

上課前，同學可以先.....

- 安裝套件
- `install.packages(c("tm" , "gutenbergr"))`
- `devtools::install_github("hrbrmstr/hottopic")`
- 想想期末報告分組與主題

R 語言 文字資料 前處理

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

課程主題重點

- 語料庫
- 常用的文字資料前處理 (preprocessing) 套件介紹
tm, stringr [quanteda]
- 正規表示法 (regular expression)

文字資料分析流程

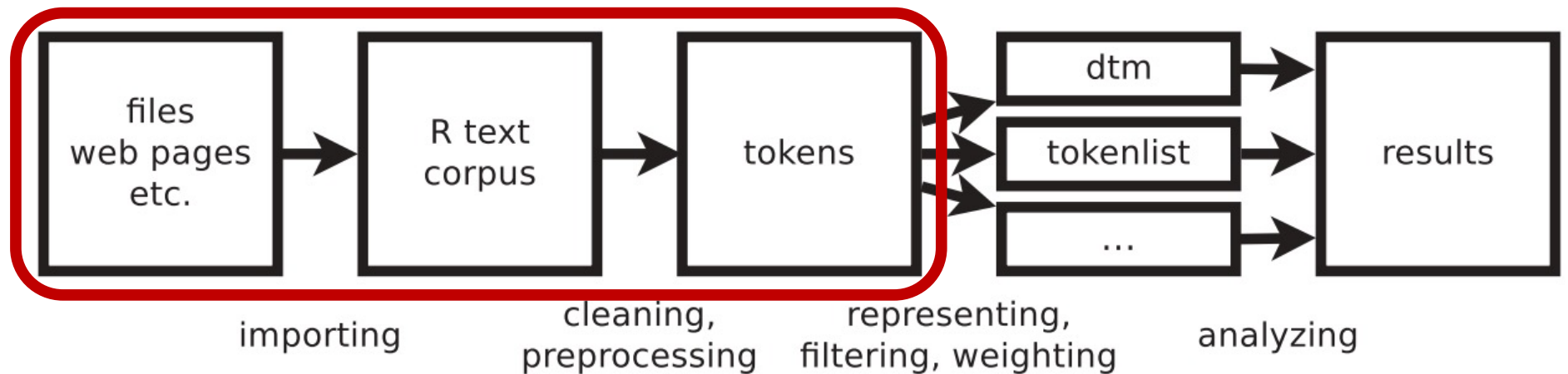
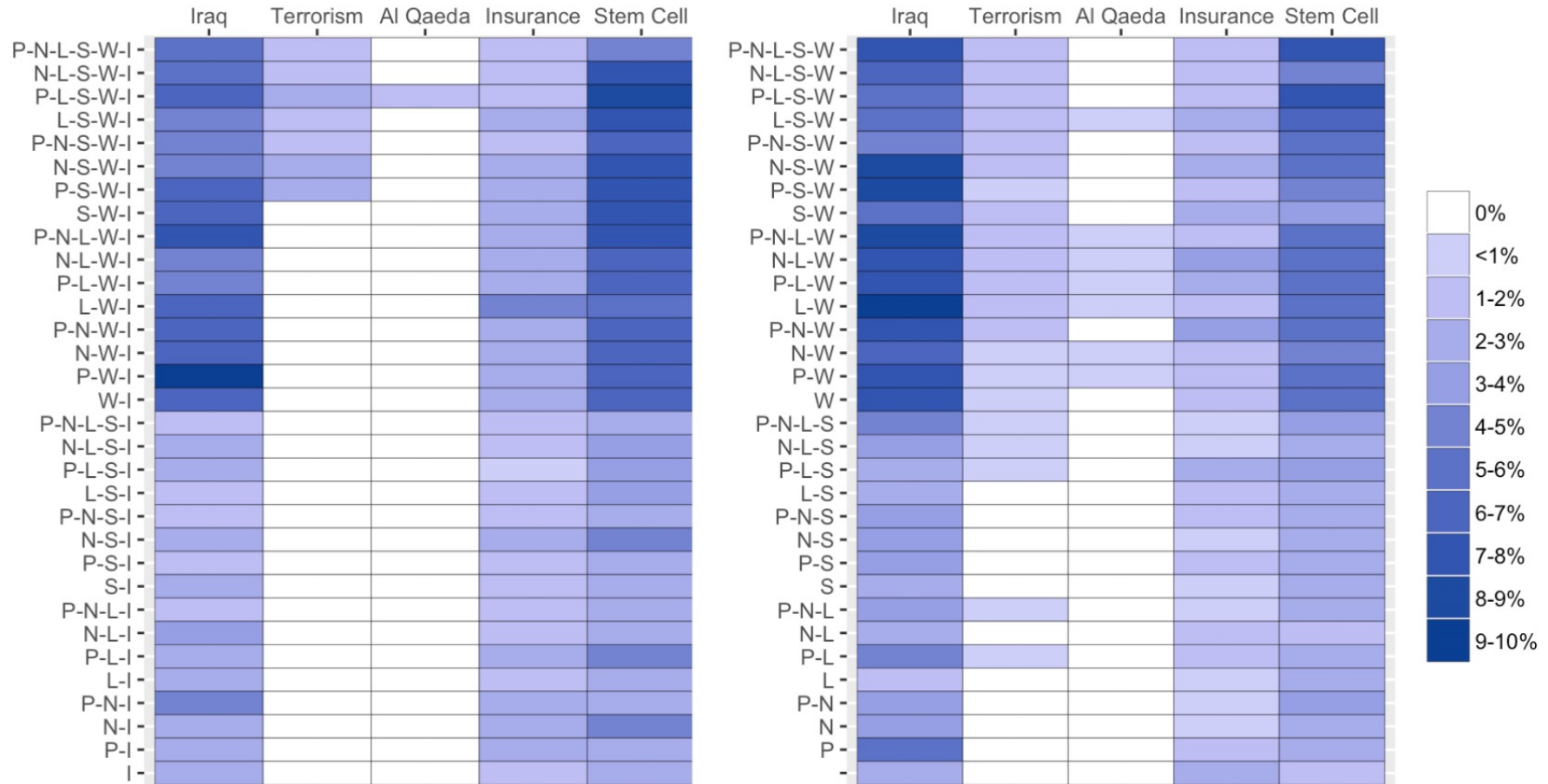


Figure 1. Order of text analysis operations for data preparation and analysis.

資料前處理為什麼重要？

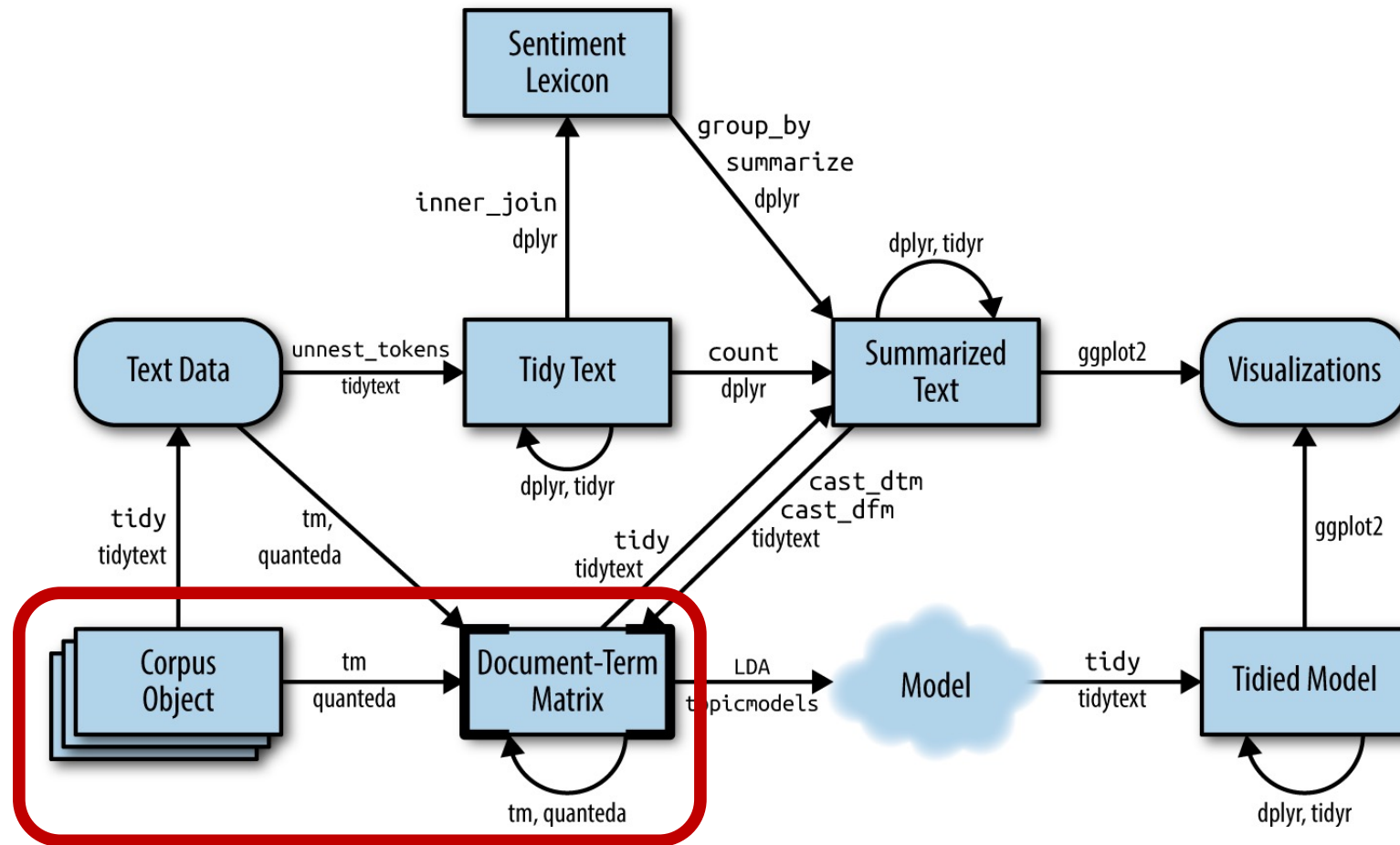


Denny, M. J., & Spirling, A. (2018). Text Preprocessing For Unsupervised Learning: Why It Matters, When It Misleads, And What To Do About It. *Political Analysis*, 26(02), 168-189.

語料庫 (corpus)

- 指大量的文本，通常經過整理，具有既定格式與標記。
- Project Gutenberg (<https://www.gutenberg.org/>)
- Manifesto Project (<https://manifesto-project.wzb.eu/>)
- HKBU Corpus of Political Speeches
(<https://digital.lib.hkbu.edu.hk/corpus/>)
- MPQA (<http://mpqa.cs.pitt.edu/>)

tm 套件



https://www.tidytextmining.com/images/tmwr_0601.png

stringr

String manipulation with stringr : : CHEAT SHEET

The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.



Detect Matches



str_detect(string, **pattern**) Detect the presence of a pattern match in a string.
`str_detect(fruit, "a")`



str_which(string, **pattern**) Find the indexes of strings that contain a pattern match.
`str_which(fruit, "a")`



str_count(string, **pattern**) Count the number of matches in a string.
`str_count(fruit, "a")`



str_locate(string, **pattern**) Locate the positions of pattern matches in a string. Also **str_locate_all**.
`str_locate(fruit, "a")`

Subset Strings



str_sub(string, start = 1L, end = -1L) Extract substrings from a character vector.
`str_sub(fruit, 1, 3); str_sub(fruit, -2)`



str_subset(string, **pattern**) Return only the strings that contain a pattern match.
`str_subset(fruit, "b")`



str_extract(string, **pattern**) Return the first pattern match found in each string, as a vector. Also **str_extract_all** to return every pattern match.
`str_extract(fruit, "[aeiou]")`



str_match(string, **pattern**) Return the first pattern match found in each string, as a matrix with a column for each () group in pattern. Also **str_match_all**.
`str_match(sentences, "(a[the] ([^]+))")`

Manage Lengths



str_length(string) The width of strings (i.e. number of code points, which generally equals the number of characters). `str_length(fruit)`



str_pad(string, width, side = c("left", "right", "both"), pad = " ") Pad strings to constant width. `str_pad(fruit, 17)`



str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...") Truncate the width of strings, replacing content with ellipsis. `str_trunc(fruit, 3)`



str_trim(string, side = c("both", "left", "right")) Trim whitespace from the start and/or end of a string. `str_trim(fruit)`

Mutate Strings



str_sub() <- value. Replace substrings by identifying the substrings with `str_sub()` and assigning into the results.
`str_sub(fruit, 1, 3) <- "str"`



str_replace(string, **pattern**, replacement) Replace the first matched pattern in each string. `str_replace(fruit, "a", "-")`



str_replace_all(string, **pattern**, replacement) Replace all matched patterns in each string. `str_replace_all(fruit, "a", "-")`



str_to_lower(string, locale = "en")¹ Convert strings to lower case.
`str_to_lower(sentences)`



str_to_upper(string, locale = "en")¹ Convert strings to upper case.
`str_to_upper(sentences)`



str_to_title(string, locale = "en")¹ Convert strings to title case. `str_to_title(sentences)`

Join and Split



str_c(..., sep = "", collapse = NULL) Join multiple strings into a single string.
`str_c(letters, LETTERS)`



str_c(..., sep = "", collapse = NULL) Collapse a vector of strings into a single string.
`str_c(letters, collapse = "")`



str_dup(string, times) Repeat strings times times. `str_dup(fruit, times = 2)`



str_split_fixed(string, **pattern**, n) Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also **str_split** to return a list of substrings.
`str_split_fixed(fruit, " ", n=2)`



str_glue(..., sep = "", envir = parent.frame()) Create a string from strings and {expressions} to evaluate. `str_glue("Pi is {pi}")`



str_glue_data(.x, ..., sep = "", envir = parent.frame(), na = "NA") Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate.
`str_glue_data(mtcars, "{rownames(mtcars)} has {hp} hp")`

Order Strings



str_order(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...) Return the vector of indexes that sorts a character vector. `x[str_order(x)]`



str_sort(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...) Sort a character vector.
`str_sort(x)`

Helpers



str_conv(string, encoding) Override the encoding of a string. `str_conv(fruit, "ISO-8859-1")`

str_view(string, **pattern**, match = NA) View HTML rendering of first regex match in each string. `str_view(fruit, "[aeiou]")`



str_view_all(string, **pattern**, match = NA) View HTML rendering of all regex matches. `str_view_all(fruit, "[aeiou]")`

str_wrap(string, width = 80, indent = 0, exdent = 0) Wrap strings into nicely formatted paragraphs. `str_wrap(sentences, 20)`

Regular Expression (regex)

- 正規表示式，又稱正則表達式。
- 表示文字的模式 (pattern)
- [RStudio cheatsheet](#)
- R regex 練習網站 (<https://regexr.com>)

Basic Regular Expressions in R Cheat Sheet

Character Classes	
<code>[[digit:]]</code> or <code>\d</code>	Digits; [0-9]
<code>\D</code>	Non-digits; [^0-9]
<code>[[lower:]]</code>	Lower-case letters; [a-z]
<code>[[upper:]]</code>	Upper-case letters; [A-Z]
<code>[[alpha:]]</code>	Alphabetic characters; [A-z]
<code>[[alnum:]]</code>	Alphanumeric characters; [A-z0-9_]
<code>\w</code>	Word characters; [A-z0-9_]
<code>\W</code>	Non-word characters
<code>[[xdigit:]]</code> or <code>\x</code>	Hexadecimal digits; [0-9A-Fa-f]
<code>[[blank:]]</code>	Space and tab
<code>[[space:]]</code> or <code>\s</code>	Space, tab, vertical tab, newline, form feed, carriage return
<code>\S</code>	Not space; [^[:space:]]
<code>[[punct:]]</code>	Punctuation characters; [!\"#\$%&'()*+,-./:;<=>?@[\]^_`{ }~]
<code>[[graph:]]</code>	Graphical characters; [[:alnum:]]{[:punct:]]
<code>[[print:]]</code>	Printable characters; [[:alnum:]]{[:punct:]]\s
<code>[[cntrl:]]</code> or <code>\c</code>	Control characters; \n, \r etc.

Special Metacharacters

<code>\n</code>	New line
<code>\r</code>	Carriage return
<code>\t</code>	Tab
<code>\f</code>	Vertical tab
<code>\v</code>	Form feed

Lookarounds and Conditionals*

<code>(?=)</code>	Lookahead (requires PERL = TRUE), e.g. <code>(?=xy)</code> : position followed by 'xy'
<code>(?!)</code>	Negative lookahead (PERL = TRUE); position NOT followed by pattern
<code>(?<=)</code>	Lookbehind (PERL = TRUE), e.g. <code>(?<=xy)</code> : position following 'xy'
<code>(?<!)</code>	Negative lookbehind (PERL = TRUE); position NOT following pattern
<code>?(i)then</code>	If-then-condition (PERL = TRUE); use lookaheads, optional char. etc in if-clause
<code>?(i)thenelse</code>	If-then-else-condition (PERL = TRUE)

*see, e.g. <http://www.regular-expressions.info/lookaround.html>
<http://www.regular-expressions.info/conditional.html>

Functions for Pattern Matching

Detect pattern

Locate pattern

Extract pattern

Replace pattern

Extract Patterns

```
regmatches(string, regex(pattern, string))  
# extract first match [1] "tam" "tim"  
regmatches(string, gregexpr(pattern, string))  
# extract all matches, outputs a list  
[[1]] "tam" [[2]] character(0) [[3]] "tim" "tom"  
stringr::str_extract(string, pattern)  
# extract first match [1] "tam" NA "tim"  
stringr::str_extract_all(string, pattern)  
# extract all matches, outputs a list  
stringr::str_match(string, pattern)  
# extract first match + individual character groups  
stringr::str_match_all(string, pattern, simplify = TRUE)  
# extract all matches, outputs a matrix  
stringr::str_locate(string, pattern)  
# extract first match + individual character groups  
stringr::str_locate_all(string, pattern)  
# extract all matches + individual character groups
```

Replace Patterns

```
sub(pattern, replacement, string)  
# replace first match  
gsub(pattern, replacement, string)  
# replace all matches  
stringr::str_replace(string, pattern, replacement)  
# replace first match  
stringr::str_replace_all(string, pattern, replacement)  
# replace all matches
```

Detect Patterns

```
grep(pattern, string)  
[1] 1 3  
grep(pattern, string, value = TRUE)  
[1] "Hipopotamus"  
[2] "time for bottomless lyrics"  
grepl(pattern, string)  
[1] TRUE FALSE TRUE  
stringr::str_detect(string, pattern)  
[1] TRUE FALSE TRUE
```

Locate Patterns

```
regexpr(pattern, string)  
# find starting position and length of first match  
gregexpr(pattern, string)  
# find starting position and length of all matches  
stringr::str_locate(string, pattern)  
# find starting and end position of first match  
stringr::str_locate_all(string, pattern)  
# find starting and end position of all matches
```

Split a String using a Pattern

```
strsplit(string, pattern) or stringr::str_split(string, pattern)
```

Character Classes and Groups

<code>.</code>	Any character except <code>\n</code>
<code>[]</code>	Or, e.g. <code>[a b]</code>
<code>[]</code>	List permitted characters, e.g. <code>[abc]</code>
<code>[a-z]</code>	Specify character ranges
<code>[^]</code>	List excluded characters
<code>(...)</code>	Grouping, enables back referencing using <code>\N</code> where N is an integer

General Modes

By default R uses *extended* regular expressions. You can switch to *PCRE* regular expressions using `PERL = TRUE` for base or by wrapping patterns with `perl()` for stringr.

All functions can be used with literal searches using `fixed = TRUE` for base or by wrapping patterns with `fixed()` for stringr.

All base functions can be made case insensitive by specifying `ignore.case = TRUE`.

Character Classes and Groups

<code>.</code>	Any character except <code>\n</code>
<code>[]</code>	Or, e.g. <code>[a b]</code>
<code>[]</code>	List permitted characters, e.g. <code>[abc]</code>
<code>[a-z]</code>	Specify character ranges
<code>[^]</code>	List excluded characters
<code>(...)</code>	Grouping, enables back referencing using <code>\N</code> where N is an integer

Escaping Characters

Metacharacters (`.`, `*`, etc.) can be used as literal characters by escaping them. Characters can be escaped using `\\` or by enclosing them in `\\Q...\\E`.

Case Conversions

Regular expressions can be made case insensitive using `(?i)`. In backreferences, the strings can be converted to lower or upper case using `\\L` or `\\U` (e.g. `\\L\\U`). This requires `PERL = TRUE`.

Quantifiers

<code>*</code>	Matches at least 0 times
<code>+</code>	Matches at least 1 time
<code>?</code>	Matches at most 1 time; optional string
<code>{n}</code>	Matches exactly n times
<code>{n,}</code>	Matches at least n times
<code>{n,m}</code>	Matches between n and m times

Greedy Matching

By default the asterisk `*` is greedy, i.e. it always matches the longest possible string. It can be used in lazy mode by adding `?`, i.e. `*?`.

Greedy mode can be turned off using `(?U)`. This switches the syntax, so that `(?U)*` is lazy and `(?U)+?` is greedy.

Note

Regular expressions can conveniently be created using e.g. the packages `rex` or `rebus`.