
第 8 章: 文字函式

8: Functions for Characters and Strings

傳統上, 統計人員較少直接處裡文字或字串資料, 多數時候是由資料管理人元處理後, 轉換成數值資料, 然後交由統計人員進行後續分析. 由於大數據時代來臨包含者多樣性的資料型態, 統計人員必須必須直接處裡文字或字串資料的機會也越來越多. R 內有許多函數可以處理文字型態的資料物件 或 文字資料 (Character Data), 常用之文字函式有 `paste()`, `substr()`, `substring()`, `grep()`, `gsub()`, `strsplit()` 等. R 套件 `stringr` 有更多處理文字或字串資料函式.

8.1 文字合併函式: `paste()`

文字合併函式 `paste()` 可以合併 2 個文字向量中的元素使其合而為 1 個文字向量.

```
1 > paste (... , sep = " ", collapse = NULL)
2 > paste0(... , collapse = NULL)
```

其中引數為

- `sep = " "` 設定 2 個字串向量合併時, 各別元素中間分隔使用的文字或符號. 回傳字串長度, 長度為原有字串向量的長度, R 內設為 1 各空白鍵.

- `collapse` 先依照 `sep` 設定回傳字串長度後，然後再次合併此字串向量的元素，後成為回傳單一字串，長度為 1。 `collapse` 設定最後元素中間分隔使用的文字或符號。

```

1 > ## paste()
2 > paste(1:5) # same as as.character(1:5)
3 [1] "1" "2" "3" "4" "5"
4 > paste("Today is", date())
5 [1] "Today is Sat Mar 11 11:37:40 2017"
6 > #
7 > paste("A", 1:5, sep = "")
8 [1] "A1" "A2" "A3" "A4" "A5"
9 > paste("A", 1:5, sep = " ")
10 [1] "A 1" "A 2" "A 3" "A 4" "A 5"
11 > paste("A", 1:5, sep = "#")
12 [1] "A#1" "A#2" "A#3" "A#4" "A#5"
13 > #
14 > paste(c("X", "Y"), 1:5)
15 [1] "X 1" "Y 2" "X 3" "Y 4" "X 5"
16 > paste(c("X", "Y"), 1:5, sep = " ")
17 [1] "X 1" "Y 2" "X 3" "Y 4" "X 5"
18 > paste(c("X", "Y"), 1:5, sep = "")
19 [1] "X1" "Y2" "X3" "Y4" "X5"
20 > paste(c("X", "Y"), 1:5, sep = "+")
21 [1] "X+1" "Y+2" "X+3" "Y+4" "X+5"
22 > paste(c("X", "Y"), 1:5, sep = "", collapse = NULL)
23 [1] "X1" "Y2" "X3" "Y4" "X5"
24 > paste(c("X", "Y"), 1:5, sep = "", collapse = " + ")
25 [1] "X1 + Y2 + X3 + Y4 + X5"

```

8.2 萃取文字函式: substr() 與 substring()

文字函式 `substr()` 與 `substring()` 可在 1 個文字向量中，萃取傳回部分字串，或替換部分字串。

```

1 > substr(x.char, start, stop)
2 > substring(x.text, first, last = 1000000L)
3 > substr(x.char, start, stop) <- new.char
4 > substring(x.text, first, last = 1000000L) <- new.char

```

其中引數

- `x.char` 為 1 個文字字串向量，

- `start`, 為 1 個正整數, 表示所要萃取文字字串的第 1 個文字位置,
- `stop`, 為 1 個正整數, 表示所要萃取文字字串的最後 1 個文字位置.
- `x.text` 為 1 個文字字串向量.
- `first` 為 1 個正整數, 表示所要萃取文字字串的第 1 個文字位置.
- `last` 為 1 個正整數, 表示所要萃取文字字串的最後 1 個文字位置.
- 使用其它文字字串來取代原先所萃取之文字字串:
 - `substr(x.char, start, stop) <- new.char`
 - `substring(x.text, first, last = 1000000) <- new.text`

可以指派新的文字字串 `new.char` 或 `new.text` 取代原先所萃取之文字字串之位置.

利用 R 內建資料框架 `state`, 關於美國 50 洲的一些相關資料, 包含向量物件, `state.abb`, `state.area`, `state.center`, `state.division`, `state.name`, `state.region`, 以及矩陣物件 `state.x77`, 在此使用 `substr()` 與 `substring()` 函式, 分析文字向量物件 `state.name`.

```

1 > ## substr()
2 > data(state)
3 > state.name
4 [1] "Alabama"      "Alaska"       "Arizona"      "Arkansas"     "California"
5 [6] "Colorado"     "Connecticut"   "Delaware"     "Florida"      "Georgia"
6 [11] "Hawaii"      "Idaho"        "Illinois"     "Indiana"      "Iowa"
7 [16] "Kansas"       "Kentucky"     "Louisiana"    "Maine"        "Maryland"
8 [21] "Massachusetts" "Michigan"    "Minnesota"    "Mississippi" "Missouri"
9 [26] "Montana"      "Nebraska"     "Nevada"       "New Hampshire" "New Jersey"
10 [31] "New Mexico"   "New York"     "North Carolina" "North Dakota"  "Ohio"
11 [36] "Oklahoma"     "Oregon"       "Pennsylvania" "Rhode Island" "South Carolina"
12 [41] "South Dakota" "Tennessee"   "Texas"        "Utah"         "Vermont"
13 [46] "Virginia"     "Washington"  "West Virginia" "Wisconsin"   "Wyoming"
14 > state.name[47:50]
15 [1] "Washington"   "West Virginia" "Wisconsin"   "Wyoming"
16 > substr(state.name[47:50], start = 1, stop = 4)
17 [1] "Wash" "West" "Wisc" "Wyom"
18 > substr(state.name[47:50], start = 1, stop = 4) <- "AAA"
19 > state.name[47:50]
```

```

20 [1] "AAAhington"    "AAAt Virginia" "AAAconsin"      "AAAming"
21 > #
22 > ## substring()
23 > data(state)
24 > state.name[47:50]
25 [1] "Washington"     "West Virginia" "Wisconsin"      "Wyoming"
26 > substring(state.name[47:50], first = 3, last = 1000)
27 [1] "shington"       "st Virginia"  "sconsin"       "oming"
28 > substring(state.name[47:50], first = 3, last = 1000) <- "BBB"
29 > state.name[47:50]
30 [1] "WaBBBngton"    "WeBBBVirginia" "WiBBBnsin"     "WyBBBng"

```

8.3 簡化文字函式: abbreviate()

文字函式 `abbreviate()` 方便用來簡化文字字串，使用如下：

```

1 > abbreviate(names.arg, minlength = 4, use.classes = TRUE,
2               dot = FALSE, strict = FALSE,
3               method = c("left.kept", "both.sides"))

```

其中引數

- `names.arg` 為 1 個向文字字串向量。
- `minlength = 4` 為簡化後文字字串的字元數目。
- `dot = FALSE` 為邏輯向量，是否須要附加“.”在簡化的文字字串內。
- `strict = FALSE` 為邏輯向量，為邏輯向量，是否須要絕對遵循 `minlength` 的設定，若 `strict = TRUE` 則可能不會形成唯一的簡化，某些元素原本是不同文字，這些元素在簡化後可能會變成相同的文字。

```

1 > ## abbreviate()
2 > data(state)
3 > state.name[47:50]
4 [1] "Washington"     "West Virginia" "Wisconsin"      "Wyoming"
5 > abbreviate(state.name[47:50], minlength = 2, dot = TRUE)
6   Washington West Virginia Wisconsin      Wyoming
7   "Wsh."           "WV."          "Wsc."         "Wy."
8 >
9 > x.chr <- c("abcd", "efgh", "abce")

```

```
10 > abbreviate(x.chr, 2)
11 abcd egh abce
12 "abcd" "ef" "abce"
13 > abbreviate(x.chr, 2, strict = TRUE) # 1st and 3rd are == "ab"
14 abcd egh abce
15 "ab" "ef" "ab"
16 > abbreviate(x.chr, 2, strict = FALSE) #
17 abcd egh abce
18 "abcd" "ef" "abce"
19 >
20 > x.abb = abbreviate(state.name[47:50], minlength = 2, method = "left.kept")
21 > y.abb = abbreviate(state.name[47:50], minlength = 2, method = "both")
22 > cbind(x.abb, y.abb)
23          x.abb y.abb
24 Washington    "Wsh" "Ws"
25 West Virginia "WV"   "WV"
26 Wisconsin     "Wsc" "Wn"
27 Wyoming       "Wy"   "Wy"
```

8.4 文字搜尋函式: grep()

文字處理必須要有一個重要的功能，就是在字串中 搜尋配對，修正 與 取代 關鍵字 或 特定型式字串，或取得特定型式字串所在的下標位置。特定型式字串 R 稱為 regular expression (Regular expression is a pattern that describes a set of strings)，程式語言 Perl 與 Python 都具備此功能。R 也可用 Perl 的搜尋配對或取代方式。regular expression 是由搜尋一個單字逐步形成搜尋特定型式字串，執行搜尋取代功能前有些事項必須先注意，尤其是 metacharacters (特殊符號字元)。

- **single character**: 搜尋或取代單字。
- **character class**: 同時搜尋或取代多個單字。
- **concatenation**: 搜尋單字且合併形成為特定型式字串。
- **repetition**: 搜尋或取代特定型式字串的次數。
- **alternation**: 交互輪流，在多個特定型式字串中，搜尋或取代某一特定型式字串。

- 使用萬用字元符號或是搜尋標點符號.
- **metacharacters**: 特殊符號字元 以及 R 保留的關鍵字. 搜尋取代特殊符號字元必須在特殊符號字元前加上雙反斜線 "\\".

```
1 ! " # $ % & àÁŽ ( ) * + , - . / : ; < = > ? @ [ \ ] ^ - ' { }
```

萬用字元符號與 特殊符號字元參見 表 8.1 - 表 8.2.

```
1 > grep(pattern, x, ignore.case = FALSE, perl = FALSE, value = FALSE,
2     fixed = FALSE, useBytes = FALSE, invert = FALSE)
3
4 > grep(pattern, x, ignore.case = FALSE, perl = FALSE,
5     fixed = FALSE, useBytes = FALSE)
6
7 > sub(pattern, replacement, x, ignore.case = FALSE, perl = FALSE,
8     fixed = FALSE, useBytes = FALSE)
9
10 > gsub(pattern, replacement, x, ignore.case = FALSE, perl = FALSE,
11    fixed = FALSE, useBytes = FALSE)
12
13 > regexpr(pattern, text, ignore.case = FALSE, perl = FALSE,
14    fixed = FALSE, useBytes = FALSE)
15
16 > gregexpr(pattern, text, ignore.case = FALSE, perl = FALSE,
17    fixed = FALSE, useBytes = FALSE)
18
19 > regexec(pattern, text, ignore.case = FALSE,
20    fixed = FALSE, useBytes = FALSE)
```

其中引數

- **pattern**: 特定型式字串.
- **x, text**: 文字向量或物件.
- **ignore.case = FALSE**: 設定大小寫有差異不可忽略.
- **perl = FALSE**: 使用 perl 程式語言的特定型式字串.
- **value = FALSE**: 回傳向量正整數下標, **value = TRUE**: 回傳向量符合特定型式的元素.
- **fixed = FALSE**: 設定特定型式字串是否須完全相符合.

- `useBytes = FALSE`: 設定搜尋取代為 `character-by-character` 而非 `byte-by-byte`.
- `invert = TRUE`: 設定搜尋取代為反向操作, 傳回不符合特定型式字串的下標位置或文字.

8.4.1 搜尋或取代單字

任何搜尋或取代方式都是由搜尋一個單字逐步形成搜尋特定型式字串. 例如, 在文字向量 `x.chr` 搜尋單字 `x`, 使用函式 `grep()` 在萃取文字向量 `x.chr` 中, 尋找的某一個特定式樣 (pattern), 並且傳回文字向量的下標 (index).

```
1 > ## single character
2 > ## grep()
3 > x.chr <- c("y", "ya bc", "abc", "ay bc", "ab yc", "ab cy")
4 > grep("y", x.chr)
5 [1] 1 2 4 5 6
6 > x.chr[grep("y", x.chr)] # index by position
7 [1] "y"      "ya bc" "ay bc" "ab yc" "ab cy"
```

特定型式字串的前一位置 (pre-pattern) 使用 `^` 字元 (caret, 脫字符號, 倒 V 符號), 則僅搜尋 1 行字串的第 1 個字元是否符合特定型式字串.

```
1 > ## ^ (caret) line first
2 > x.chr <- c("y", "ya bc", "abc", "ay bc", "ab yc", "ab cy")
3 > grep("^y", x.chr)
4 [1] 1 2
5 > x.chr[grep("^y", x.chr)]
6 [1] "y"      "ya bc"
```

特定型式字串的後一位置 (pattern-post) 使用 `$` 字元 (dollar sign, 金錢符號), 則僅搜尋 1 行字串的最後 1 個字元是否符合特定型式字串. `metacharacters`.

```
1 > ## $ (dollar sign) line last
2 > x.chr <- c("y", "ya bc", "abc", "ay bc", "ab yc", "ab cy")
3 > grep("y$", x.chr)
4 [1] 1 6
5 > x.chr[grep("y$", x.chr)]
6 [1] "y"      "ab cy"
```

特殊符號字元 `"."` (句點, period), 搜尋任何符合特定型式字串的字串, 包含空格.

`^` 字元, `$` 字元 與 `"."` 字元 是所謂的特殊符號字元

```
1 > x.chr[grep(".bc", x.chr)]
2 [1] "ya bc" "abc"    "ay bc"
```

僅搜尋文字 (word) 的第 1 個字元是否符合特定型式字串, 但不要搜尋 1 行字串的第 1 個字元, 則特定型式字串的前一位置增加 1 空格. 同樣的, 僅搜尋文字 (word) 的最後 1 個字元是否符合特定型式字串, 但不要搜尋 1 行字串的最後 1 個字元, 則特定型式字串的後一位置增加 1 空格.

```
1 > ## word first, except line first
2 > x.chr <- c("y", "ya bc", "abc", "ay bc", "ab yc", "ab cy")
3 > grep(" y", x.chr)
4 [1] 5
5 > x.chr[grep(" y", x.chr)]
6 [1] "ab yc"
7 > #
8 > ## word last, except line last
9 > x.chr <- c("y", "ya bc", "abc", "ay bc", "ab yc", "ab cy")
10 > grep("y ", x.chr)
11 [1] 4
12 > x.chr[grep("y ", x.chr)]
13 [1] "ay bc"
```

8.4.2 同時搜尋或取代多個單字

R 在 [regular expression](#) 中定義 character class 是由中括號包圍的特定型式字串 [\[and\]](#), 且 [grep\(\)](#) 搜尋時, 只要中括號包圍內的字元中的任一個單字符合, 都會回傳. 例如, 若特定型式字串為 "[cgr]" 1 個字元只要與 "c", "g", "r" 任一個單字符合, 都會回傳. 若特定型式字串為 "^ [cgr]", 則文字向量內的 1 行的第 1 個字只要與 "c", "g", "r" 任一個單字符合, 都會回傳. 反之, 若 "^[^cgr]" 則文字向量內的 1 行的第 1 個字元不含任一個單字 "c", "g", "r" 才會回傳. 參見 表 8.1 - 表 8.2.

```
1 > ## character class
2 > x.chr <- c("cat", "bar", "rat", "sac", "mug", "gate")
3 > grep("[cgr]", x.chr)
4 [1] 1 2 3 4 5 6
5 > x.chr[grep("[cgr]", x.chr)]
6 [1] "cat"  "bar"  "rat"  "sac"  "mug"  "gate"
7 > #
8 > x.chr <- c("cat", "bar", "rat", "sac", "mug", "gate")
```

```

9 > grep("^[cgr]", x.chr)
10 [1] 1 3 6
11 > x.chr[grep("^[cgr]", x.chr)]
12 [1] "cat" "rat" "gate"
13 > #
14 > x.chr <- c("cat", "bar", "rat", "sac", "mug", "gate")
15 > grep("^[^cgr]", x.chr)
16 [1] 2 4 5
17 > x.chr[grep("^[^cgr]", x.chr)]
18 [1] "bar" "sac" "mug"

```

例如, "`^.[a].+`" 搜尋任何第 1 個單字, 接著是非 a 單字, 接續為至少 1 個或多個單字.

```

1 > x.chr <- c("cat", "bar", "rat", "sac", "mug", "gate")
2 > grep("^.[a].+", x.chr)
3 [1] 5
4 > x.chr[grep("^.[a].+", x.chr)]
5 [1] "mug"

```

文字搜尋時, 也可使用文字範圍設定, 例如,

- `[0-9]`: 搜尋數字 0-9.
- `[A-Z]`: 搜尋大寫字母 A-Z.
- `[a-z]`: 搜尋小寫字母 a-z.
- `[0-9A-Za-z]`: 同時搜尋數字與字母.

表 8.1 顯示更多的用法.

8.4.3 合併成特定字串

搜尋單一字母可以搜尋特定字合併使用, 例如, "`^[cgr]at$`", 搜尋 1 行的第 1 個字元為任 1 個 `c`, `g`, `r`, `$` 搜尋 1 行字串的最後字元為 `at`, 因此可搜尋出 `"cat"`, `"rat"` 等字. 若要 3 個字元, 可使用 "`^[ct]a[br]$`".

```

1 > ##
2 > x.chr <- c("cat", "car", "bar", "rat", "sac", "mug", "gate", "tab")
3 > grep("^[ct]a[br]$", x.chr)
4 [1] 2 8

```

```

5 > x.chr[grep("^bct]a[br]$", x.chr)]
6 [1] "car" "bar" "tab"
7 > #
8 > x.chr <- c("cat", "car", "bar", "rat", "sac", "mug", "gate", "tab")
9 > grep("^[ct]a[br]$", x.chr)
10 [1] 2 8

```

句點 (period) "." 為 **metacharacters**, 作為萬用字元, 可以搜尋配對任 1 個字元. 例如, "**c.t**" 搜尋配對啟始為 **c** + 任何字元 + **t**.

```

1 > ## period (.)
2 > x.chr <- c("cat", "cut", "bar", "cattle", "sat", "cup", "city")
3 > grep("c.t", x.chr)
4 [1] 1 2 4 7
5 > x.chr[grep("c.t", x.chr)]
6 [1] "cat"     "cut"     "cattle"  "city"

```

8.5 搜尋或取代特定型式字串次數

搜尋特定型式字串, 並不受限 1 次, 可以設定在一字符串中多次重覆搜尋取代, 或設定重覆搜尋取代的次數, 參見 表 8.2. 例如, "+" 表示搜尋配對同一元素內字符串至少 1 次以上, "^ [cC].+e\$" 搜尋配對啟始為 **c** 或 **C** + 任何字元 + 最後字元為 **e**.

```

1 > ## Repetition
2 > x.chr <- c("cat", "cute", "bar", "cattle", "sat", "cup", "city", "Centigrade")
3 > grep("^[cC].+e$", x.chr)
4 [1] 2 4 8
5 > x.chr[grep("^[cC].+e$", x.chr)]
6 [1] "cute"      "cattle"    "Centigrade"

```

設定重覆搜尋取代的次數時, 可以使用大括號 "{}" 包含設定的次數. 例如, 搜尋 1 個字元 "m" 至少 2 次, 可以使使用 "m{2,}". 若要搜尋多個字元 "tom" 至少 2 次, 可以使使用 "(tom){2,}".

```

1 > ## Repetition
2 > x.chr <- c("Tom", "Thomas", "Tomas", "Tommy", "tomas", "tommas", "Tomtom")
3 > grep("m{2,}", x.chr)
4 [1] 4 6
5 > x.chr[grep("m{2,}", x.chr)]
6 [1] "Tommy"    "tommas"
7 > #

```

```

8 > x.chr <- c("Tom", "Thomas", "Tomas", "Tommy", "tomas", "tommas", "Tomtom")
9 > grep("[Tt]om{2,}", x.chr)
10 [1] 4 6
11 > x.chr[grep("[Tt]om{2,}", x.chr)]
12 [1] "Tommy"  "tommas"
13 > #
14 > x.chr <- c("Tom", "Thomas", "Tomas", "Tommy", "tomas", "tommas", "Tomtom")
15 > grep("[Tt]om{2,}", x.chr)
16 [1] 7
17 > x.chr[grep("[Tt]om{2,}", x.chr)]
18 [1] "Tomtom"

```

8.5.1 交互輪流搜尋或取代

資料處裡有時須同時搜尋配對 2 種以上特定型式字串，此使可以使用 " | " (vertical bar) 結合 2 種以上特定型式字串。例如，B 型肝炎 ICD-10 疾病診斷碼，

```

1 B16, B16.0, B16.1, B16.2, B16.3, B16.4,
2 B16.5, B16.7, B16.8, B16.9, B17, B17.0, B17.2, B17.8,
3 B18, B18.0, B18.1, B18.8, B18.9

```

但是，保險資料庫內的 B 型肝炎疾病診斷碼並未包含 "." (period)，因此呈現為

```

1 B16, B160, B161, B162, B163, B164, B165, B167, B168, B169,
2 B17, B170, B172, B178,
3 B18, B180, B181, B188, B189

```

若要將保險資料庫內的疾病診斷碼與正式的 ICD 疾病診斷碼搜尋配對，可使用

"^B16[0-9]?|\$|^B17[0,2,8]?|\$|^B18[0,1,8,9]?\$"

上述搜尋配對 "¹^B16[0-9]?\$" 或 "²^B17[0,2,8]?\$" 或 "³^B18[0,1,8,9]?\$>".

```

1 > ## Alternation
2 > hbv.chr <- c("B16", "B160", "B161", "B162", "B163", "B164",
3           "B165", "B167", "B168", "B169",
4           "B17", "B170", "B172", "B178",
5           "B18", "B180", "B181", "B188", "B189")
6 > grep("B16[0-9]?$", hbv.chr)      # B16
7 [1] 1 2 3 4 5 6 7 8 9 10
8 > grep("B17[0,2,8]?$", hbv.chr)    # B17
9 [1] 11 12 13 14
10 > grep("B18[0,1,8,9]?$", hbv.chr) # B18
11 [1] 15 16 17 18 19
12 > #
13 > grep("B16[0-9]?|^B17[0,2,8]?|^B18[0,1,8,9]?$", hbv.chr) # any

```

```

14 [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19
15 > ## misspelling name
16 > x.chr <- c("Tom", "Thomas", "Tomas", "Tommy", "tomas")
17 > grep("^[Tt]omm?y?$^*[Tt]homas$+tomas$", x.chr)
18 [1] 1 2 4
19 > x.chr[grep("^[Tt]omm?y?$^*[Tt]homas$+tomas$", x.chr)]
20 [1] "Tom"     "Thomas" "Tommy"

```

8.5.2 特殊字元與標點符號

R 在搜尋或取代時會用到 metacharacters (特殊符號字元) 或 萬用字元 方便搜尋或取代, 以下為常用萬用字元的搜尋或取代.

```

1 "[0-9]" - Digits
2 "[a-z]" - Lower-case letters
3 "[A-Z]" - Upper-case letters
4 "[a-zA-Z]" - Alphabetic characters
5 "[^a-zA-Z]" - Non-alphabetic characters
6 "[a-zA-Z0-9]" - Alphanumeric characters
7 "[ \t\n\r\f\v]" - Space characters
8 "[\$*+.?[^{(\\"#\%&\sim_/< = > !,:;\'")}]@-]" - Punctuation Characters

1 > ## any string containing any of the characters
2 > x.chr <- c("nose", "letter38", "window9", "apple0")
3 > grep("[nco]", x.chr, value = TRUE)
4 [1] "nose"      "window9"
5 > grep("[01234567]", x.chr, value = TRUE)
6 [1] "letter38"  "apple0"
7 > grep("[0-7]", x.chr, value = TRUE)
8 [1] "letter38"  "apple0"

```

這些特殊符號字元也是標點符號, 因此在搜尋或取代標點符號必須特殊處理. 通常是使用雙反斜線 \\, 例如, 使用 `gsub()` 搜尋 "\$", 取代 "." (period).

```

1 > ## metacharacters
2 > s.chr = "gsub$uses$regular$expressions"
3 > s.err <- gsub(pattern = "$", replacement = ".", s.chr)
4 > s.err
5 [1] "gsub$uses$regular$expressions."
6 > s.correct <- gsub(pattern = "\\\"", replacement = ".", s.chr)
7 > s.correct
8 [1] "gsub.uses.regular.expressions"

```

在 character class 搜尋或取代這些字元的處理規則, 比較誤用與正確方式.

```
1 > ## metacharacters
2 > ## example: wrong methods
3 > metaChar = c("$", "*", "+", ".", "?", "[", "^", "{", "", "(", "\\")
4 > grep(pattern = "$", x = metaChar, value = TRUE)
5 > grep(pattern = "\\", x = metaChar, value = TRUE)
6 > grep(pattern = "(", x = metaChar, value = TRUE)
7 > gsub(pattern = "", replacement = ".",
8       "gsubusesregularexpressions")
9 > strsplit(x = "strsplit.aslo.uses.regular.expressions",
10        split = ".")
11 > ## examples: correct methods
12 > metaChar = c("$", "*", "+", ".", "?", "[", "^", "{", "", "(", "\\")
13 > grep(pattern = "\\\$", x = metaChar, value = TRUE)
14 [1] "$"
15 > grep(pattern = "$", x = metaChar, value = TRUE, fixed = TRUE)
16 [1] "$"
17 > grep(pattern = "\\\\", x = metaChar, value = TRUE)
18 [1] "\\"
19 > grep(pattern = "\\\\", x = metaChar, value = TRUE, fixed = TRUE)
20 [1] "\\"
21 > grep(pattern = "\\\\()", x = metaChar, value = TRUE)
22 [1] "("
23 > grep(pattern = "()", x = metaChar, value = TRUE, fixed = TRUE)
24 [1] "("
25 > #
26 > gsub(pattern = "\\\\", replacement = ".",
27       "gsubusesregularexpressions")
28 [1] "gsub.uses.regular.expressions"
29 > strsplit(x = "strsplit.aslo.uses.regular.expressions",
30        split = "\\.")
31 [[1]]
32 [1] "strsplit"      "aslo"          "uses"          "regular"        "expressions"
33
34 > strsplit(x = "strsplit.aslo.uses.regular.expressions",
35        split = ".", fixed = TRUE)
36 [[1]]
37 [1] "strsplit"      "aslo"          "uses"          "regular"        "expressions"
38 > #
39 > ## "\\\\+" at least once
40 x.chr <- c("id...of....patient", "patient....age")
41 gsub(pattern = "\\\\.+", replacement = ".", x = x.chr)
```

上述在搜尋或取代特殊字元時，在標點符號前必須使用雙反斜線 \\，但是以下字元仍須例外處裡。

```
1 ] ^ - \
```

- 搜尋 "]" 須放在 character class 第 1 位置

- “-”須放在 character class 第 1 位置 或 最後位置
- “^”不可放 character class 在第 1 位置 (因第 1 位置代表排除 character class 內的字元)
- 除了 “\”之外, 搜尋其他特殊字元時 可放在 character class 任何位置.

```

1 > ## metacharacters
2 > x.str <- "my\\strin\\g\\with\\slashes"
3 > x.str
4 [1] "my\\strin\\g\\with\\slashes"
5 > print(x.str)
6 [1] "my\\strin\\g\\with\\slashes"
7 > cat(x.str)
8 my\strin\g\with\slashes
9 > gsub("(\\J)", "", x.str)
10 [1] "mystringwithslashes"
11 > gsub("\\", "", x.str, fixed = TRUE)
12 [1] "mystringwithslashes"
13 > #

```

表 8.1: 常見萬用字元設定 |

萬用字元	替代簡寫	說明
[:lower:]	[a-z]	英文小寫字母 lower-case letters
[:upper:]	[A-Z]	英文大寫字母 upper-case letters
[:alpha:]	[A-Za-z]	英文大小寫字母 alphabetic characters
[:digit:]	[0-9]	數字字元 (Digits)
[:alnum:]	[A-Za-z0-9]	英文字母 + 數字 alphanumeric characters
[:punct:]	require \	標點符號 punctuation characters ! " # \$ % & , () * + , - . / : ; < = > ? @ [\] ^ _ ‘ { } ~
[:space:]	"[\t\n\r\f\v]"	space, newline vertical tab, ... form feed, carriage return,
[:graph:]	[:alnum:][:punct:]	graphical characters
[:print:]	[:alnum:][:punct:][:space:]	printable characters
[:xdigit:]	[0-9A-Fa-f]	hexadecimal digits

表 8.2: 常見萬用字元設定 II

萬用字元	說明
^	搜尋 1 行字串的第 1 個字元 " <code>^m</code> "
	搜尋 1 行字串的第 1 個字元不是特定字元 " <code>[^m]</code> "
?	搜尋配對同一元素內字串最多 1 次
*	搜尋配對同一元素內字串 0 次 或 無限多次
+	搜尋配對同一元素內字串至少 1 次以上
{n}	搜尋配對同一元素內字串最多 n 次
{n,}	搜尋配對同一元素內字串至少 n 次
{n, m}	搜尋配對同一元素內字串至少 n 次, 最多 m 次
<code>[anyone]</code>	character class, 中括號包圍的字串內的任一個單字符合 " <code>[]</code> " <code>[nco]</code> any string containing any of the characters: "n", "c", "o"
(group)	小包圍的字串群組 "([Tt]om){2,}" "[(" or ")"
	搜尋特定字合併使用 "Tomas Luis", " <code>[]</code> "
{	不適用在 R 的 regular expressions " <code>[]</code> "

8.6 文字分割函式: `strsplit()`

文字分割函式 `strsplit()` 可以根據字串內特定字元或文字將長字串分割成小字串。

```
1 strsplit(x, split, fixed = FALSE, perl = FALSE, useBytes = FALSE)
```

其中引數

- `x, text`: 分割前的字串向量或物件.
- `split`: 分割依據的特定字元, 可以是特定型式字串 regular expression.
`split = character(0)` 與 `split = ""` 是相同意義.
- `fixed = FALSE`: 設定特定型式字串是否須完全相符合.
- `perl = FALSE`: 使用 perl 程式語言的特定型式字串.

- `useBytes = FALSE`: 設定搜尋取代為 `character-by-character` 而非 `byte-by-byte`.

若分割依據的特定字元在字串的第 1 位置 或 最後位置, 則分割後會產生第一個 或 最後位置 元素為 "".

```
1 > ## strsplit()
2 > x.chr <- "asfef"
3 > strsplit(x.chr, "e") # # split x on the letter e
4 [[1]]
5 [1] "ASF" "f"
6 > #
7 > x.chr <- "Display a scentence with spaces"
8 > strsplit(x.chr, NULL)
9 [[1]]
10 [1] "D" "i" "s" "p" "l" "a" "y" " " "a" " " "s" "c" "e" "n" "t" "e" "n" "c" "e"
11 [20] " " "w" "i" "t" "h" " " "s" "p" "a" "c" "e" "s"
12
13 > noquote(strsplit(x.chr, NULL)[[1]])
14 [1] D i s p l a y   a   s c e n t e n c e   w i t h   s p a c e s
15 > #
16 > x.chr <- c(as = "asfef", qu = "qwerty", "yuiop[", "b", "stuff.blah.yech")
17 > strsplit(x.chr, "e") # # split x on the letter e
18 $as
19 [1] "ASF" "f"
20
21 $qu
22 [1] "qw" "rty"
23
24 [[3]]
25 [1] "yuiop["
26
27 [[4]]
28 [1] "b"
29
30 [[5]]
31 [1] "stuff.blah.y" "ch"
32
33 > #
34 > x.chr <- "a.b.c"
35 > strsplit(x.chr, ".")
36 [[1]]
37 [1] "a" "b" "c"
38
39 > unlist(strsplit(x.chr, "."))
40 [1] "a" "b" "c"
41 > strsplit(x.chr, "[.]")
42 [[1]]
43 [1] "a" "b" "c"
44
```

```
45 > unlist(strsplit(x.chr, "[.]"))
46 [1] "a" "b" "c"
47 > unlist(strsplit(x.chr, ".", fixed = TRUE))
48 [1] "a" "b" "c"
49 > #
50 > x.chr <- paste(c("", "a", ""), collapse = "#")
51 > x.chr
52 [1] "#a#"
53 > strsplit(x.chr, split = "#")
54 [[1]]
55 [1] "" "a"
56
57 > strsplit(x.chr, split = "#")[[1]]
58 [1] "" "a"
59 > #
60 > strsplit("", " ")[[1]]      # character(0)
61 character(0)
62 > strsplit(" ", " ")[[1]]    # [1] ""
63 [1] ""
```