

# 一般數據結構與格式part 2

報告者：盧森 & 阿舍

2017.06.26

## Factor 與 Levels

```
1 factor( v ) # v 必須是vector或是integer
2 wdays <- c("Wed","Fri","Wed","Tue","Tue","Fri","Wed","Thu")
3 f.wdays <- factor(wdays)
4 f.wdays
5 [1] Wed Fri Wed Tue Tue Fri Wed Thu
6 Levels: Fri Thu Tue Wed
7
8 #levels會告訴你有哪些內容（一種分類的感覺）
9 > levels(f.wdays)
10 [1] "Fri" "Thu" "Tue" "Wed"
11
12 #nlevels會告訴你這個factor中levels的數量
13 > nlevels(f.wdays)
14 [1] 4
```

```
1 #用tapply搭配factor格式來分類資料 有個matrix格式的Data如下：
2 Data
3      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
4 score.A "20" "40" "24" "56" "32" "53" "42" "43"
5 score.B "82" "59" "19" "74" "93" "52" "84" "74"
6 wdays  "Wed" "Fri" "Wed" "Tue" "Tue" "Fri" "Wed" "Thu"
7
8 f.wdays <- factor(wdays)
9 tapply(as.numeric(Data[1,]),f.wdays,summmary)
10 $Fri
11      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
12  40.00   43.25   46.50   46.50   49.75   53.00
13
14 $Thu
15      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
16      43      43      43      43      43      43
17
18 $Tue
19      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
20      32      38      44      44      50      56
21
22 $Wed
```

23	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
24	20.00	22.00	24.00	28.67	33.00	42.00

## 一些簡單好用的函數（vector）

```
1 #產生序列的函數 seq()
2 seq(1:6)
3 [1] 1 2 3 4 5 6
4
5 #產生重複字串的函數 rep()
6 rep("Hello",3)
7 [1] "Hello" "Hello" "Hello"
8
9 #字串黏合
10 paste(2016,month.abb[5],20) #未指定分隔符號
11 [1] "2016 May 20"
12 paste(2016,month.abb[5],20,sep = '-') #指定分隔符號
13 [1] "2016-May-20"
14
```

## vector 與 matrix

```
1 #從vector變成matrix的方法
2 seq(1:6) #產生一個vector，內容是序列 1到6
3 [1] 1 2 3 4 5 6
4
5 #方法一(上次有講過) matrix()
6 matrix(seq(1:6),nrow=6,ncol=1)
7      [,1]
8 [1,]    1
9 [2,]    2
10 [3,]    3
11 [4,]    4
12 [5,]    5
13 [6,]    6
14
15 matrix(seq(1:6),nrow = 1,ncol = 6)
16      [,1] [,2] [,3] [,4] [,5] [,6]
17 [1,]    1    2    3    4    5    6
18
19 #方法二 對 vector 使用cbind() 或 rbind()
20 cbind(seq(1:6))
21      [,1]
22 [1,]    1
23 [2,]    2
24 [3,]    3
25 [4,]    4
```

```

26 [5,]    5
27 [6,]    6
28
29 rbind(seq(1:6))
30      [,1] [,2] [,3] [,4] [,5] [,6]
31 [1,]    1    2    3    4    5    6

```

## 單行或單列的matrix是vector

```

1  matrix.example
2      [,1] [,2] [,3]
3 [1,]    2    1    5
4 [2,]    4    2   10
5 [3,]    6    3   15
6 [4,]    8    4   20
7 [5,]   10    5   25
8 [6,]   12    6   30
9
10 matrix.example[,1]
11 [1]  2  4  6  8 10 12
12
13 matrix.example[3,]
14 [1]  6  3 15
15
16 #因為是vector格式所以沒有維度(dimension)
17 dim(matrix.example[3,])
18 NULL
19
20 #有些函數只能用在matrix,所以當對vector格式使用該函數就會出現警告或是錯誤
21 #對matrix加上drop=F 可以讓單行或單列的matrix保留matrix的格式
22 matrix.example[3,,drop=F]
23      [,1] [,2] [,3]
24 [1,]    6    3   15
25
26 #加上drop=F後，因為格式仍然是matrix所以仍有維度資訊
27 dim(matrix.example[3,,drop=F])
28 [1] 1 3

```

## do.call()

```

1  #當想對vector格式的資料作資料合併時
2  A <- 1:10
3  B <- 100:91
4  C <- 78:87
5
6  #cbind()可以輕鬆幫你把他們合併在一起 （格式變成matrix）
7  cbind(A,B,C)

```

```

8      A    B    C
9  [1,]   1 100  78
10 [2,]   2  99  79
11 [3,]   3  98  80
12 [4,]   4  97  81
13 [5,]   5  96  82
14 [6,]   6  95  83
15 [7,]   7  94  84
16 [8,]   8  93  85
17 [9,]   9  92  86
18 [10,] 10  91  87
19
20 #但如果是list格式的資料呢？
21 list(A,B,C)
22 [[1]]
23 [1]  1  2  3  4  5  6  7  8  9 10
24
25 [[2]]
26 [1] 100  99  98  97  96  95  94  93  92  91
27
28 [[3]]
29 [1] 78 79 80 81 82 83 84 85 86 87
30
31 #我們會發現對list使用cbind()，會得到的結果並不是你想要的
32 cbind(list(A,B,C))
33      [,1]
34 [1,] Integer,10
35 [2,] Integer,10
36 [3,] Integer,10
37
38 #這時候 do.call()就是一個非常好用的幫手
39 do.call(cbind, list(A,B,C))
40      [,1] [,2] [,3]
41 [1,]    1 100   78
42 [2,]    2  99   79
43 [3,]    3  98   80
44 [4,]    4  97   81
45 [5,]    5  96   82
46 [6,]    6  95   83
47 [7,]    7  94   84
48 [8,]    8  93   85
49 [9,]    9  92   86
50 [10,]   10  91   87

```

## data.frame的小應用

```

1 #這是上次R basic我們分享的 data.frame 例子
2 x1 <- c("father","mother","brother","sister")

```

```
3 x2 <- c("Canon","Pentax","Olympus","Nikon")
4 x3 <- c("gold","red","green","blue")
5 x4 <- c(2,1,1,2)
6 camera <- data.frame(member = x1, brand = x2, color = x3, amount = x4)
7 camera
8     member  brand color amount
9 1 father   Canon  gold      2
10 2 mother  Pentax  red      1
11 3 brother Olympus green     1
12 4 sister   Nikon  blue     2
13
14 #只取用部分資料 subset()
15 #從camera這個資料框架中 挑選amount>1的資料 只顯示member欄位
16 subset(camera, select=member,subset=(amount>1))
17     member
18 1 father
19 4 sister
20
21 # select也可以用c()選擇顯示的多項欄位資料
22 # subset可以有很多條件，以&符號分隔選擇條件
```

修改資料（不建議使用）

```
1 # fix() 會顯示函數內完整資料，並可點選指定欄位作修改
2 fix(camera)
3
4 # 這個方法非常快就可以更改資料的內容
5 # 但一般都會建議不要修改原始資料（非必要少用）
```

從data.frame中拿掉遺失值

```
1 # 假設有一個名為dfm的data.frame格式資料
2 dfm
3     V1  V2
4 1 0.9 0.8
5 2  NA 2.7
6 3 0.3 0.1
7 4 2.4  NA
8 5 5.2 2.0
9
10 #有遺失值的資料會難以做計算 （範例已累積加總為例）
11 cumsum(dfm)
12     V1  V2
13 1 0.9 0.8
14 2  NA 3.5
15 3  NA 3.6
16 4  NA  NA
17 5  NA  NA
18
19 #加上na.omit() 會拿掉含遺失值的整列資料唷！！
20 cumsum(na.omit(dfm))
```

```
21      V1  V2
22  1 0.9 0.8
23  3 1.2 0.9
24  5 6.4 2.9
25
26 #如果只想不要計算有遺失值的欄，也可以分開計算
27 cumsum(na.omit(dfm[,1]))
28 [1] 0.9 1.2 3.6 8.8
29 cumsum(na.omit(dfm[,2]))
30 [1] 0.8 3.5 3.6 5.6
```

## 字串處理

### substr 擷取部分字串

```
1 x<-'My name is Emma.'
2 substr(x,3,10)  #取第3~10之字元
3 [1] " name is"
```

### substring 擷取部分字串

```
1 x.strings <- c(
2   "abcdefghij",
3   "ABCDEFGHIJ",
4   "1234567890"
5 )
6 substr(x.strings, 1:4, 8)
7 [1] "abcdefgh" "BCDEFGH"  "345678"
8 substring(x.strings, 1:4, 8)
9 [1] "abcdefgh" "BCDEFGH"  "345678"   "defgh"
```

### stringr 套件

可搭配正規表達法使用

### str\_split 資料剖析

```
1 fruits <- c( "apples and oranges and pears and bananas", "pineapples and man
2 gos and guavas" )
3 str_split(fruits, " and ") #可搭配正規表示法
4 [[1]]
5 [1] "apples"  "oranges" "pears"   "bananas"
6
```

```
7 [[2]]
8 [1] "pineapples" "mangos"      "guavas"
```

## str\_mach & str\_match\_all 查找匹配字符

```
1 Name <- c("王小名","李美美","董小丁和方小姬","王阿枝","陳小雲")
2 str_match(Name,".小.") #傳回第一個匹配值  #傳回matrix
3      [,1]
4 [1,] "王小名"
5 [2,] NA
6 [3,] "董小丁"
7 [4,] NA
8 [5,] "陳小雲"
9
10 str_match_all(Name,".小.") #傳回所有匹配值  #傳回list
11 [[1]]
12      [,1]
13 [1,] "王小名"
14
15 [[2]]
16      [,1]
17
18 [[3]]
19      [,1]
20 [1,] "董小丁"
21 [2,] "方小姬"
22
23 [[4]]
24      [,1]
25
26 [[5]]
27      [,1]
28 [1,] "陳小雲"
```

## str\_extract & str\_extract\_all 提取匹配字符

```
1 str_extract(Name,".小.") #返回第一個匹配值  #返回向量
2 [1] "王小名" NA      "董小丁" NA      "陳小雲"
3
4 str_extract_all(Name,".小.") #返回所有匹配值  #返回list
5 [[1]]
```

```
6  [1] "王小名"
7
8  [[2]]
9  character(0)
10
11 [[3]]
12 [1] "董小丁" "方小姬"
13
14 [[4]]
15 character(0)
16
17 [[5]]
18 [1] "陳小雲"
```

## str\_replace & str\_replace\_all 取代

```
1  > str_replace(Name, "小", "阿")
2  [1] "王阿名"          "李美美"          "董阿丁和方小姬" "王阿枝"          "陳阿雲"

3  > str_replace_all(Name, "小", "阿")
4  [1] "王阿名"          "李美美"          "董阿丁和方阿姬" "王阿枝"          "陳阿雲"
```

Thank you!