

Lecture 4 Data management: Part I

- Control structure
- Missing value
- Dates
- Useful functions
- How to write our own functions

4.1 Control structure

See ?Control

Avoid using loops in R. I taught a workshop on efficient coding and computing and explained why. Here is the link: <https://github.com/ly129/MiCM>
(<https://github.com/ly129/MiCM>)

4.1.1 For loop

In [1]:

```
df <- data.frame(names = c("Lucy", "John", "Mark", "Candy"),  
                  score = c(67, 56, 87, 91))  
df
```

A data.frame: 4 ×

2

names	score
-------	-------

<fct>	<dbl>
-------	-------

Lucy	67
------	----

John	56
------	----

Mark	87
------	----

Candy	91
-------	----

In [2]:

```
x <- NULL
for (i in 1:5){
  x[i] = 2*i
}
x
```

2 4 6 8 10

4.1.2 While loop

Two useless operators in R that I found useful for teaching: modulus and integer division.

In [3]:

```
9 %% 2    # 9 mod 2
9 %/% 2
```

1

4

Can we write a while loop to do the two operations at the same time?

In [4]:

```
# y %% x
i <- 0
y <- 9
x <- 2
while (y>=x){
  y <- y - x
  i <- i + 1
}
y    # modulus
i    # integer division
# why?
```

1

4

4.1.3 If, else, ifelse

Not a loop. `ifelse` is the amazing vectorized alternative to `if ... else,`

Once is enough -

```
-- "Honey, on your way home, buy 6 oranges at the supermar  
ket. If they have watermelons, get 1."
```

```
-- Mr. Programmer came home with 1 orange.
```

```
-- Furious girlfriend, "Why the [--beep--] did you get onl  
y 1 orange?"
```

```
-- "Because they have watermelons."
```

In [5]:

```
watermelon <- FALSE
no.orange <- if (watermelon == TRUE){
  "Buy 1 orange"
} else {
  print("Buy 6 oranges") # As seen in class, print() is useful here.
}
no.orange
```

[1] "Buy 6 oranges"

'Buy 6 oranges'

In [6]:

```
# I prefer a simple function, ifelse(test, yes, no)
watermelon <- F
ifelse(watermelon == TRUE, yes = "Buy 1 orange", no = "Buy 6 oranges")
```

'Buy 6 oranges'

In [7]:

```
# ifelse is vectorized
df$pass <- ifelse(test = df$score >= 65, yes = TRUE, no = FALSE)
df
```

A data.frame: 4 × 3

names	score	pass
<fct>	<dbl>	<lgl>
Lucy	67	TRUE
John	56	FALSE
Mark	87	TRUE
Candy	91	TRUE

4.1.4 Repeat loop

In [8]:

```
i <- 0
# repeat {system("say Because they have watermelons!")
#       i <- i + 1
#       if (i>=3){
#         break
#       }
# }
```

Exercise: use the repeat loop to calculate $9 \% \% 2$ and $9 \% / \% 2$.

In []:

Are there any situations that loops cannot be replaced by vector operations?

4.2 Missing values

- NA

In [9]:

```
# Using indices from last lecture to change specific entries in R objects
df.copy <- df
df.copy$score[2] <- df.copy$names[3] <- NA
df.copy
```

A data.frame: 4 × 3

names	score	pass
<fct>	<dbl>	<lgl>
Lucy	67	TRUE
John	NA	FALSE
NA	87	TRUE
Candy	91	TRUE

In [10]:

```
is.na(df.copy)
```

A matrix: 4 × 3 of type lgl

names	score	pass
FALSE	FALSE	FALSE
FALSE	TRUE	FALSE
TRUE	FALSE	FALSE
FALSE	FALSE	FALSE

In [11]:

```
# Total number of cells with missing values
sum(is.na(df.copy))
```

In [12]:

```
# Whether a data point (row) is complete  
complete.cases(df.copy)
```

TRUE FALSE FALSE TRUE

In [13]:

```
!complete.cases(df.copy)
```

FALSE TRUE TRUE FALSE

In [14]:

```
# Incomplete data points  
df.copy[!complete.cases(df.copy), ]  
# Recall the logical operator "!"
```

A data.frame: 2 × 3

	names	score	pass
	<fct>	<dbl>	<lgl>
2	John	NA	FALSE
3	NA	87	TRUE

In [15]:

```
# Taking the average score  
mean(df.copy$score)
```

<NA>

In [16]:

```
mean(df.copy$score, na.rm = TRUE)
```

81.66666666666667

In [17]:

```
sum(df.copy$score)
sum(df.copy$score, na.rm = T)
```

<NA>

245

In [18]:

```
na.omit(df.copy)
```

A data.frame: 2 × 3

	names	score	pass
	<fct>	<dbl>	<lgl>
1	Lucy	67	TRUE
4	Candy	91	TRUE

4.3 Dates

In [19]:

```
Sys.Date()
# Note the standard date format in R
```

2019-09-25

In [20]:

```
Sys.time() # Eastern Daylight Time
```

```
[1] "2019-09-25 18:56:28 EDT"
```


In [21]:

```
date()
```

'Wed Sep 25 18:56:28 2019'

In [22]:

```
first.hw.post <- as.Date("Oct 4, 2018", tryFormats = "%b %d, %Y"  
)  
first.hw.post
```

2018-10-04

In [23]:

```
first.hw.due <- as.Date("2018년10월11일", tryFormats = "%Y년%m월%d일"  
)  
first.hw.due  
# Just want to show you that any format can be recognized.  
# As long as you can let R know how to read it.
```

2018-10-11

In [24]:

```
# Help file: Date-time Conversion Functions to and from Character  
# ?strptime
```

In [25]:

```
first.hw.due - Sys.Date()
```

Time difference of -349 days

In [26]:

```
as.numeric(Sys.Date())
```

18164

In [27]:

```
# Time origin of R  
Sys.Date() - as.numeric(Sys.Date())
```

1970-01-01

In [28]:

```
# How long does it take R to load the survival package  
time0 <- proc.time()  
library(survival)  
proc.time() - time0
```

user	system	elapsed
0.765	0.053	0.822

In [29]:

```
format(Sys.Date(), format = "%A %B %d %Y")
```

'Wednesday September 25 2019'

4.4 Useful functions

4.4.1 Numeric functions

In [30]:

```
# Absolute value  
abs(-3)
```

3

In [31]:

```
ceiling(3.14159)
```

4

In [32]:

```
floor(3.14159)
```

3

In [33]:

```
trunc(3.14159)
```

3

In [34]:

```
signif(3.14159, 3)
```

3.14

In [35]:

```
# ?round
```

Use these functions to calculate $9 \% \% 2$ and $9 \% / \% 2$.

In []:

4.4.2 Character functions

- `paste()` and `expression()`
 - `paste()` put text and variable values together into a text string.
 - `expression()` can be used to display math symbols when needed, e.g. in plot titles.

Few situations where you have to deal with text in R

- Data frame entries
- Plot title, labels, legends, etc...

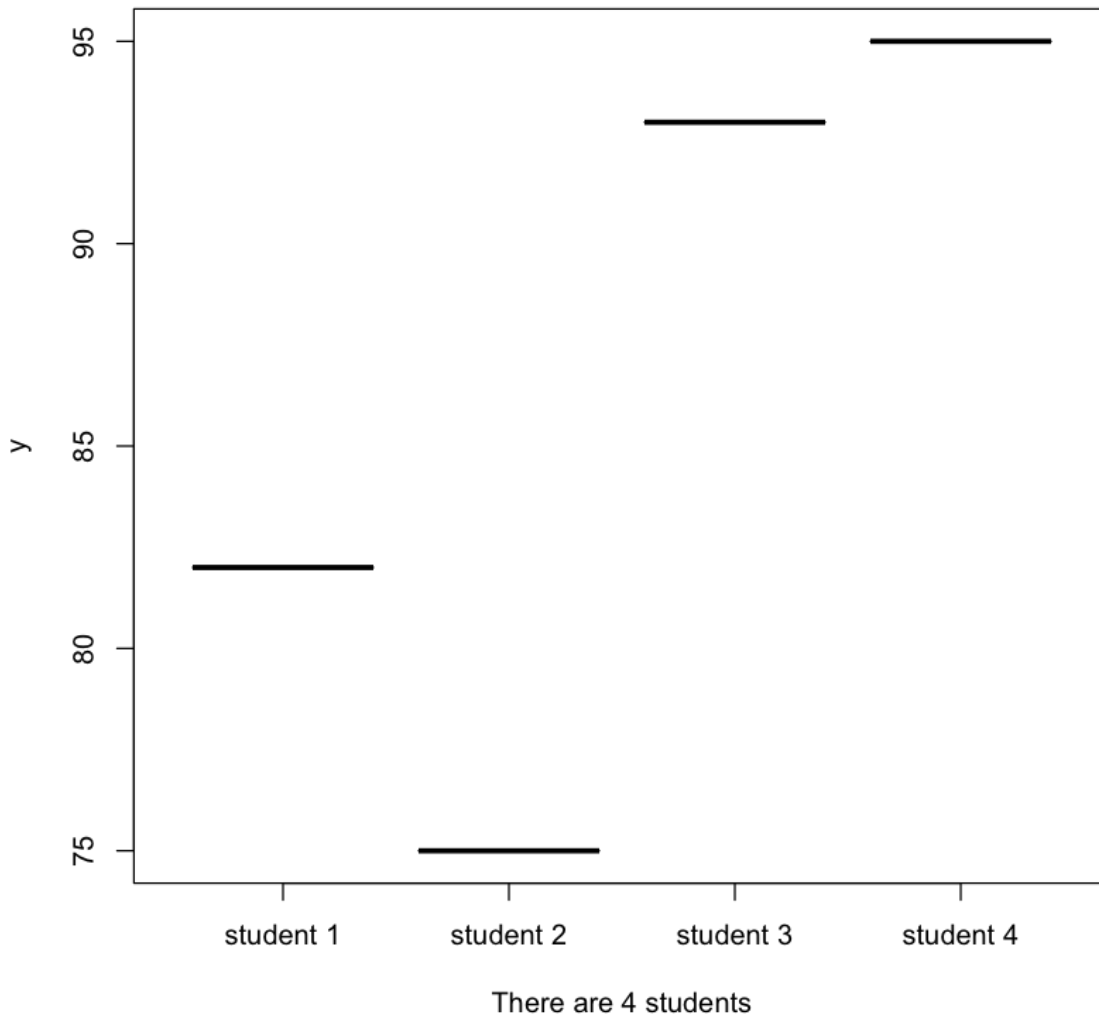
In [36]:

```
for (i in 1:4){
  df$student.no[i] <- paste("student", i)
  df$curved.score[i] <- round(sqrt(df$score[i]) * 10)
}
str(df)

n <- nrow(df)
plot(as.factor(df$student.no), df$curved.score,
     # Math symbols in text
     main = expression(paste("Score is ", alpha, ", curved score
is ", sqrt(alpha)%*%10)),
     # Variable value in text
     xlab = paste("There are", n, "students"))
```

```
'data.frame':   4 obs. of  5 variables:
 $ names      : Factor w/ 4 levels "Candy","John",.
.: 3 2 4 1
 $ score      : num  67 56 87 91
 $ pass       : logi  TRUE FALSE TRUE TRUE
 $ student.no : chr   "student 1" "student 2" "stude
nt 3" "student 4"
 $ curved.score: num  82 75 93 95
```

Score is α , curved score is $\sqrt{\alpha} \times 10$



4.4.3 `apply` family functions

Some say that `apply()` family functions distinguish R experts and newbies.

Again, much more in my workshop <https://github.com/ly129/MiCM>
(<https://github.com/ly129/MiCM>).

`apply()`

In [37]:

```
df.scores <- df[, c("score", "curved.score")]; df.scores
```

A data.frame: 4 × 2

score	curved.score
<dbl>	<dbl>
67	82
56	75
87	93
91	95

In [38]:

```
apply(df.scores, MARGIN = 2, FUN = mean)
```

```
score
75.25
curved.score
86.25
```

In [39]:

```
apply(df.scores, MARGIN = 1, FUN = diff)  # diff() calculates the difference - see Section 4.4.4
```

```
15 19 6 4
```

In [40]:

```
myarray <- array(1:12, dim = c(2,3,2)); print(myarray)
```

, , 1

	[,1]	[,2]	[,3]
[1,]	1	3	5
[2,]	2	4	6

, , 2

	[,1]	[,2]	[,3]
[1,]	7	9	11
[2,]	8	10	12

In [41]:

```
apply(myarray, MARGIN = c(2, 3), sum)
```

A matrix:

3 × 2 of

type int

3 15

7 19

11 23

sapply()

In [42]:

```
sapply(df, is.numeric)
```

names

FALSE

score

TRUE

pass

FALSE

student.no

FALSE

curved.score

TRUE

There is also lapply(), tapply(), etc...

And their parallel versions mclapply(), parLapply() in the 'parallel' package for parallel computing.

4.4.4 Other useful functions

In [43]:

```
age=c(1,6,4,5,8,5,4,3)
weight=c(45,65,34)
age
```

1 6 4 5 8 5 4 3

In [44]:

```
mean(age)
```

4.5

In [45]:

```
prod(age)
```

57600

In [46]:

```
median(age)
```

4.5

In [47]:

```
var(age)  
sd(age)
```

4.28571428571429

2.07019667802706

In [48]:

```
max(age)  
min(age)  
range(age)
```

8

1

1 8

In [49]:

```
which.max(age)    #returns the index of the greatest element of x  
which.min(age)    #returns the index of the smallest element of x
```

5

1

In [50]:

```
seq(from = 0, to = 1, by = 0.25)
quantile(age, probs = seq(from = 0, to = 1, by = 0.25))
# Returns the specified quantiles.
```

0 0.25 0.5 0.75 1

0%

1

25%

3.75

50%

4.5

75%

5.25

100%

8

In [51]:

```
unique(age) # Gives the vector of distinct values
```

1 6 4 5 8 3

In [52]:

```
diff(age) # Replaces a vector by the vector of first differences
```

5 -2 1 3 -3 -1 -1

In [53]:

```
sort(age) # Sorts elements into order
```

1 3 4 4 5 5 6 8

In [54]:

```
order(age)
age[order(age)]    # x[order(x)] orders elements of x
```

```
1  8  3  7  4  6  2  5
```

```
1  3  4  4  5  5  6  8
```

In [55]:

```
cumsum(age)    # Cumulative sums
cumprod(age)    # Cumulative products
```

```
1  7  11  16  24  29  33  36
```

```
1  6  24  120  960  4800  19200  57600
```

In [56]:

```
age
cat <- cut(age, breaks = 2); cat    # Divide continuous variable
in factor with n levels
table(cat)    # Cross tabulation and table creation
```

```
1  6  4  5  8  5  4  3
```

```
(0.993,4.5]  (4.5,8.01]  (0.993,4.5]  (4.5,8.01]  (4.5,8.01]  (4.5,8.01]
```

```
(0.993,4.5]  (0.993,4.5]
```

► Levels:

```
cat
(0.993,4.5]  (4.5,8.01]
      4              4
```

In [57]:

```
# Split the variable into categories  
age.cat <- split(age, cut(age,2))  
age.cat
```

`$`[0.993,4.5]``

1 4 4 3

`$`[4.5,8.01]``

6 5 8 5

In [58]:

```
# split() gives a list  
str(age.cat)
```

List of 2

`$ (0.993,4.5]: num [1:4] 1 4 4 3`

`$ (4.5,8.01] : num [1:4] 6 5 8 5`

In [59]:

```
# lapply: list apply  
lapply(age.cat, mean)
```

`$`[0.993,4.5]``

3

`$`[4.5,8.01]``

6

4.5 Write our own functions

- `function()`

In [60]:

```
# The structure
```

```
func_name <- function(argument){  
  statement  
}
```

Write my own function of x^y :

In [61]:

```
X.to.the.power.of.Y <- function(y, x){  
  x^y  
}  
X.to.the.power.of.Y(x = 3, y = 2)  
X.to.the.power.of.Y(3, 2)      # Following a question in class, note the difference.
```

9

8

Uses:

- If we need to do some operation a lot later.
- Work with apply() family.
 - The 'FUN' argument in apply() family functions only take the name of the functions only.
 - No arguments, operators or combinations of these allowed.

Example: calculate the square of the score

In [62]:

```
df.scores
```

A data.frame: 4 × 2

score	curved.score
<dbl>	<dbl>
67	82
56	75
87	93
91	95

In [63]:

```
# The following code does not work  
# apply(df.scores, MARGIN = 2, FUN = ^2)
```

In [64]:

```
# Instead we can do  
my.fun <- function(x){x^2}  
apply(df.scores, MARGIN = 2, FUN = my.fun)
```

A matrix: 4 × 2 of type

dbl

score	curved.score
4489	6724
3136	5625
7569	8649
8281	9025

Exercise: write our own function to calculate $x \% y$ and $x \%/\% y$.

- Note how to return the output in `function()` and assess the results correspondingly.

In [65]:

```
# Two inputs, y and x, so two arguments

# Option 1 - use %% and %/% operators
modulus1 <- function(y, x){
  mod <- y %% x
  int.div <- y %/% x
  return(list(modulus=mod, integer.division=int.div))
}
out1 <- modulus1(y = 9, x = 2)
print(out1)
str(out1)
```

```
$modulus
```

```
[1] 1
```

```
$integer.division
```

```
[1] 4
```

```
List of 2
```

```
 $ modulus          : num 1
```

```
 $ integer.division: num 4
```

In [66]:

```
out1$modulus
out1$integer.division
```

```
1
```

```
4
```

In [67]:

```
# Option 2 - use trunc() or floor()

modulus2 <- function(y, x){
  mod <- trunc(y/x)      # or floor(y/x)
  int.div <- y - x * mod
  return(c(modulus=mod, integer.division=int.div))
}
out2 <- modulus2(9, 2)
print(out2)
str(out2)
```

```
      modulus integer.division
      4          1
Named num [1:2] 4 1
- attr(*, "names")= chr [1:2] "modulus" "integer.division"
```

In [68]:

```
out2[1]
out2[2]
```

modulus: 4

integer.division: 1

In [69]:

```
attr(out2, "names")
```

'modulus' 'integer.division'

In [70]:

```
# Option 3 - use loops
```

```
modulus3 <- function(y, x){  
  i <- 0  
  while (y>=x){  
    y <- y - x  
    i <- i + 1  
  }  
  return(cat("modulus=", y, ", Integer division=", i)) # modulus  
}
```

```
# I want modulus(y, x) to give me 'y mod x' for any integers y and x.
```

```
out3 <- modulus3(9, 2)
```

```
# Note that without printing out3, the result is already shown.
```

```
modulus= 1 , Integer division= 4
```

In []: