# 中文文字探勘

國立臺灣大學共同教育中心

蔡芸玎

# 套件安裝

library(tmcn)：https://r-forge.r-project.org/R/?group_id=1571

library(NLP)

library(tm)

library(jiebaRD)

library(jiebaR)

library(RColorBrewer)
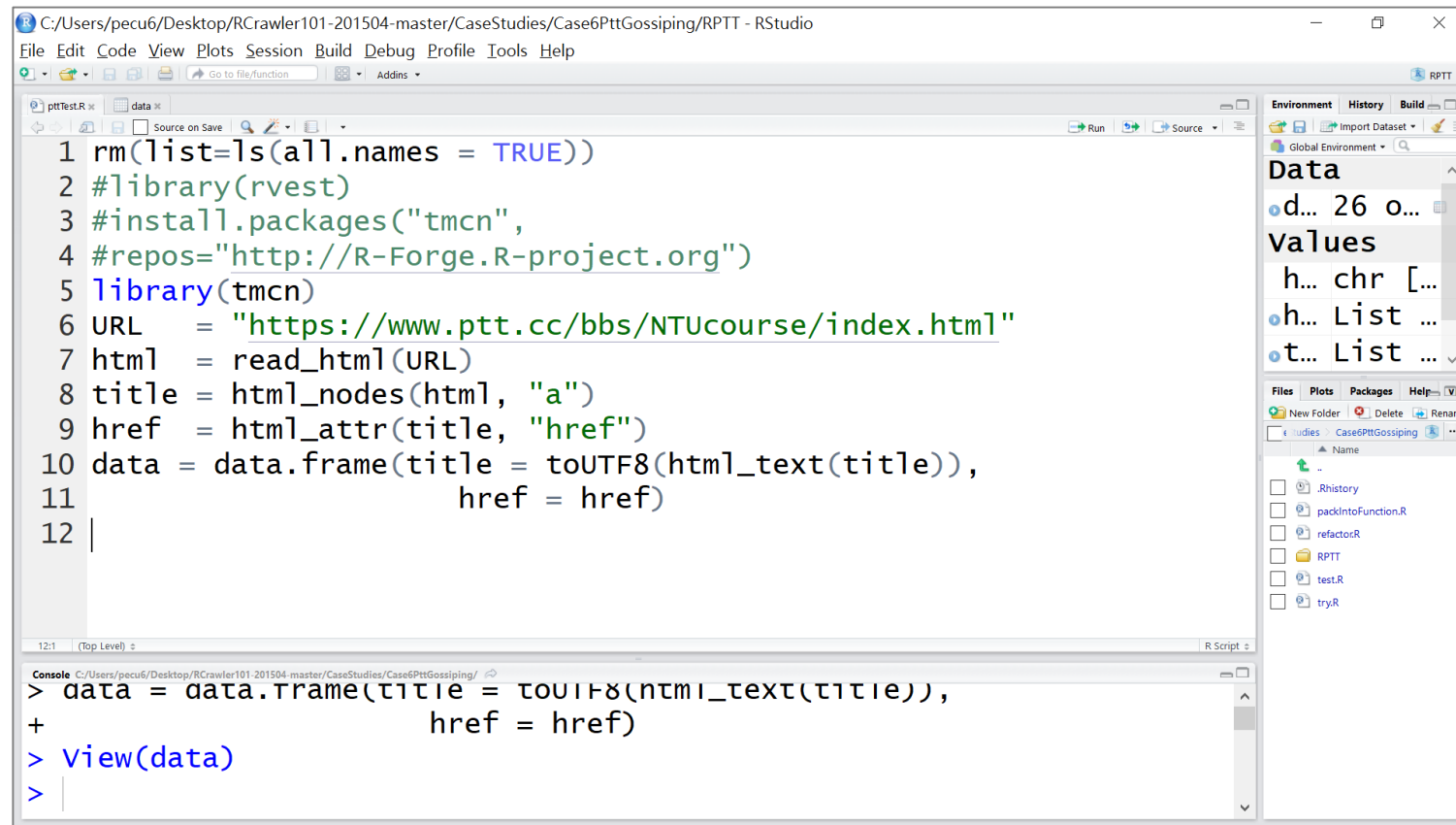
library(wordcloud)

library(rvest)

# 文本蒐集

1. 網路爬蟲，範例程式：
   https://github.com/pecu/RCrawler101-201504
2. 本機端純文字檔

# 網路爬蟲文本蒐集

# 網路爬蟲文本蒐集

# 網路爬蟲文本蒐集

# 網路爬蟲文本蒐集

# 網路爬蟲文本蒐集

# 網路爬蟲文本蒐集



```r
 3 #install.packages("tmcn",
 4 #repos="http://R-Forge.R-project.org")
 5 library(tmcn)
 6 URL    = "https://www.ptt.cc/bbs/NTUcourse/index.html"
 7 html   = read_html(URL)
 8 title = html_nodes(html, "a")
 9 href   = html_attr(title, "href")
10 data = data.frame(title = toUTF8(html_text(title)),
11                           href = href)
12 data = data[-c(1:10),]
13 getContent <- function(x) {
14   url   = paste0("https://www.ptt.cc", x)
15   tag   = html_node(read_html(url), 'div#main-content.bbs-screen.bb
16   text = toUTF8(html_text(tag))
17 }
18 #getContent(data$href[1])
19 allText = sapply(data$href, getContent)
20 allText
```

搭配 sapply

# 網路爬蟲文本蒐集

C:/Users/pecu6/Desktop/RCrawler101-201504-master/CaseStudies/Case6PttGossiping/RPTT - RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Source

Console C:/Users/pecu6/Desktop/RCrawler101-201504-master/CaseStudies/Case6PttGossiping/

19:08\... <truncated>
[15] "作者benck (小倫)看板NTUcourse標題[整理] 常見Q&A 發文前必看時間Tue Feb 23 14:20:02 2010\n初選\n\nQ：二階初選可以選衝堂的課嗎？ (2011/1/25 update)\nA：不可以。\n\n\n開學後加退選\n\nQ：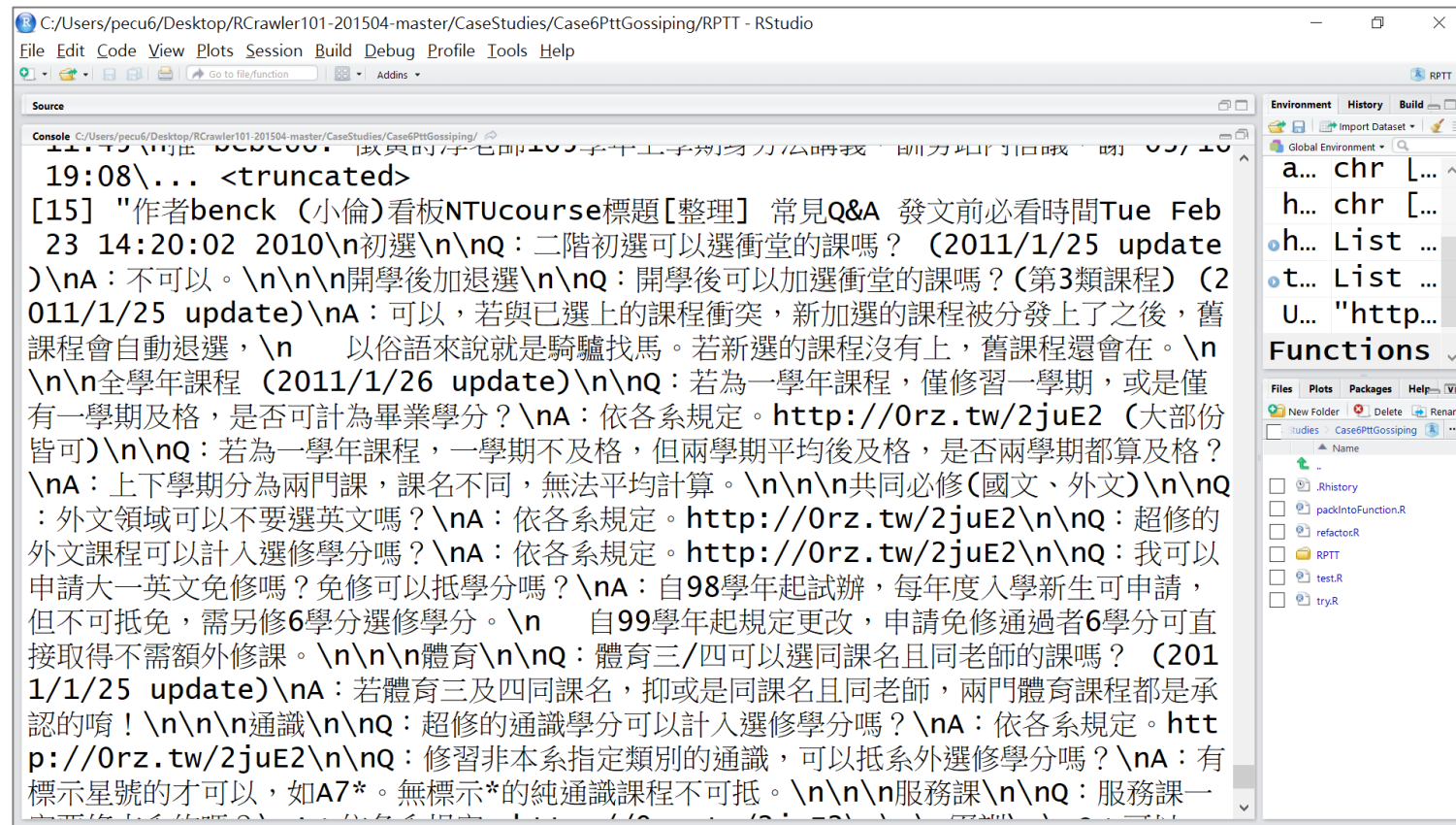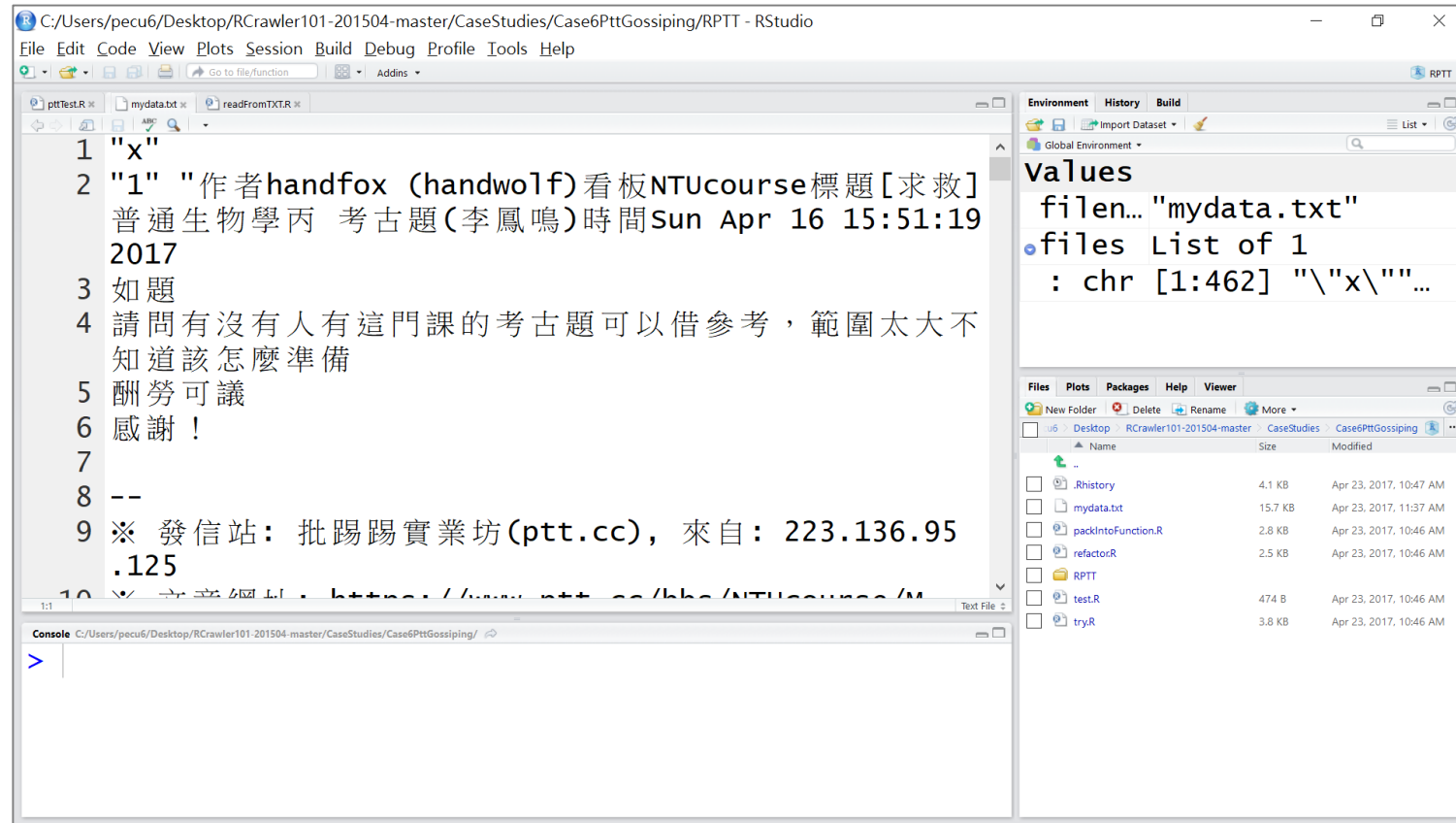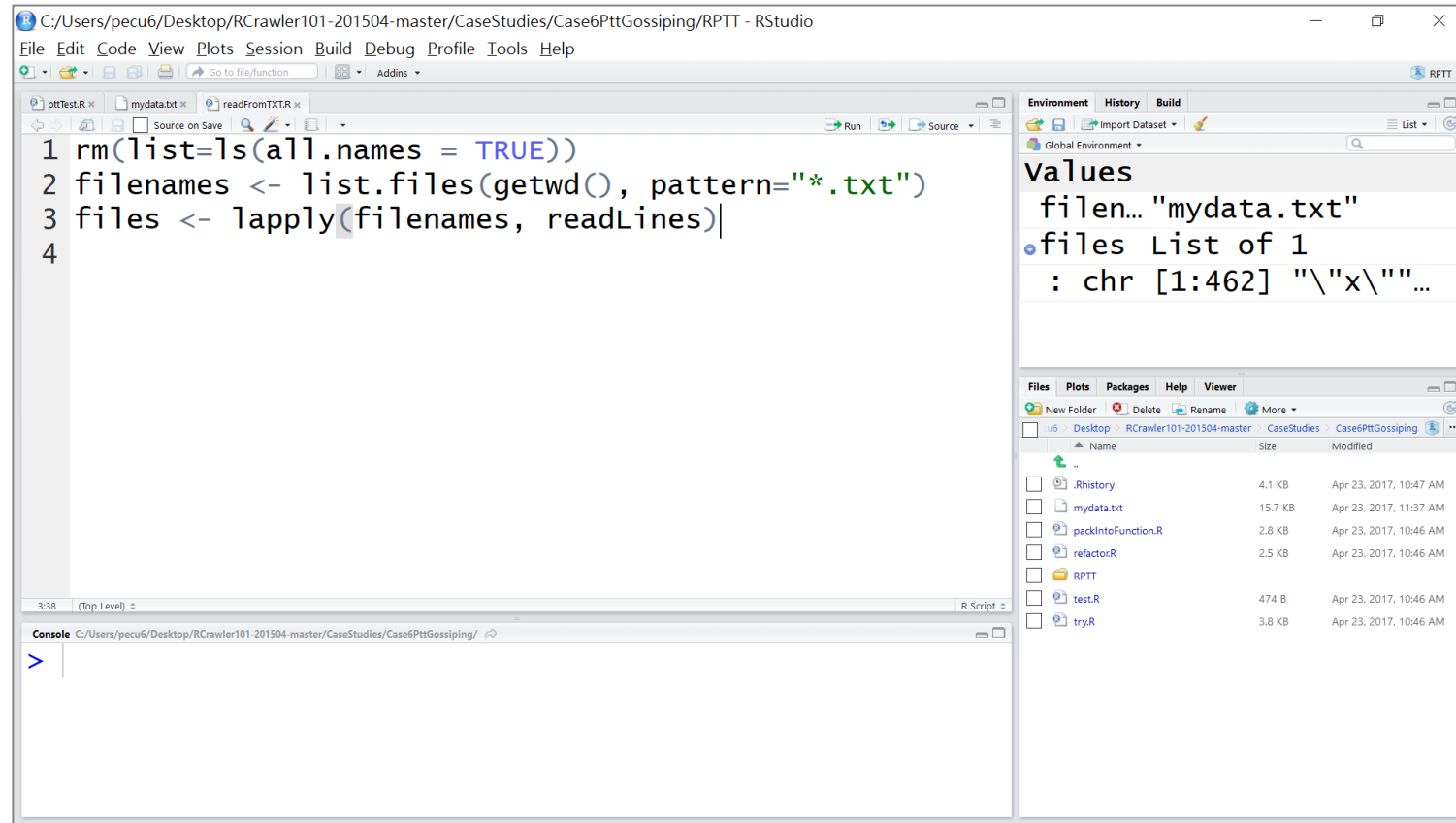開學後可以加選衝堂的課嗎？(第3類課程) (2011/1/25 update)\nA：可以，若與已選上的課程衝突，新加選的課程被分發上了之後，舊課程會自動退選，\n 以俗語來說就是騎驢找馬。若新選的課程沒有上，舊課程還會在。\n\n\n全學年課程 (2011/1/26 update)\n\nQ：若為一學年課程，僅修習一學期，或是僅有一學期及格，是否可計為畢業學分？\nA：依各系規定。http://0rz.tw/2juE2 （大部份皆可)\n\nQ：若為一學年課程，一學期不及格，但兩學期平均後及格，是否兩學期都算及格？\nA：上下學期分為兩門課，課名不同，無法平均計算。\n\n\n共同必修(國文、外文)\n\nQ：外文領域可以不要選英文嗎？\nA：依各系規定。http://0rz.tw/2juE2\n\nQ：超修的外文課程可以計入選修學分嗎？\nA：依各系規定。http://0rz.tw/2juE2\n\nQ：我可以申請大一英文免修嗎？免修可以抵學分嗎？\nA：自98學年起試辦，每年度入學新生可申請，但不可抵免，需另修6學分選修學分。\n 自99學年起規定更改，申請免修通過者6學分可直接取得不需額外修課。\n\n\n體育\n\nQ：體育三/四可以選同課名且同老師的課嗎？ (2011/1/25 update)\nA：若體育三及四同課名，抑或是同課名且同老師，兩門體育課程都是承認的唷！\n\n\n通識\n\nQ：超修的通識學分可以計入選修學分嗎？\nA：依各系規定。http://0rz.tw/2juE2\n\nQ：修習非本系指定類別的通識，可以抵系外選修學分嗎？\nA：有標示星號的才可以，如A7*。無標示*的純通識課程不可抵。\n\n\n服務課\n\nQ：服務課一

Environment History Build
Import Dataset
Global Environment

a... chr [...
h... chr [...
h... List ...
t... List ...
U... "http...

Functions

Files Plots Packages Help Vie
New Folder Delete Rename
Studies Case6PttGossiping
Name
..
.Rhistory
packIntoFunction.R
refactor.R
RPTT
test.R
try.R

# 本機端純文字檔本蒐集

# 本機端純文字檔本蒐集

# 課堂練習

連續產生 10 頁以上的純文字檔
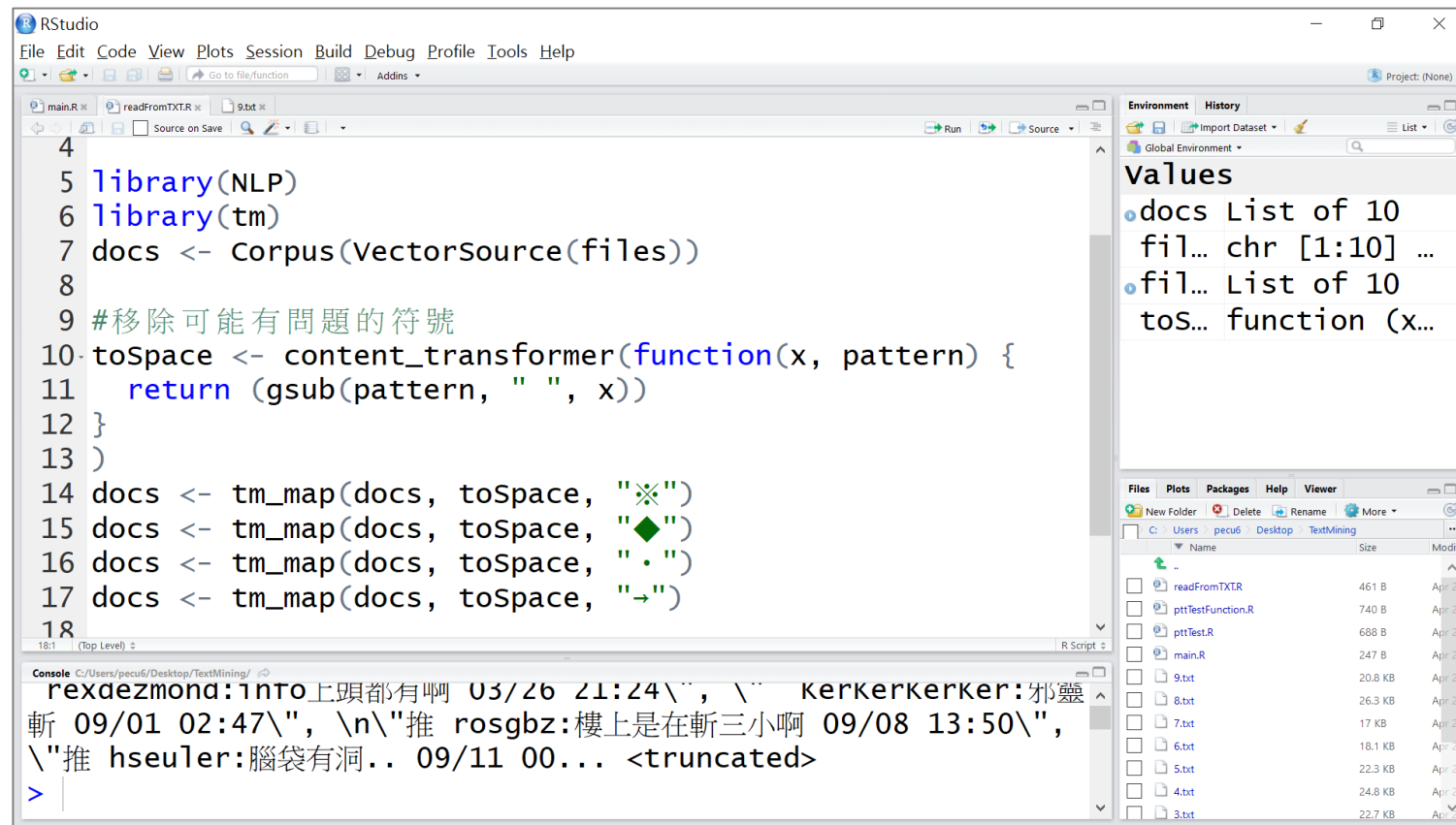
1. 人工複製貼上？
2. 網路爬蟲？
3. 現成的資料庫？

# 課堂練習



自訂 Function
搭配 mapply

自動生成檔案

# 文本清理

- 參考資料：https://cran.r-project.org/doc/contrib/de_Jonge+van_der_Loo-Introduction_to_data_cleaning_with_R.pdf
- 大小寫轉換
- 標點符號、數字移除
- URLs 移除
- 表情符號、停用詞移除

# 文本清理

# 斷詞處理

產生正確中文詞頻矩陣
1. 使用 cutter=worker() 產生切詞器。
2. 使用 new_user_word 將新詞彙加入詞庫。
3. 使用 cutter=worker("tag") 可切割出詞彙與提供詞彙的詞性。

# 詞頻矩陣

# 詞頻矩陣

# 文字雲

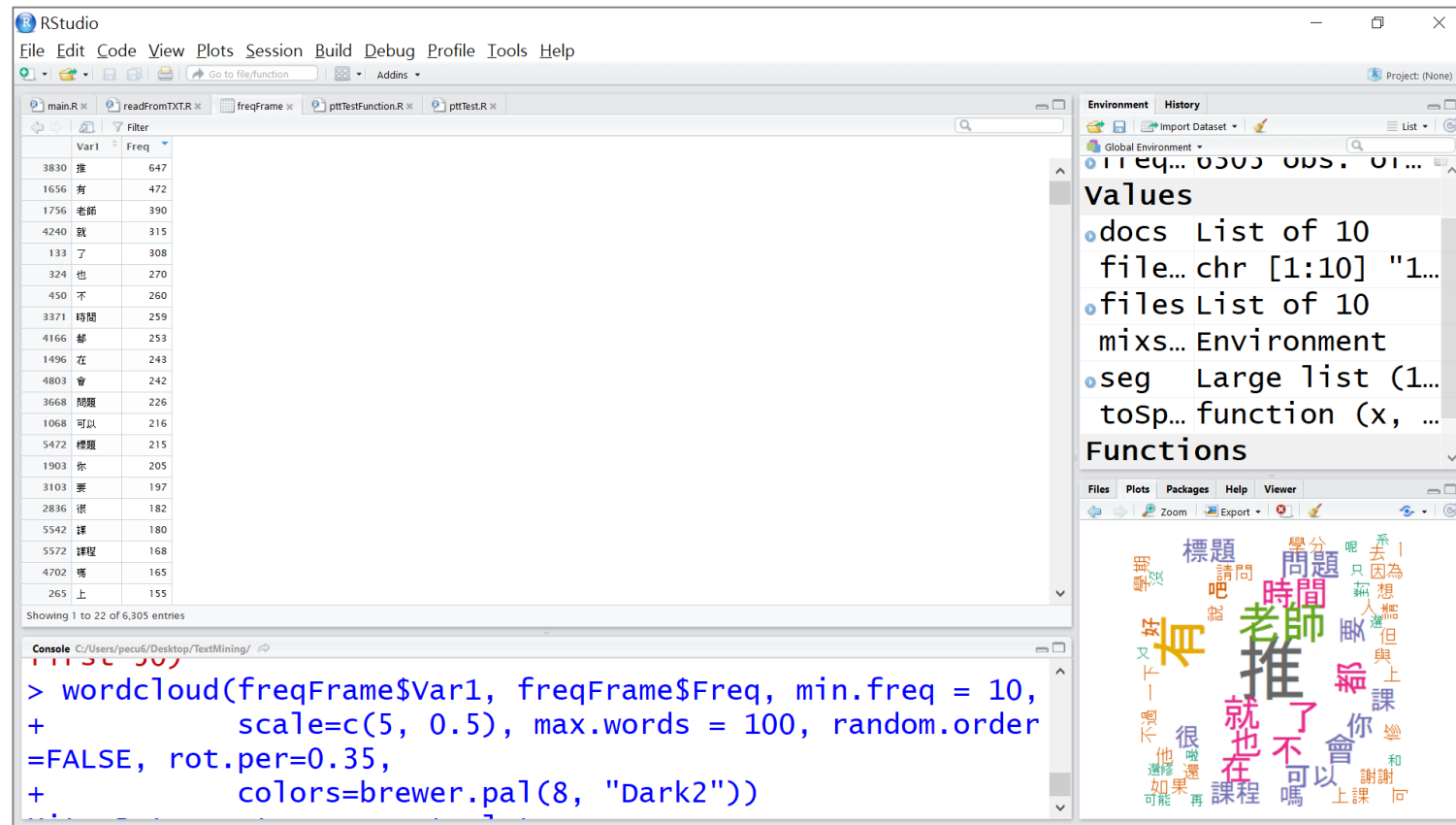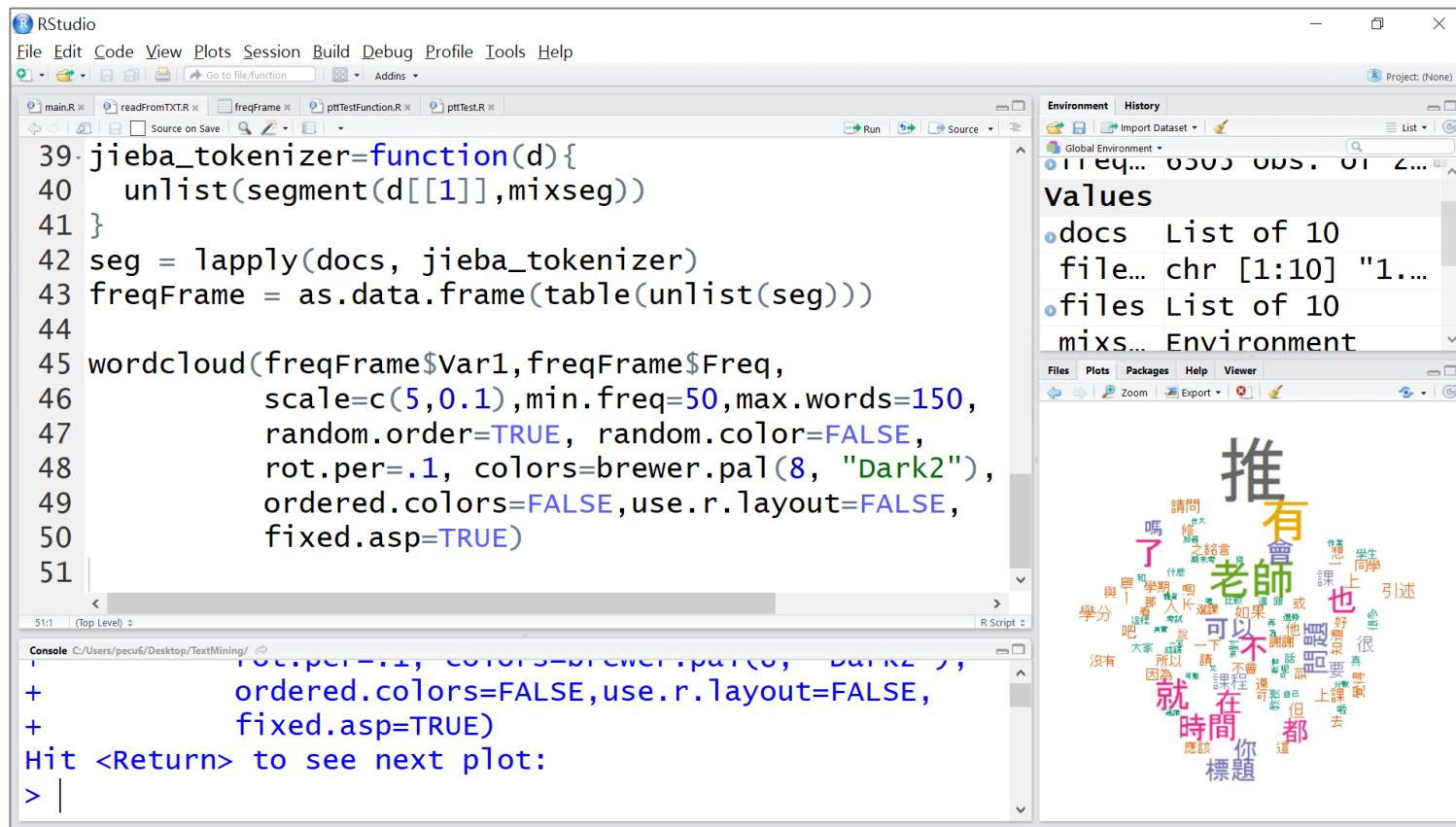https://cran.r-project.org/web/packages/wordcloud/wordcloud.pdf

wordcloud(words,freq,scale=c(4,.5),min.freq=3,max.words=Inf,

      random.order=TRUE, random.color=FALSE, rot.per=.1,

      colors="black",ordered.colors=FALSE,use.r.layout=FALSE,

      fixed.asp=TRUE, ...)

# 文字雲

# 文字探勘目標

https://buzzorange.com/techorange/2017/04/13/data-to-pxmart-hsu/

1. 文字雲
2. 詞語詞間的關係
3. 文本關聯