
第 15 章：資料處理

15: Data Management

統計分析的第一步就是資料清理與資料處理，通常可佔據整體統計分析的 70%-80% 時間。資料清理包含同一個體重覆輸入觀測值，錯誤輸入觀測值，缺失值檢視與處理，連續變數語類別變數的區分，選取適當的個體，資料排序與分組等等，R 有許多資料處理的函式與套件，但是 R 資料儲存有許多不同類型與格式，包含 向量，矩陣，陣列，列表等，且 R 個別函式與套件，統計模型接受的資料類型與格式各不相同，因此讓使用者容易混淆而造成困擾。R 提供 `apply()` 系列函式，功能強大，但因接受的資料類型與格式各不相同，讓初學者相當困惑，參見 表 15.1. 套件 `plyr`, `dplyr`, `reshape`, `reshape2`, `tidyverse` 等，提供較 R 容易理解的資料處理函式，但是計算速度有差異，且經常更新與改名，也讓初學者相當困惑。

表 15.1: R Basic `apply()` Function

Input object	Output Object			
	array	data frame	list	nothing
array	apply	.	.	.
data frame	.	aggregate	by	.
list	sapply	.	lapply	.
replicates	replicate	.	replicate	.
function	arguments	mapply	mapply	.

資料物件通常以 矩陣, 陣列, 列表 與 資料框架 型式儲純數值, 統計資料物件內個別維度的成分 (變數) 通常是向量物件, 在統計分析資料時, 必須對資料物件內各別的成分 (變數) 進行運算, 在 R 中, 許多函式通常利用向量做運算, 但在統計分析資料時, 更多時候必須對許多向量使用同一函式做重複運算, 例如, 資料框架內的每一變數, 欄向量 (行向量), 計算平均值與變異數. 在統計分析資料時, 有時候會針對一個變數, 依照另外一個類別變數不同的類別水準, 分別使用同一函式做運算, 例如. 根據不同治療的組別, 分別計算血壓的樣本平均值與樣本變異數. 在上述的情形, 許多程式語言常常須用到迴圈 (loop), 或是自行寫作向量操作函式. 在 R 中, 回圈演算法的執行非常無效率, 應儘量避免, 因此, 在 R 中許多函式直接利用向量做運算, 且在 R 中有些函式, 例如, `apply()`, `lapply()`, `sapply()`, `tapply()`, `mapply()`, `Vectorize()`, `by()`, `sweep()`, `aggregate()`, `replicate()` 等, 可以對不同向量進行重複的運算, 如此避免以回圈演算法進行重複的運算, 可以更有效率的執行類似迴圈的函式指令.

表 15.2: 常見資料物件操作函式

函式	輸入	輸出
<code>by</code>	<code>df</code>	<code>list</code>
<code>aggregate</code>	<code>df</code>	<code>df</code>
<code>replicate</code>	<code>r</code>	<code>array, list</code>
<code>apply</code>	<code>matrix, array, df</code>	<code>vector, array, list</code>
<code>dendrapply</code>	<code>dendrogram</code>	<code>dendrogram</code>
<code>eapply</code>	<code>environment</code>	<code>list</code>
<code>lapply</code>	<code>list, df, vector</code>	<code>list</code>
<code>mapply</code>	<code>list, df, vectors</code>	<code>vector, matrix, list</code>
<code>rapply</code>	<code>list</code>	<code>vector, list,</code>
<code>sapply</code>	<code>list, df, vector</code>	<code>vector, matrix, list</code>
<code>tapply</code>	<code>df, categories</code>	<code>array, list</code>
<code>vapply</code>	<code>list, df, vector</code>	<code>vector, matrix, list</code>

15.1 資料物件操作函式: `apply()` 函式

使用函式 `apply()` 對一個 資料框架 (data frame), 矩陣 (matrix) 陣列物件 (array) 內之任一維度的邊際 (margins of an array) 執行同一個函式運算, 並回傳一向量, 陣列或列表. 函式指令為

```
1 > apply(X, MARGIN, FUN, ...)
```

其中的主要引數分別為:

- `X` 為所要執行之陣列物件, 若資料物件 `X` 為矩陣 或 資料框架, 則函式 `apply()` 會將該 矩陣 或 資料框架 轉換成 2-維度 陣列 (array).
- `MARGIN` 為陣列物件 `X` 的邊際維度序號.
 - `MARGIN = 1` 會對陣列 `X` 的第 1 個維度 (row) 執行函式, 例如, 對 矩陣 或 資料框架 的每一列 (row) 做運算.
 - `MARGIN = 2` 會對陣列 `X` 的第 2 個維度 (column) 執行函式, 例如, 對 矩陣 或 資料框架 的每一欄 (行, column) 做運算等等.
 - `MARGIN = k` 會對陣列 `X` 的第 k 個維度 執行函式.
- `FUN` 為執行同一函式的名稱.

例 15.1. 安眠藥臨床試驗

醫學中心一位醫師針對一新開發之 安眠藥 T 的藥效進行一研究, 這位醫師想證明 安眠藥 T 比 安慰劑 C 可以有效地治療失眠患者. 因此這位醫師進行一個先導研究, 共 10 位失眠患者同意參予研究. 這 10 位受試者, 在第 1 天給予 安慰劑 C , 並測量第 1 天當夜的連續睡眠時間. 然後第 2 天給予安眠藥 T , 再次測量第 2 天當夜的連續睡眠時間. 研究的主要反應變數為受試者的連續睡眠的時間, 若其第 2 天睡眠的時間遠大於第 1 天連續睡眠的時間, 表示新開發之安眠藥 T 可以有效地治療失眠患者.

表 15.3 顯示研究結果 (單位: 小時). 研究的結果是否有足夠的證據力支持研究者的想法?

表 15.3: 安眠藥的藥效臨床試驗研究: 服藥後的睡眠時間

受試者	1	2	3	4	5	6	7	8	9	10
C	2.0	3.6	4.0	3.9	4.0	1.3	3.5	3.9	3.8	1.5
T	5.0	4.8	2.5	6.2	3.6	7.3	4.9	2.6	3.7	7.4

```

1 > ## apply() + matrix
2 > x.mat <- matrix(1:12, nrow = 4, byrow = T)
3 > x.mat
4      [,1] [,2] [,3]
5 [1,]    1    2    3
6 [2,]    4    5    6
7 [3,]    7    8    9
8 [4,]   10   11   12
9 > apply(X = x.mat, MARGIN = 1, FUN = mean)
10 [1] 2 5 8 11
11 > apply(x.mat, 2, mean)
12 [1] 5.5 6.5 7.5
13 > #
14 > ## apply() + data.frame
15 > setwd("C:/RData/")
16 > sleep <- read.table(file = "sleepwidePaired.csv",
17                         header = TRUE, sep = ",",
18                         row.names = NULL)
19 > head(sleep)
20   id Tdrug Cdrug
21 1  1   5.0   2.0
22 2  2   4.8   3.6
23 3  3   2.5   4.0
24 4  4   6.2   3.9
25 5  5   3.6   4.0
26 6  6   7.3   1.3
27 > apply(sleep[, 2:3], 1, mean)
28 [1] 3.50 4.20 3.25 5.05 3.80 4.30 4.20 3.25 3.75 4.45
29 > apply(sleep[, 2:3], 2, mean)
30 Tdrug Cdrug
31 4.80 3.15
32 > #
33 > ## more apply() + matrix
34 > set.seed(100)
35 > x.mat <- matrix(rnorm(1000), 200, 5)
36 > apply(x.mat, 2, quantile, probs = c(0.25, 0.75))
37      [,1]      [,2]      [,3]      [,4]      [,5]

```

```
38 25% -0.5056721 -0.6637388 -0.6670268 -0.7151312 -0.6835946
39 75%  0.5273765  0.6685950  0.6429354  0.7748576  0.8415970
40 >
41 > y.mat <- matrix(rnorm(10000), 2000, 5)
42 > apply(y.mat, 2, quantile, probs = c(0.25, 0.75))
43      [,1]      [,2]      [,3]      [,4]      [,5]
44 25% -0.7058193 -0.6644742 -0.6559814 -0.6828271 -0.6520564
45 75%  0.6618302  0.6722488  0.6733151  0.6641776  0.6773927
46 > #
47 > z.mat <- matrix(rnorm(10000), 5, 2000)
48 > apply(z.mat, 1, quantile, probs = c(0.25, 0.75))
49      [,1]      [,2]      [,3]      [,4]      [,5]
50 25% -0.6717053 -0.6963323 -0.6643946 -0.7375793 -0.6294552
51 75%  0.6154957  0.6828310  0.6504301  0.6262656  0.7088369
52 >
53 >
54 > ## apply() + array
55 > b.arr <- array(1:24, dim = c(4, 3, 2),
56                  dimnames = list(c("x1", "x2", "x3", "x4"),
57                               c("y1", "y2", "y3"),
58                               c("z1", "z2")))
59 > b.arr
60 , , z1
61   y1 y2 y3
62 x1  1  5  9
63 x2  2  6 10
64 x3  3  7 11
65 x4  4  8 12
66
67 , , z2
68   y1 y2 y3
69 x1 13 17 21
70 x2 14 18 22
71 x3 15 19 23
72 x4 16 20 24
73
74 > #
75 > b.arr[1, , ]
76   z1 z2
77 y1  1 13
78 y2  5 17
79 y3  9 21
80 > mean(b.arr[1, , ])
81 [1] 11
82 > #
83 > apply(b.arr, 1, mean) # average by x
84 x1 x2 x3 x4
85 11 12 13 14
86 > rowMeans(b.arr, dims = 1)
87 x1 x2 x3 x4
```

```
88 11 12 13 14
89 > #
90 > apply(b.arr, c(2, 3), mean) # average over x
91     z1   z2
92 y1  2.5 14.5
93 y2  6.5 18.5
94 y3 10.5 22.5
95 > colMeans(b.arr, dims = 1)
96     z1   z2
97 y1  2.5 14.5
98 y2  6.5 18.5
99 y3 10.5 22.5
100 > #
101 > apply(b.arr, c(1, 2), mean)
102    y1 y2 y3
103 x1  7 11 15
104 x2  8 12 16
105 x3  9 13 17
106 x4 10 14 18
107 > rowMeans(b.arr, dims = 2)
108    y1 y2 y3
109 x1  7 11 15
110 x2  8 12 16
111 x3  9 13 17
112 x4 10 14 18
113 > #
114 > # array average
115 > a.arr <- array(rnorm(2*2*10), c(2, 2, 10))
116 > apply(a.arr, c(1, 2), mean)
117     [,1]      [,2]
118 [1,] 0.2426878 -0.41321519
119 [2,] 0.1213732  0.09303838
120 > rowMeans(a.arr, dims = 2)
121     [,1]      [,2]
122 [1,] 0.2426878 -0.41321519
123 [2,] 0.1213732  0.09303838
124 > #
125 > apply(a.arr, c(1, 3), mean)
126     [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
127 [1,] 0.4570282 -0.4108178 -0.5218928 -0.4031859 -0.1281254 -0.2967553
128 [2,] 0.3604596  0.1664159  0.2622479  0.2454029 -0.4436532  0.7604331
129     [,7]      [,8]      [,9]      [,10]
130 [1,] -0.1859059 -0.5653128  0.5192811  0.6830499
131 [2,] -1.0262100 -0.5717037  1.1088832  0.2097824
132 > rowMeans(a.arr, dims = 1)
133 [1] -0.08526367  0.10720581
134 > colMeans(a.arr, dims = 1)
135     [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
136 [1,] 0.2399494  0.1379897  0.7800541  0.1544253 -0.3909275  0.2618471
137 [2,] 0.5775384 -0.3823917 -1.0396990 -0.3122082 -0.1808511  0.2018306
```

```
138      [,7]      [,8]      [,9]      [,10]
139 [1,] -0.1266432 -0.4104021 0.4224371 0.7515755
140 [2,] -1.0854728 -0.7266143 1.2057272 0.1412568
141 > #
142 > colMeans(a.arr, dims = 2)
143 [1] 0.40874388 -0.12220098 -0.12982247 -0.07889147 -0.28588930 0.23183887
144 [7] -0.60605798 -0.56850821 0.81408215 0.44641616
```

15.2 資料物件操作函式: lapply() 與 sapply()

使用函式 `lapply()` 對一個列表物件 (list) 內的每一成分 (component), 執行同一個函式做運算, 並回傳一個列表物件 (list), 且此回傳的列表物件長度與原有的列表物件長度相同.

函式 `sapply()` 是相對於函式 `lapply()` 較容易使用的替代函式, 使用函式 `sapply()` 對一個列表物件 (list) 內的每一成分 (component), 執行同一個函式運算, 並且回傳一向量或矩陣. 例如, 使用函式指令 `sapply(X.df)` 可以對資料框架 (列表物件) `X.df` 執行同一個函式做運算, 是指對每一個變數 (欄位), 做運算. 函式 `vapply()` 是類似於函式 `sapply()`. 函式指令為

```
1 > lapply(X, FUN, ...)
2 > sapply(X, FUN, ..., simplify = TRUE, USE.NAMES = TRUE)
3 > vapply(X, FUN, FUN.VALUE, ..., USE.NAMES = TRUE)
```

其中的主要引數分別為:

- `X` 為所要執行之 列表物件 或 資料框架.
- 對列表每一個成分執行同一函式.
- `FUN` 為執行同一函式的名稱.
- `FUN.VALUE` 為 執行函數 `FUN` 回傳的物件名稱.
- `USE.NAMES = TRUE`: 若 `X` 物件成分已有命名, 內設直接使用 `X` 命名.

```
1 > ## lapply(), sapply() + list
2 > set.seed(100)
3 > x.list <- list(a = 1:5, b = rnorm(10, mean = 0, sd = 1))
4 > lapply(x.list, mean)
5 $a
6 [1] 3
7
8 $b
9 [1] -0.01795716
10 > #
11 > sapply(x.list, mean)
12      a          b
13  3.00000000 -0.01795716
14 > #
15 > x.list <- list(a = 11:15, b = rnorm(10), c = rnorm(1000, mean = 5, sd = 2))
16 > lapply(x.list, mean)
17 $a
18 [1] 13
19
20 $b
21 [1] 0.2336915
22
23 $c
24 [1] 5.033755
25 > #
26 > sapply(x.list, mean)
27      a          b          c
28 13.0000000  0.2336915  5.0337552
29 > #
30 > x <- 1:2
31 > lapply(x, runif)
32 [[1]]
33 [1] 0.5156047
34
35 [[2]]
36 [1] 0.4649941 0.5249516
37 > #
38 > sapply(x, runif)
39 [[1]]
40 [1] 0.1502433
41
42 [[2]]
43 [1] 0.7110009 0.3702982
44 > #
45 > x <- 1:2
46 > lapply(x, runif, min = 0, max = 10)
47 [[1]]
48 [1] 4.533344
49
50 [[2]]
```

```
51 [1] 6.179135 2.836679
52 > #
53 > sapply(x, runif, min = 0, max = 10)
54 [[1]]
55 [1] 2.069829
56
57 [[2]]
58 [1] 6.077585 6.662716
59 > #
60 > x.list <- list(a = matrix(1:4, 2, 2), b = matrix(1:6, 3, 2))
61 > x.list
62 $a
63     [,1] [,2]
64 [1,]    1    3
65 [2,]    2    4
66
67 $b
68     [,1] [,2]
69 [1,]    1    4
70 [2,]    2    5
71 [3,]    3    6
72 > #
73 > lapply(x.list, function(elt) elt[,1])
74 $a
75 [1] 1 2
76
77 $b
78 [1] 1 2 3
79 > #
80 > sapply(x.list, function(elt) elt[,1])
81 $a
82 [1] 1 2
83
84 $b
85 [1] 1 2 3
86 > #
87 > c.list <- list(a = 1:20,
88                     beta = exp(-2:2),
89                     logic = c(TRUE, FALSE, FALSE, TRUE, FALSE))
90 > lapply(c.list, mean)
91 $a
92 [1] 10.5
93
94 $beta
95 [1] 2.322111
96
97 $logic
98 [1] 0.4
99 > #
100 > sapply(c.list, mean)
```

```
101      a      beta      logic
102 10.500000 2.322111 0.400000
103 > #
104 > # sapply() + list = simplify
105 > sapply(c.list, mean, simplify = FALSE)
106 $a
107 [1] 10.5
108
109 $beta
110 [1] 2.322111
111
112 $logic
113 [1] 0.4
114 > #
115 > sapply(c.list, mean, simplify = TRUE)
116      a      beta      logic
117 10.500000 2.322111 0.400000
118 > sapply(c.list, mean)
119      a      beta      logic
120 10.500000 2.322111 0.400000
121 > #
122 > ## lapply() + data.frame
123 > setwd("C:/RData/")
124 > sleep <- read.table(file = "sleepwidePaired.csv",
125                         header = TRUE, sep = ",", dec = ".",
126                         row.names = NULL)
127 > #
128 > lapply(sleep, mean)
129 $id
130 [1] 5.5
131
132 $Tdrug
133 [1] 4.8
134
135 $Cdrug
136 [1] 3.15
137 > #
138 > sapply(sleep, mean)
139   id Tdrug Cdrug
140 5.50 4.80 3.15
141 > #
142 > ## vapply() + list + FUN
143 > d.list <- sapply(3:5, seq) # list of vectors
144 > d.list
145 [[1]]
146 [1] 1 2 3
147
148 [[2]]
149 [1] 1 2 3 4
150
```

```
151 [[3]]  
152 [1] 1 2 3 4 5  
153 > #  
154 > sapply(d.list, fivenum)  
155      [,1] [,2] [,3]  
156 [1,] 1.0 1.0 1  
157 [2,] 1.5 1.5 2  
158 [3,] 2.0 2.5 3  
159 [4,] 2.5 3.5 4  
160 [5,] 3.0 4.0 5  
161 > #  
162 > vapply(d.list, fivenum,  
163             c(Min = 0, "Q1" = 0,  
164              Median = 0, "3rd Qu." = 0, "Max" = 0))  
165      [,1] [,2] [,3]  
166 Min     1.0 1.0 1  
167 Q1     1.5 1.5 2  
168 Median 2.0 2.5 3  
169 3rd Qu. 2.5 3.5 4  
170 Max     3.0 4.0 5
```

15.3 資料物件操作函式: tapply()

使用函式 `tapply()` 可以對一參差不齊的陣列 (ragged array), 依照類別變數的類別水準或組別的指標值, 分別使用同一函式做運算. 例如, 計算分組平均值. 函式指令為

```
1 > tapply(X, INDEX, FUN = NULL, ..., simplify = TRUE)
```

其中的主要引數分別為:

- `X` 是一個原型物件 (atomic object), 通常向量.
- `INDEX` 通常是一個 類別變數, 包含 類別水準或組別 所形的 因子向量 (或 陣列), 長度與 `X` 相同.

例 15.2. 抽菸, 膽固醇與血壓調查研究

一個早期大規模的抽樣調查, 以問卷方式調查 抽菸, 膽固醇 與 血壓 的關聯性研究, 隨機抽取了近 1.5 萬位成年人進行問卷訪問, 資料變數在表 ??, 部分個體資料在

表 ??, 檔案在 BPsurvey.csv. 研究主要描述一般族群的健康狀況, 探討 疾病 與 發生原因 之關聯性, 例如, 分析 抽菸, 膽固醇 與 血壓 的關聯性.

```

1 > ## tapply() + data.frame
2 > setwd("C:/RData/")
3 > bp.df <- read.table(file = "BPsurvey.csv",
4                         header = TRUE, sep = ",",
5                         dec = ".",
6                         row.names = NULL)
6 > dim(bp.df)
7 [1] 15643     10
8 > head(bp.df)
9   id sex race smoke age wt ht sbp dbp chol
10 1  3    1    1      1 21 81 179 120  67 268
11 2  4    0    1      1 32 61 162 126  86 160
12 3  9    0    1      2 48 68 157 131  73 236
13 4 10    1    1      1 35 92 177 130  82 225
14 5 11    1    1      2 48 70 168 120  70 260
15 6 19    1    2      3 44 86 178 133  85 187
16 >#
17 > tapply(bp.df$wt, bp.df$sex, mean)
18
19 0          1
19 70.56511 80.18642
20 > tapply(bp.df$wt, bp.df$sex, mean, simplify = FALSE)
21 $'0'
22 [1] 70.56511
23
24 $'1'
25 [1] 80.18642
26 > #
27 > # sapply() by pass arguments into tapply()
28 > sapply(bp.df[, 5:10], FUN = tapply, bp.df$sex, mean)
29
30      age      wt      ht      sbp      dbp      chol
30 0 47.64641 70.56511 160.1776 123.7655 72.19050 207.8440
31 1 48.44739 80.18642 173.3462 128.4596 76.79267 203.6927

```

15.4 資料物件操作函式: mapply()

函式 `mapply()` 是函式 `apply()` 的多變量型式, 對一個列表物件 (list) 內的每一成分 (component), 執行同一個函式運算, `mapply()` 回傳一向量, 矩陣 或 列表.

函式 `mapply()` 是函式 `Map` 的簡化. 函式 `mapply()` 指令為

```

1 > mapply(FUN, ...,
2           MoreArgs = NULL, SIMPLIFY = TRUE,
2           USE.NAMES = TRUE)

```

其中的主要引數分別為:

- `FUN()` 為執行同一函式的名稱, 且 `FUN()` 本身可以含有引數. `FUN()` 的引數以 `FUN(ARG1 = arg1, ARG2 = arg2, ...)` 等放入原有函式 `mapply()` 之中.
- 執行 `FUN` 函式的運算時, 使用 `ARG1 = arg1` 的第 1 元素, `ARG2 = arg2` 的第 2 元素, ..., 等等, 作為 `FUN` 的引數.
- `MoreArgs = list()`, 為執行 `FUN` 函式所需要用列表物件, 且引數可以循環使用.

```
1 > ## mapply()
2 > list(rep(1, 3), rep(2, 2), rep(3, 1))
3 [[1]]
4 [1] 1 1 1
5
6 [[2]]
7 [1] 2 2
8
9 [[3]]
10 [1] 3
11 > #
12 > mapply(rep, 1:3, 3:1)
13 [[1]]
14 [1] 1 1 1
15
16 [[2]]
17 [1] 2 2
18
19 [[3]]
20 [1] 3
21 > #
22 > mapply(rep, times = 1:3, x = 3:1)
23 [[1]]
24 [1] 3
25
26 [[2]]
27 [1] 2 2
28
29 [[3]]
30 [1] 1 1 1
31 > #
32 > mapply(rep, times = 1:3, MoreArgs = list(x = 7))
33 [[1]]
```

```
34 [1] 7
35
36 [[2]]
37 [1] 7 7
38
39 [[3]]
40 [1] 7 7 7
41 > ##
42 > noise <- function(n, mean, sd) {
43   rnorm(n, mean, sd)
44 }
45 > set.seed(100)
46 > noise(n = 3, mean = 1, sd = 2)
47 [1] -0.004384701 1.263062331 0.842165820
48 > noise(n = 1:3, mean = 1:3, sd = 2)
49 [1] 2.773570 2.233943 3.637260
50 > mapply(noise, n = 1:3, mean = 1:3, sd = 2)
51 [[1]]
52 [1] -0.1635814
53
54 [[2]]
55 [1] 3.4290654 0.3494811
56
57 [[3]]
58 [1] 2.280276 3.179772 3.192549
59 > #
60 > list(noise(1, 1, 2), noise(2, 2, 2),
61 +       noise(3, 3, 2), noise(4, 4, 2),
62 +       noise(5, 5, 2))
63 [[1]]
64 [1] 0.5967321
65
66 [[2]]
67 [1] 3.479681 2.246759
68
69 [[3]]
70 [1] 2.941367 2.222292 4.021713
71
72 [[4]]
73 [1] 2.172372 8.620594 3.123820 5.528121
74
75 [[5]]
76 [1] 5.523923 6.546809 3.371242 4.123099 3.559557
```

15.5 資料物件操作函式: Vectorize()

函式 `Vectorize()` 是 `mapply()` 的的包裝封套, 回傳一新函式, 有如呼叫 `mapply()` 函式. 函式指令為

```
1 > Vectorize(FUN, vectorize.args = arg.names, SIMPLIFY = TRUE,  
2           USE.NAMES = TRUE)
```

其中的主要引數分別為:

- `FUN()` 為執行同一函式的名稱, 且 `FUN()` 本身可以含有引數. `FUN()` 必須是正式存在的函式名.
- `vectorize.args` 為文字向量, 須將 `FUN()` 引數轉換為文字向量.

```
1 > ## Vectorize()  
2 > mapply(rep, 1:3, 3:1)  
3 [[1]]  
4 [1] 1 1 1  
5  
6 [[2]]  
7 [1] 2 2  
8  
9 [[3]]  
10 [1] 3  
11 > #  
12 > vrep <- Vectorize(rep)  
13 > vrep(1:3, 3:1)  
14 [1] 1 1 1 2 2 3  
15 > #  
16 > mapply(rep, times = 1:3, x = 3:1)  
17 [[1]]  
18 [1] 3  
19  
20 [[2]]  
21 [1] 2 2  
22  
23 [[3]]  
24 [1] 1 1 1  
25 > #  
26 > vrep(times = 1:3, x = 3:1)  
27 [1] 3 2 2 1 1 1  
28 > #  
29 > vrep <- Vectorize(rep.int)  
30 > vrep(times = 1:3, x = 45)
```

```

31 [[1]]
32 [1] 45
33
34 [[2]]
35 [1] 45 45
36
37 [[3]]
38 [1] 45 45 45
39 > #
40 > vrep <- Vectorize(rep.int, "times")
41 > vrep(times = 1:3, x = 45)
42 [[1]]
43 [1] 45
44
45 [[2]]
46 [1] 45 45
47
48 [[3]]
49 [1] 45 45 45
50 > #
51 > # error
52 > # not a formal function name
53 > vrep <- Vectorize(rep, "times")
54 Error in Vectorize(rep, "times") :
55   must specify names of formal arguments for 'vectorize'
56 > help(rep) # rep.int as a formal function mane

```

15.6 資料物件操作函式: by()

使用函式 `by()` 可將 資料框架 依不同的組別因子，計算統計量或執行同一個函式做運算。函式指令為

```
1 > by(data, INDICES, FUN, ..., simplify = TRUE)
```

其中的主要引數分別為：

- `data` 為 資料框架 或 矩陣。
- `INDICES` 為因子向量代表類別變數的組別，其長度與資料框架變數的長度 (`nrows(data)`) 相同。

- `FUN()` 是指所要計算的統計量或執行同一函式的名稱, `FUN()` 的引數必須是資料框架 或 矩陣 進行操作的函式.
- `by()` 回傳列表物件 (list).

```
1 > ## by()
2 > setwd("C:/RData/")
3 > bp.df <- read.table(file = "BPsurvey.csv",
4                         header = TRUE, sep = ",", dec = ".",
5                         row.names = NULL)
6 > head(bp.df)
7   id sex race smoke age wt ht sbp dbp chol
8 1 3   1    1     1   21 81 179 120  67 268
9 2 4   0    1     1   32 61 162 126  86 160
10 3 9   0    1     2   48 68 157 131  73 236
11 4 10  1    1     1   35 92 177 130  82 225
12 5 11  1    1     2   48 70 168 120  70 260
13 6 19  1    2     3   44 86 178 133  85 187
14 > #
15 > by(data = bp.df[, c(8:10)], INDICES = bp.df$sex, FUN = summary)
16 bp.df$sex: 0
17   sbp          dbp          chol
18 Min. :69.0  Min. :20.00  Min. : 74.0
19 1st Qu.:108.0 1st Qu.:65.00 1st Qu.:175.0
20 Median :119.0 Median :72.00 Median :203.0
21 Mean   :123.8 Mean   :72.19 Mean   :207.8
22 3rd Qu.:136.0 3rd Qu.:78.00 3rd Qu.:235.0
23 Max.  :237.0 Max.  :142.00 Max.  :676.0
24 -----
25 bp.df$sex: 1
26   sbp          dbp          chol
27 Min. :80.0  Min. :40.00  Min. : 59.0
28 1st Qu.:116.0 1st Qu.:70.00 1st Qu.:174.0
29 Median :125.0 Median :76.00 Median :201.0
30 Mean   :128.5 Mean   :76.79 Mean   :203.7
31 3rd Qu.:137.0 3rd Qu.:83.00 3rd Qu.:229.0
32 Max.  :244.0 Max.  :134.00 Max.  :702.0
33 > # by(bp.df[, c(5:10)], bp.df$sex, summary)
34 > #
35 > ##
36 > tmp <- by(bp.df, bp.df[, "sex"],
37               function(x) lm(sbp ~ age, data = x))
38 > sapply(tmp, coef)
39
40 (Intercept) 88.5392872 105.5821478
41 age          0.7393261  0.4722124
42 > #
43 > # NOT WORK ANY MORE
44 > # BECAUSE MEAN FUN DID NOT TAKE df OBJECT
```

```

45 > by(data = bp.df[, c(8:10)], bp.df$sex, mean)
46 bp.df$sex: 0
47 [1] NA
48 -----
49 bp.df$sex: 1
50 [1] NA
51 Warning messages:
52 .....

```

15.7 資料物件操作函式: aggregate()

函式 `aggregate()` 可將 資料框架 分成不同的組別, 分別計算統計量或執行同一個函式做運算. 函式 `aggregate()` 與函式 `by()` 有類似功能. 函式 `aggregate()` 指令為

```

1 > ## S3 method for class 'data.frame'
2 > aggregate(x, by, FUN, ..., simplify = TRUE, drop = TRUE)
3 > ## S3 method for class 'formula'
4 > aggregate(formula, data, FUN, ...,
5   subset, na.action = na.omit)

```

其中的主要引數分別為:

- `x` 為資料框架.
- `formula` 為統計模型.
- `by = f.list`, 是 列表物件 (list), 長度與 資料框架 `X.df` 相同的 因子變數.
可用來分組分別計算統計量.
- `FUN` 是指所要計算的統計量或執行同一個函式名稱.
- `aggregate()` 中 `FUN()` 的限制較 `by()` 少, 但 `aggregate()` 回傳 資料框架 (`data.frame`) 與 `by()` 回傳列表 (list) 不同.

```

1 > ## aggregate()
2 > aggregate(bp.df[, 5:10], by = list(SEX = bp.df$sex), mean)
3 SEX      age      wt      ht      sbp      dbp      chol
4 1    0 47.64641 70.56511 160.1776 123.7655 72.19050 207.8440

```

```
5 2 1 48.44739 80.18642 173.3462 128.4596 76.79267 203.6927
6 > class(aggregate(bp.df[, 5:10], by = list(SEX = bp.df$sex), mean))
7 [1] "data.frame"
8 > #
9 > # transpose a data.frame
10 > t(aggregate(bp.df[, 9:10], by = list(SEX = bp.df$sex), summary))
11          [,1]   [,2]
12 SEX          0.00  1.00
13 dbp.Min.    20.00 40.00
14 dbp.1st Qu. 65.00 70.00
15 dbp.Median  72.00 76.00
16 dbp.Mean    72.19 76.79
17 dbp.3rd Qu. 78.00 83.00
18 dbp.Max.   142.00 134.00
19 chol.Min.   74.00 59.00
20 chol.1st Qu. 175.00 174.00
21 chol.Median 203.00 201.00
22 chol.Mean   207.80 203.70
23 chol.3rd Qu. 235.00 229.00
24 chol.Max.  676.00 702.00
25 > #
26 > ## compare aggregate() and sapply()
27 > aggregate(bp.df[, 8:10], by = list(SEX = bp.df$sex), mean)
28   SEX      sbp      dbp      chol
29 1  0 123.7655 72.19050 207.8440
30 2  1 128.4596 76.79267 203.6927
31 > sapply(bp.df[, c(8:10)], tapply, bp.df$sex, mean)
32      sbp      dbp      chol
33 0 123.7655 72.19050 207.8440
34 1 128.4596 76.79267 203.6927
```

15.8 資料物件操作函式: sweep()

使用函式 `sweep()` 可以對一個陣列 物件, 針對邊際維度, 計算一個統計量或執行同一個函式做運算. 函式指令為

```
1 > sweep(x, MARGIN, STATS, FUN = "-", check.margin = TRUE, ...)
```

其中的主要引數分別為:

- `X` 為一個 陣列 或 矩陣 物件.
- `FUN` 是指所要計算的統計量或執行同一函式的名稱.

- 函式 `sweep()` 自動內設為 `(FUN = "-")` (減法運算), 若執行特定符號運算, 例如, 此特定符號運算為 `+` (加號), 則符號須加用雙引號成為 `"+"`.

```

1 > ## sweep()
2 > b.data <- as.data.frame(matrix(c(1:24), nrow = 6, byrow = T))
3 > (b.mean <- apply(b.data, 2, mean))
4 V1 V2 V3 V4
5 11 12 13 14
6 > sweep(b.data, 2, b.mean)
7     V1 V2 V3 V4
8 1 -10 -10 -10
9 2 -6 -6 -6
10 3 -2 -2 -2
11 4 2 2 2
12 5 6 6 6
13 6 10 10 10
14 > #
15 > (b.mean <- apply(b.data, 1, mean))
16 [1] 2.5 6.5 10.5 14.5 18.5 22.5
17 > sweep(b.data, 1, b.mean)
18     V1 V2 V3 V4
19 1 -1.5 -0.5 0.5 1.5
20 2 -1.5 -0.5 0.5 1.5
21 3 -1.5 -0.5 0.5 1.5
22 4 -1.5 -0.5 0.5 1.5
23 5 -1.5 -0.5 0.5 1.5
24 6 -1.5 -0.5 0.5 1.5

```

15.9 資料物件操作函式: replicate()

使用函式 `replicate()` 與函式 `sapply()` 類似, 為函式 `sapply()` 的包裝封套, 可以重覆執行同一函式做運算. 函式指令為

```
1 > replicate(n, expr, simplify = "array")
```

其中的主要引數

- `n` 所要重覆執行的次數.
- `expr` 為所要重覆執行的同一運算式.
- 函式 `replicate()` 通常是執行產生隨機變數.

```
1 > ## replicate random variables from exponential distribution
2 > replicate(5, rexp(10))
3      [,1]     [,2]     [,3]     [,4]     [,5]
4 [1,] 0.06436724 1.49382223 0.1169804 0.39905662 0.0312094139
5 [2,] 4.05790204 0.08366908 2.2931382 0.81557964 0.7972878694
6 [3,] 0.52927286 3.49740175 0.2327129 0.75613704 0.4220017293
7 [4,] 0.08287350 0.35182542 1.3666231 0.74424509 1.1743398280
8 [5,] 0.14477947 0.53206901 1.1084216 0.33239908 0.8897708109
9 [6,] 1.57175982 1.86628848 4.9939808 0.38849803 2.0421579144
10 [7,] 0.02001488 0.16144714 1.5584840 1.37843723 0.2155170352
11 [8,] 0.52199820 1.34512130 0.4799255 0.47659392 0.3325432939
12 [9,] 0.08017182 0.55292033 2.6603417 1.65020506 1.9962155669
13 [10,] 0.23903690 1.18049660 1.1706547 0.04541023 0.0007613422
14 > replicate(5, mean(rexp(10)))
15 [1] 1.0176443 0.7254602 1.0603147 1.0598029 1.4094343
16 > #
17 > hist(replicate(100, mean(rexp(10))))
```

15.10 資料框架分割與合併

許多時候需要對資料框架分割 (split) 與合併 (merge), 可以分別使用函式 `split()` 與函式 `merge()` 進行分割與合併.

15.10.1 資料框架分割函式: `split()`

使用函式 `split()` 可以對資料框架分割, 傳回列表物件, 使用函式指令如下:

```
1 > split(x, f, drop = FALSE, ...)
2 > split(x, f, drop = FALSE, ...) <- value
3 > unsplit(value, f, drop = FALSE)
```

其中的主要引數分別為:

- `x` 為 1 個所要分割之資料框架.
- `f` 為 1 個類別變數, 因子物件 或 因素物件 (`factor`), 或包含因子物件的交互作用項. 可以使用用 `as.factor(f)` 將變數 `f` 轉換成因子變數.

- `drop = FALSE` 為 1 個邏輯指令，自動內設為 `FALSE`，是否要排除那些未在 `f` 因子之水準內的資料。
- `value` 為 1 個分割後之資料框架或列表物件向量。

例 15.3. R 內建資料 Puromycin

R 內建資料 `data(Puromycin)`，比較 Puromycin 與 Placebo (`state = "treated"` 或 “untreated”) 在不同 `substrate` 濃度 (`conc`) 下的酵素化學反應速率。

```

1 > ## split()
2 > data(Puromycin)
3 > head(Puromycin)
4 conc rate state
5 1 0.02 76 treated
6 2 0.02 47 treated
7 3 0.06 97 treated
8 4 0.06 107 treated
9 5 0.11 123 treated
10 6 0.11 139 treated
11 > #
12 > (df.split <- split(Puromycin, Puromycin$state))
13 $treated
14   conc rate state
15 1 0.02 76 treated
16 2 0.02 47 treated
17 3 0.06 97 treated
18 .....
19
20 $untreated
21   conc rate state
22 13 0.02 67 untreated
23 14 0.02 51 untreated
24 15 0.06 84 untreated
25 .....
26 > #
27 > lapply(df.split, summary)
28 $treated
29   conc          rate          state
30 Min.   :0.020   Min.   :47.0   treated  :12
31 1st Qu.:0.060   1st Qu.:104.5  untreated: 0
32 Median :0.165   Median :145.5
33 Mean   :0.345   Mean   :141.6
34 3rd Qu.:0.560   3rd Qu.:193.2
35 Max.   :1.100   Max.   :207.0
36
37 $untreated

```

```
38      conc          rate          state
39 Min.   :0.0200  Min.   : 51.0  treated  : 0
40 1st Qu.:0.0600  1st Qu.: 85.0  untreated:11
41 Median :0.1100  Median :115.0
42 Mean   :0.2764  Mean   :110.7
43 3rd Qu.:0.3900  3rd Qu.:137.5
44 Max.   :1.1000  Max.   :160.0
45 > #
46 > unsplit(df.split, Puromycin$state)
47   conc rate state
48 1  0.02  76  treated
49 2  0.02  47  treated
50 3  0.06  97  treated
51 .....
52 13 0.02  67  untreated
53 14 0.02  51  untreated
54 15 0.06  84  untreated
55 .....
56 > #
57 > split(Puromycin$rate, Puromycin$state)
58 $treated
59 [1] 76 47 97 107 123 139 159 152 191 201 207 200
60
61 $untreated
62 [1] 67 51 84 86 98 115 131 124 144 158 160
63 > #
64 > p.split <- split(Puromycin$rate, Puromycin$state)
65 > lapply(p.split, mean)
66 $treated
67 [1] 141.5833
68
69 $untreated
70 [1] 110.7273
71 > #
72 > unsplit(p.split, Puromycin$state)
73 [1] 76 47 97 107 123 139 159 152 191 201 207 200 67 51 84
74 [16] 86 98 115 131 124 144 158 160
```

15.10.2 合併資料框架函式: merge()

使用函式 `merge()` 可以對資料框架合併, 傳回列表物件. 函式指令使用如下:

```
1 > ## S3 method for class 'data.frame'
2 > merge(x, y, by = intersect(names(x), names(y)),
3         by.x = by, by.y = by, all = FALSE, all.x = all, all.y = all,
4         sort = TRUE, suffixes = c(".x", ".y"),
5         incomparables = NULL, ...)
```

其中引數為：

- `x, y` 為要合併的 資料框架.

- `by, by.x = "x.name", by.y = "y.name"`

設定資料框架 `x, y` 合併後的資料框架中，所要依據之共同之變數的欄位名稱 (common columns)，例如，`by.x = "x.name", by.y = y.name"`，內設相同的變數名 `by`.

- `all = FALSE, all.x = all, all.y = all`

`all.x` 為 1 個邏輯指令，當 `all.x = TRUE` 時，在所要合併之資料框架 `x`，若與資料框架 `y` 沒有相配之 `by.y` 列位 (rows)，會附加在合併後之資料框架最後，並設定在相對應之 `by.y` 變數為 `NA`. `all.y` 為 1 個邏輯指令，作用與 `all.x` 類似. `all` 表示 `all.x, all.y` 有相同之邏輯指令.

- `sort = TRUE` 回傳結果是否須依照 `column` 排序.

- `suffixes = c(".x", ".y")`]

在 `x, y` 合併之資料框架中，非合併所依據之共同之變數的其餘變數欄位名稱，若有仍有相同的變數欄位名稱，則會附加文字在其餘仍有相同名稱變數欄位.

- `incomparables` 設定一些數值不能被配對.

```

1 > ## merge()
2 > x <- data.frame(k1 = c(NA, NA, 3, 4, 5), k2 = c(1, NA, NA, 4, 5), id = 1:5, id1 = 11:15)
3 > y <- data.frame(k1 = c(NA, 2, NA, 4, 5), k2 = c(NA, NA, 3, 4, 5), id = 1:5, id2 = 11:15)
4 > cbind(x, y)
5   k1 k2 id id1 k1 k2 id id2
6  1 NA  1  1  11 NA NA  1  11
7  2 NA NA  2  12  2 NA  2  12
8  3  3 NA  3  13 NA  3  3  13
9  4  4  4  4  14  4  4  4  14
10 5  5  5  5  15  5  5  5  15
11 > #
12 > merge(x, y, by = c("id"))
13   id k1.x k2.x id1 k1.y k2.y id2
14 1    1     NA    1    11    NA    NA    11
15 2    2     NA    2    12    2     NA    12
16 3    3     NA    3    13    NA     3    13

```

```
17 4 4 4 4 14 4 4 14
18 5 5 5 5 15 5 5 15
19 > #
20 > merge(x, y, by.x = c("id1"), by.y = c("id2"))
21   id1 k1.x k2.x id.x k1.y k2.y id.y
22 1 11 NA 1 1 NA NA 1
23 2 12 NA NA 2 2 NA 2
24 3 13 3 NA 3 NA 3 3
25 4 14 4 4 4 4 4 4
26 5 15 5 5 5 5 5 5
27 > #
28 > merge(x, y, by = c("k1", "k2")) # NA's match
29   k1 k2 id.x id1 id.y id2
30 1 4 4 4 14 4 14
31 2 5 5 5 15 5 15
32 3 NA NA 2 12 1 11
33 > #
34 > merge(x, y, by = "k1") # NA's match, so 6 rows
35   k1 k2.x id.x id1 k2.y id.y id2
36 1 4 4 4 14 4 4 14
37 2 5 5 5 15 5 5 15
38 3 NA 1 1 11 NA 1 11
39 4 NA 1 1 11 3 3 13
40 5 NA NA 2 12 NA 1 11
41 6 NA NA 2 12 3 3 13
42 > #
43 > merge(x, y, by = "k2", incomparables = NA) # 2 rows
44   k2 k1.x id.x id1 k1.y id.y id2
45 1 4 4 4 14 4 4 14
46 2 5 5 5 15 5 5 15
47 > #
48 > merge(x, y, by = "k1", all = TRUE)
49   k1 k2.x id.x id1 k2.y id.y id2
50 1 2 NA NA NA NA 2 12
51 2 3 NA 3 13 NA NA NA
52 3 4 4 4 14 4 4 14
53 4 5 5 5 15 5 5 15
54 5 NA 1 1 11 NA 1 11
55 6 NA 1 1 11 3 3 13
56 7 NA NA 2 12 NA 1 11
57 8 NA NA 2 12 3 3 13
58 > #
59 > merge(x, y, by = "k1", all = FALSE)
60   k1 k2.x id.x id1 k2.y id.y id2
61 1 4 4 4 14 4 4 14
62 2 5 5 5 15 5 5 15
63 3 NA 1 1 11 NA 1 11
64 4 NA 1 1 11 3 3 13
65 5 NA NA 2 12 NA 1 11
66 6 NA NA 2 12 3 3 13
```

```
67 > #
68 > merge(x, y, by = "k1", all.x = TRUE)
69   k1 k2.x id.x id1 k2.y id.y id2
70 1 3    NA    3 13    NA    NA  NA
71 2 4    4    4 14    4    4 14
72 3 5    5    5 15    5    5 15
73 4 NA   1    1 11    NA    1 11
74 5 NA   1    1 11    3    3 13
75 6 NA   NA   2 12    NA    1 11
76 7 NA   NA   2 12    3    3 13
77 > #
78 > merge(x, y, by = "k1", all.x = FALSE)
79   k1 k2.x id.x id1 k2.y id.y id2
80 1 4    4    4 14    4    4 14
81 2 5    5    5 15    5    5 15
82 3 NA   1    1 11    NA    1 11
83 4 NA   1    1 11    3    3 13
84 5 NA   NA   2 12    NA    1 11
85 6 NA   NA   2 12    3    3 13
86 > #
87 > merge(x, y, by = "k1", all.y = TRUE)
88   k1 k2.x id.x id1 k2.y id.y id2
89 1 2    NA    NA  NA    NA    2 12
90 2 4    4    4 14    4    4 14
91 3 5    5    5 15    5    5 15
92 4 NA   1    1 11    NA    1 11
93 5 NA   1    1 11    3    3 13
94 6 NA   NA   2 12    NA    1 11
95 7 NA   NA   2 12    3    3 13
96 > #
97 > merge(x, y, by = "k1", all.y = FALSE)
98   k1 k2.x id.x id1 k2.y id.y id2
99 1 4    4    4 14    4    4 14
100 2 5   5    5 15    5    5 15
101 3 NA  1    1 11    NA    1 11
102 4 NA  1    1 11    3    3 13
103 5 NA  NA   2 12    NA    1 11
104 6 NA  NA   2 12    3    3 13
105 > #
106 > x.id <- c(1:6)
107 > x.f1 <- rep(c(1, 2), 3)
108 > x.f2 <- rep(c("m", "f"), each = 3)
109 > x.f3 <- rep(c("L", "M", "H"), length.out = 6)
110 > x.df <- data.frame(x.id = x.id, x.f1 = x.f1, x.f2 = x.f2, x.f3 = x.f3)
111 > #
112 > y.id <- c(3:8)
113 > x.f1 <- rep(c(1, 2), 3)
114 > y.f2 <- rep(c("m", "f"), each = 3)
115 > y.f3 <- rep(c("L", "M", "H", "E"), length.out = 6)
116 > y.df <- data.frame(y.id = y.id, x.f1 = x.f1, y.f2 = y.f2, y.f3 = y.f3)
```

```
117 > #
118 > cbind(x.df, y.df)
119   x.id x.f1 x.f2 x.f3 y.id x.f1 y.f2 y.f3
120 1    1    1    m    L    3    1    m    L
121 2    2    2    m    M    4    2    m    M
122 3    3    1    m    H    5    1    m    H
123 4    4    2    f    L    6    2    f    E
124 5    5    1    f    M    7    1    f    L
125 6    6    2    f    H    8    2    f    M
126 > #
127 > merge(x.df, y.df, by.x = "x.id", by.y = "y.id", sort = FALSE)
128   x.id x.f1.x x.f2 x.f3 x.f1.y y.f2 y.f3
129 1    3    1    m    H    1    m    L
130 2    4    2    f    L    2    m    M
131 3    5    1    f    M    1    m    H
132 4    6    2    f    H    2    f    E
133 > #
134 > merge(x.df, y.df, by.x = "x.id", by.y = "y.id", all = TRUE, sort = FALSE)
135   x.id x.f1.x x.f2 x.f3 x.f1.y y.f2 y.f3
136 1    3    1    m    H    1    m    L
137 2    4    2    f    L    2    m    M
138 3    5    1    f    M    1    m    H
139 4    6    2    f    H    2    f    E
140 5    1    1    m    L    NA <NA> <NA>
141 6    2    2    m    M    NA <NA> <NA>
142 7    7    NA <NA> <NA>    1    f    L
143 8    8    NA <NA> <NA>    2    f    M
144 > #
145 > merge(y.df, x.df, by.x = "y.id", by.y = "x.id", all = TRUE, sort = TRUE)
146   y.id x.f1.x y.f2 y.f3 x.f1.y x.f2 x.f3
147 1    1    NA <NA> <NA>    1    m    L
148 2    2    NA <NA> <NA>    2    m    M
149 3    3    1    m    L    1    m    H
150 4    4    2    m    M    2    f    L
151 5    5    1    m    H    1    f    M
152 6    6    2    f    E    2    f    H
153 7    7    1    f    L    NA <NA> <NA>
154 8    8    2    f    M    NA <NA> <NA>
155 > #
156 > merge(x.df, y.df, by.x = "x.id", by.y = "y.id",
157           sort = FALSE, all = TRUE, suffixes = c(".X", ".Y"))
158   x.id x.f1.X x.f2 x.f3 x.f1.Y y.f2 y.f3
159 1    3    1    m    H    1    m    L
160 2    4    2    f    L    2    m    M
161 3    5    1    f    M    1    m    H
162 4    6    2    f    H    2    f    E
163 5    1    1    m    L    NA <NA> <NA>
164 6    2    2    m    M    NA <NA> <NA>
165 7    7    NA <NA> <NA>    1    f    L
166 8    8    NA <NA> <NA>    2    f    M
```

15.11 資料框架函式: with()

使用資料框架函式 `with()` 直接使用 資料框架 `data.name` 內的變數, 可以在使用資料框架 `data.name` 內的變數 `var.name` 時, 避免使用 `data.name$var.name` 的繁複形式, 也可避免使用 `attach()` 所帶來不可預測的變數名稱衝突. 函式 `with()` 的指令為

```
1 > with(data, expr, ...)
2 > within(data, expr, ...)
```

其中引數為:

- `data` 為 資料框架名.
- `expr` 對 資料框架 `data` 進行的運算式.

```
1 > ## with ()
2 > data(Puromycin)
3 > with(Puromycin, mean(rate))
4 [1] 126.8261
5 > with(Puromycin, fivenum(rate))
6 [1] 47.0 91.5 124.0 158.5 207.0
7 > with(Puromycin, t.test(rate ~ state))
8
9      Welch Two Sample t-test
10 data: rate by state
11 t = 1.6375, df = 19.578, p-value = 0.1175
12 alternative hypothesis: true difference in means is not equal to 0
13 95 percent confidence interval:
14 -8.504864 70.216985
15 sample estimates:
16 mean in group treated mean in group untreated
17             141.5833                  110.7273
18 > #
19 > ## try
20 > with(Puromycin, plot(x = conc, y = rate))
21 > plot(rate ~ conc, data = Puromycin)
```