

```

1 # The apply family
2
3 # Apply functions allow you to apply a function across different subsets of data.
4 # For example, you may want to obtain the mean (max, min, or median) for each column of
  a data.frame.
5 # It Offers an alternative to explicit iteration using for() loop. Can be faster.
6
7 # The apply family includes
8 # apply(): apply a function to rows or columns of a matrix or data.frame
9 # lapply(): apply a function to elements of a list or data.frame
10 # sapply(): same as lapply(), but simplifies the output
11 # tapply(): apply a function to levels of a factor vector.
12
13 # 1. apply()
14 # apply(x, MARGIN=1, FUN=my.fun): apply my.fun() across rows of a matrix or data.frame x.
15 # apply(x, MARGIN=2, FUN=my.fun): apply my.fun() across columns of a matrix or
  data.frame x.
16
17 # Examples
18
19 x <- matrix(rnorm(9), 3, 3)
20 x
21
22 apply(x, MARGIN = 1, FUN=min)
23
24 for (i in 1:nrow(x)) {
25   print(min(x[i,]))
26 }
27
28 apply(x, MARGIN = 2, FUN=sum)
29 colSums(x)
30
31 # Applying a function that takes extra arguments
32 # Sometimes we want to use a function over rows or columns of a matrix, that takes extra
33 # arguments (besides the row or column itself). We can put these as inputs to apply().
34
35 setwd("C:/Users/elsect_main")
36 rev_exp0 <- read.csv("district_rev_exp.csv", na.strings = "-")
37 head(rev_exp0)
38
39 maximum <- apply(rev_exp0[, -c(1,3,4)], MARGIN=2, FUN=max, na.rm=TRUE) # na.rm=TRUE is
  an argument of max.
40 maximum
41
42 first <- apply(rev_exp0[, -c(1,3,4)], MARGIN=2, FUN=which.max)
43 first
44 rev_exp0[first[1],]
45
46 second <- apply(rev_exp0[-first[1], -c(1,3,4)], MARGIN=2, FUN=which.max)
47 second
48 rev_exp0[second[1],]
49
50 # We can use this apply() to our own functions.
51 my.fun <- function(x) {
52   m1 <- median(x, na.rm=TRUE)
53   m2 <- mean(x, na.rm=TRUE)
54   return(c(m1,m2))
55 }
56
57 apply(rev_exp0[, -c(1,3,4)], MARGIN=2, FUN=my.fun)
58
59 # lapply(), elements of a list (data.frame)
60
61 my.list = list(nums=seq(0.1,0.6,by=0.1), chars=letters[1:12],
62               TF=sample(c(TRUE,FALSE), 6, replace=TRUE))
63 my.list
64
65 lapply(my.list, FUN=mean) # The type of outcome is always list.
66

```

```

67 # lapply() with extra arguments
68
69 mean.omitting.one = function(i,vec) {
70   return(mean(vec[-i]))
71 }
72
73 my.vec = rev_exp0[, "TOTALEXP"]
74 n = length(my.vec)
75 my.vec.jack = lapply(1:n, FUN=mean.omitting.one, vec=my.vec)
76
77 head(my.vec.jack) # It's a list, and here are the first 6 elements
78
79 # sapply(), elements of a list or data.frame
80 # The sapply() function works just like lapply(), but tries to simplify the return
  value whenever possible.
81 # E.g., most common is the conversion from a list to a vector
82
83 my.vec.jack1 = sapply(1:n, FUN=mean.omitting.one, vec=my.vec)
84 length(my.vec.jack1)
85
86 sqrt((n-1)/n * sum((my.vec.jack1 - mean(my.vec.jack1))^2)) # Jackknife standard error
87 sqrt((n-1)^2/n * sd(my.vec.jack1)
88
89 sd(my.vec)/sqrt(n) # conventional standard error
90
91 # tapply(), levels of a factor vector
92 # tapply(x, INDEX=my.index, FUN=my.fun): apply my.fun() to subsets of elements in x
  that share a common level in my.index
93
94 str(rev_exp0)
95
96 # tapply is designed to work on a single vector (column), while aggregate can work with
  multiple columns.
97 tapply(rev_exp0[, "TOTALREV"], INDEX=list(ST=rev_exp0$STATE, EnR = rev_exp0$ENROLL >
  1000), FUN=mean, na.rm=TRUE)
98 aggregate(rev_exp0[, c("ENROLL", "TOTALREV", "TOTALEXP")], by=list(ST = rev_exp0$STATE,
  EnR = rev_exp0$ENROLL > 1000), FUN=mean, na.rm=TRUE)
99
100 ## Exercise 1
  #####
101
102 president <- c("Obama:2009-2007", "Bush:2001-2009", "Clinton:1993-2001",
  "Bush:1989-1993", "Reagan:1981-1989")
103
104 # Using strsplit(), split names from terms.
105
106
107 # Using lapply(), convert to uppercase strings.
108
109
110
111 select_first <- function(x) {
112   x[1]
113 }
114
115 # Using lapply() and select_first, select each president name.
116
117
118
119 ## Exercise 2
  #####
120
121 mtcars
122 n <- nrow(mtcars)
123
124 # Using the Jackknife procedure introduced above, calculate the standard error of the
  slope coefficient.

```

```
125 # Create a function that runs a regression with omitting one observation.
126 # Define x=mtcars$wt, y=mtcars$mpg
127 # Using sapply(), apply this function to all possible subsamples leaving one
    observation out.
128 # Calculate the Jackknife standard error, using the formula introduced above.
129
130 # Compare this with the conventional standard error.
131
132
133
134
```