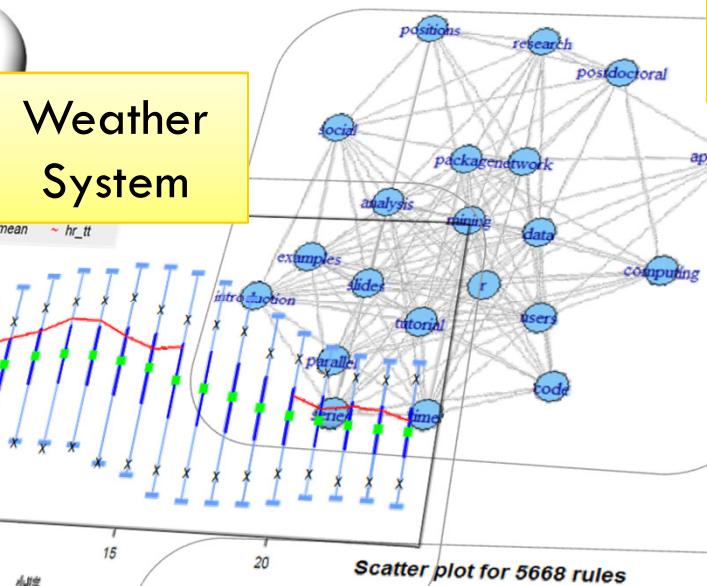
3. Big Data Application

Big Data Application

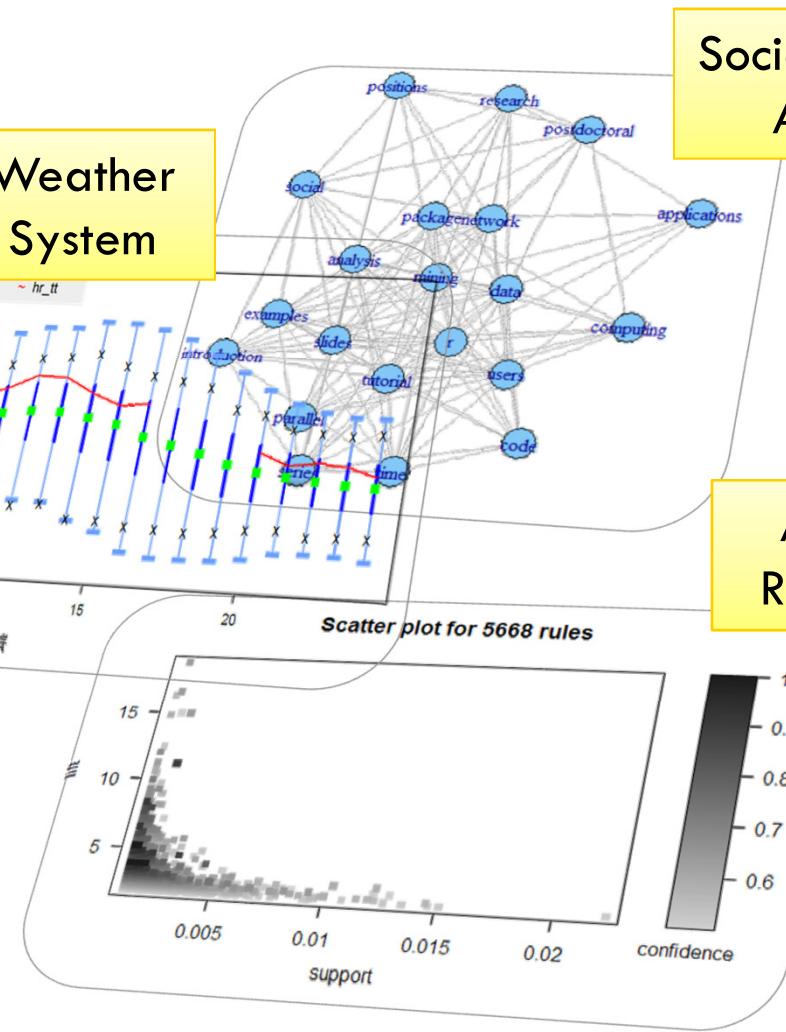
GIS



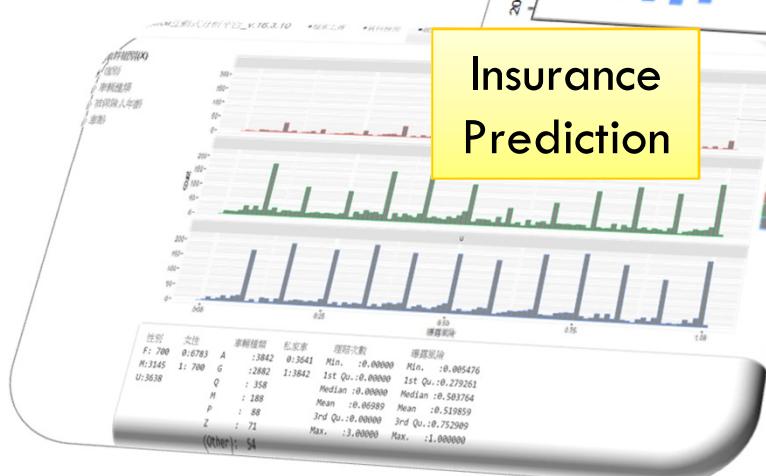
Weather System



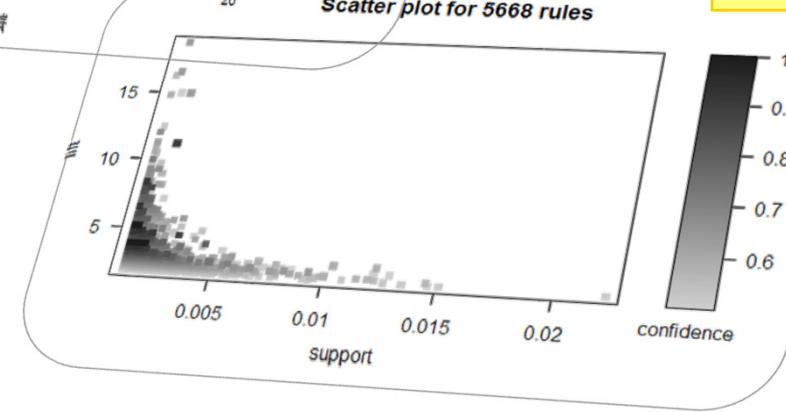
Social Network Analysis



Insurance Prediction



Association Rule Analysis



Weather System (data size=16 Million)



Insurance Forecasting

• Threshold

機率模型閾值 1

• Probability Result

預測資料上傳
檢視結果

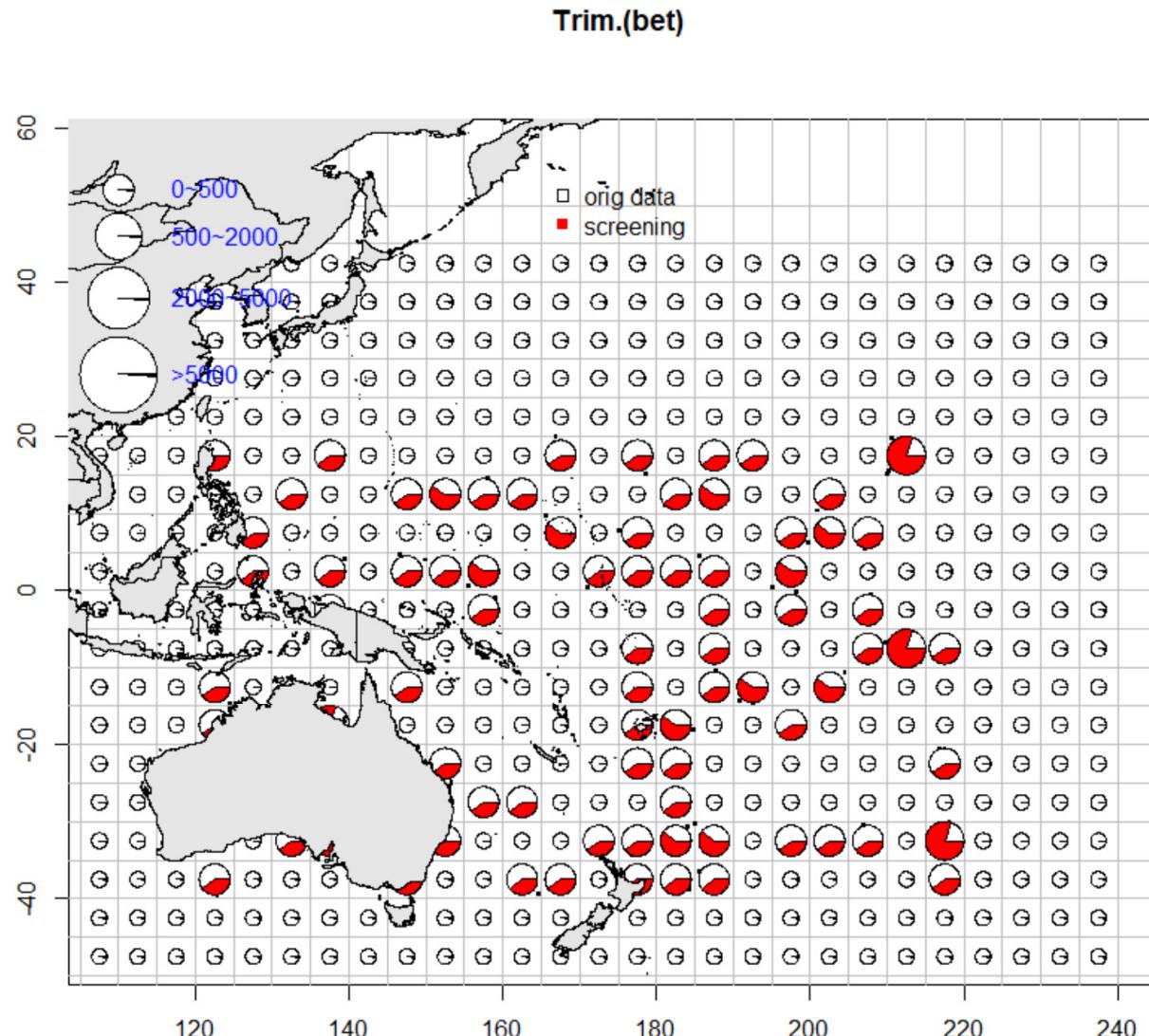
編號	性別	女性	車輛種類	私家車	曝露風險	曝露風險對數	無索償折扣	被保險人年齡	私家車-車齡0	私家車-車齡1	私家車-車齡2	私家車-車齡0_1_2組合	車齡	車齡0_1_2組合	預測機率	理賠
1	M	0	A	1	0.9144422	-0.08944106	50	4	1	0	0	1	0	2	0.1069	有
2	M	0	A	1	0.8158795	-0.20348856	20	4	0	0	1	1	2	2	0.1441	有
3	M	0	A	1	0.8377823	-0.17699695	50	3	0	0	1	1	2	2	0.1866	有
4	M	0	A	1	0.4325804	-0.83798702	50	6	0	1	0	1	1	2	0.0944	無
5	M	0	A	1	0.7173169	-0.33223755	50	4	0	0	1	1	2	2	0.1218	有
6	M	0	A	1	0.8377823	-0.17699695	50	4	0	0	1	1	2	2	0.1495	有
7	M	0	A	1	0.8487337	-0.16400975	50	5	0	0	1	1	2	2	0.1422	有
8	F	1	A	1	0.8268309	-0.19015503	10	3	0	0	1	1	2	2	0.1733	有
9	M	0	A	1	0.7145791	-0.33606164	0	5	1	0	0	1	0	2	0.0694	無
10	M	0	A	1	0.3340178	-1.09656101	0	3	0	0	1	1	2	2	0.0783	無

Showing 1 to 10 of 12 entries

Previous 1 2 Next

127.0.0.1:6177/#tab-9487-2

Space Pie Chart Detector



Open Data

- OPEN DATA TAIWAN
 - <https://data.gov.tw/>
- OPEN DATA PHILIPPINES
 - <https://data.gov.ph/>
- UCI Machine Learning Repository
 - <https://archive.ics.uci.edu/ml/datasets.html>

References

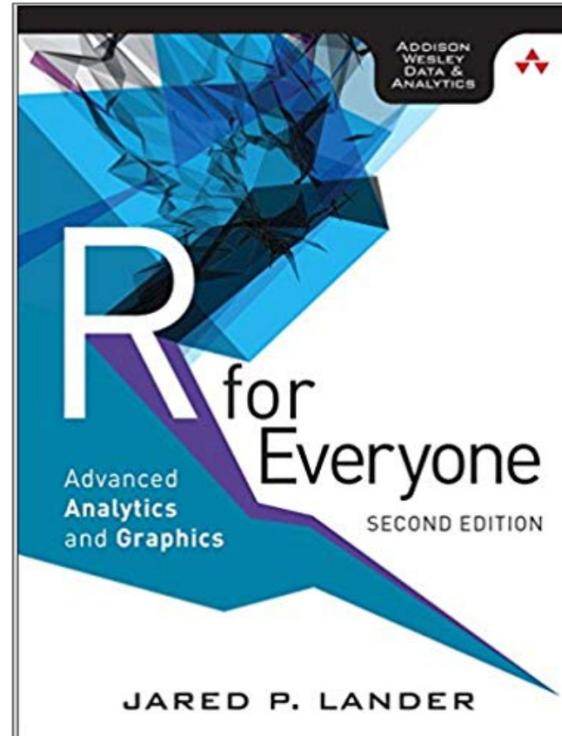
Using R for Data Analysis and Graphics

Introduction, Code and Commentary

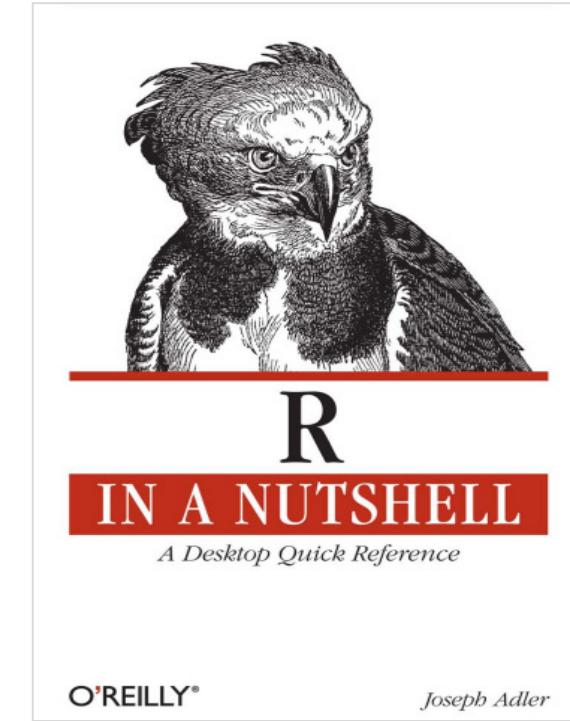
J H Maindonald

Centre for Mathematics and Its Applications,
Australian National University.

(Basic)



(Intermediate)



(Advanced)

Thank you

alan9956@gmail.com

<http://rwepa.blogspot.tw/>