

STAT 2005 – PROGRAMMING LANGUAGES FOR STATISTICS

TUTORIAL 4 ATTRIBUTES

2020

LIU Ran

Department of Statistics, The Chinese University of Hong Kong

1 Define Functions

Aim: Write $f(x) = x^2$ and $g(x) = x^2 + 2x + 1$.

```

1 > f = function(x) {
2 +   y = x^2
3 +   return(y)
4 + }
5 > f(-3:3) # Only if the calculating process in your function allows the vector input.
6 [1] 9 4 1 0 1 4 9
7 > g = function(x) {
8 +   f(x)+2*x+1 # x^2+2*x+1
9 + }
10 > g(-1) # you can call the defined function in another function
11 [1] 0

```

Remark 1.1. Previous defined functions can be used in new functions' definitions. Therefore, it is better to separate one huge function into several small functions. It will be helpful for debugging and reading.

You can also use the function itself(recursion): $f(n) = n! = n(n-1)(n-2)\dots * 2 * 1$, so we can get the recurrence formula $f(n) = n * f(n-1) = n * (n-1) * f(n-2) = \dots$

```

1 # Calculate the factorial
2 Factorial<-function(n) {
3   if(n == 1){
4     return(1)
5   }
6   else{
7     # each time the argument will minus one,
8     # and the result will be multiplied by the current value
9     return(n*Factorial(n-1))
10  }
11 }
12 # end at the point one when you call Factorial(1)

```

Please note the difference between the case that the argument is a vector and the case that there are several arguments. It is quite misleading.

```

1 > sum_of_three<-function(a,b,c) {
2 +   return(a+2*b+3*c)
3 + }
4 > sum_of_three(1:3) # input a vector
5 Error in sum_of_three(1:3) : argument "b" is missing, with no default
6 # Use the comma to separate the arguments
7 > sum_of_three(1,2,3) # input three arguments by order, the first is a, the second is b ...
8 [1] 14
9 > sum_of_three(a=1,c=3,b=2) # you can also use the keyword to input arguments
10 [1] 14

```

If there are the default values, it is better to use the keyword to input values.

```

1 # add the default value in the middle
2 > sum_of_three<-function(a,b=2,c) {
3 +   return(a+2*b+3*c)
4 + }
5 > sum_of_three(1,3)
6 Error in sum_of_three(1, 3) : argument "c" is missing, with no default
7 > sum_of_three(a=1,c=3)

```

```

8 [1] 14
9 # The function will receive the keyword first,
10 # and the remaining arguments will receive the input by order.
11 > sum_of_three(1,a=3,3) # a=3,b=1,c=3; keyword is 'a', remain 'b' and 'c'
12 [1] 14
13 > sum_of_three(1,b=5,3) # a=1,b=5,c=3; keyword is 'b', remain 'a' and 'c'
14 [1] 20
15 # not recommend the mixture input (order and keyword)

```

The three dots ellipsis arguments `...`: it means that the function is designed to take any number of named or unnamed arguments.

```

1 > s<-function(multiplier, ...){
2 +   print(..1) # print the first argument in the ellipsis arguments
3 +   # instead of using print(...) directly
4 +   print(list(...)) # list(a=1,b=2,c=3), transform the ellipsis arguments into a list
5 +   return(multiplier*sum(...))
6 + }
7 > s(2, a=1, b=2, c=3) # 2*sum(a=1, b=2, c=3) = 12
8 [1] 1
9 $a
10 [1] 1
11
12 $b
13 [1] 2
14
15 $c
16 [1] 3
17
18 [1] 12
19
20 # Also we can use the keyword, the remaining arguments will be included in the ellipsis
21 > s(a=1, multiplier = 3, 2) # list(a=1, 2)
22 [1] 1
23 $a
24 [1] 1
25
26 [[2]]
27 [1] 2
28
29 [1] 9

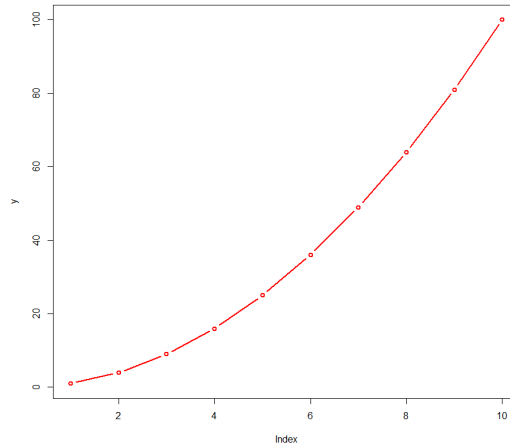
```

The ellipsis arguments are often used in the function for drawing. Because there are many parameter settings in the plot command. It is convenient to use three dots arguments to receive the settings.

```

1 # f(x) = x^2
2 > draw_square<-function(data, ...){
3 +   y = data^2 # get the square of data
4 +   plot(y, ...) # transfer the parameters to the function 'plot'
5 + }
6 > draw_square(seq(1,10), col='red', type='b', lwd=2)
7 # no need for defining arguments 'col', 'type', 'lwd'

```



Remark 1.2. Use `>>>` to insert indents in functions/loops, which is helpful for understanding and debugging.

Remark 1.3. `return` is only useful inside functions. It can be omitted if the last line in the function is the desired output.

Remark 1.4. The `plot` function's arguments also have the three dots argument. Getting familiar with the three dots will deepen the understanding of some complex functions.

2 if

Aim: Write an indicator function $\mathbb{1}(x, b) = \begin{cases} 1 & , \text{ if } x \geq b \\ 0 & , \text{ if } x < b \end{cases}$.

```

1 > I = function(x, b){
2 +   if (x>=b){
3 +     out = 1
4 +   } else {
5 +     out = 0
6 +   }
7 +   out
8 + }
9 > I(1, 3)
10 [1] 0
11 # input a vector
12 > I(1:5, 2)
13 [1] 0
14 Warning message:
15 In if (x >= b) { :
16   the condition has length > 1 and only the first element will be used
17 > I = Vectorize(I)
18 > I(1:5, 2)
19 [1] 0 1 1 1 1
20 > I(c(-1, 4, 9), 4:6)
21 [1] 0 0 1

```

For the vectorized version of if-else, you can also use `ifelse`:

```

1 > I = function(x, b){
2 +   ifelse(x>=b, 1, 0)
3 + }
4 > I(1:5, 2)
5 [1] 0 1 1 1 1
6 > I(c(-1, 4, 9), 4:6)
7 [1] 0 0 1

```

Remark 2.1. The 'and' operator `&` is vectorized but `&&` is not.

Remark 2.2. `Vectorize` can enable a function to support vector inputs. Though it is quite convenient, it can cause serious mistakes sometimes. A better way is to define a function that can support vector inputs at the very beginning.

3 Switch case

```
switch(expression, case1, case2, case3....)
```

Here, the expression is evaluated and based on this value, the corresponding item in the list is returned. If the value evaluated from the expression matches with more than one item of the list, `switch()` function returns the first matched item.

```
1 > switch(2, "red", "green", "blue")
2 [1] "green"
3 > switch('a', a="red", b="green", c="blue")
4 [1] "red"
```

It is useful for the case which has several complex pipelines(functions):

```
1 # simple pipelines example
2 FUNC<-function(data, pipeline_sel){
3   switch(pipeline_sel,
4     '1' = mean(data), # Function1
5     '2' = sd(data),   # Function2
6     '3' = sum(data)   # Function3
7   )
8   FUNC(c(1,2,3), '3') # return 6
```

4 Loops

4.1 for Loop

Aim: Compute $S = \sum_{i=1}^{100} i$.

```
1 > S = 0
2 > for(i in 1:100){
3 +   # already set the beginning, the ending and each step
4 +   S = S+i
5 + }
6 > S
7 [1] 5050
```

4.2 while Loop

Aim: Compute $S = \sum_{i=1}^{100} i$.

```
1 > S = 0
2 > i = 1
3 > while (i<=100){
4 +   S = S+i
5 +   i = i+1 # should update the i
6 + }
7 > S
8 [1] 5050
```

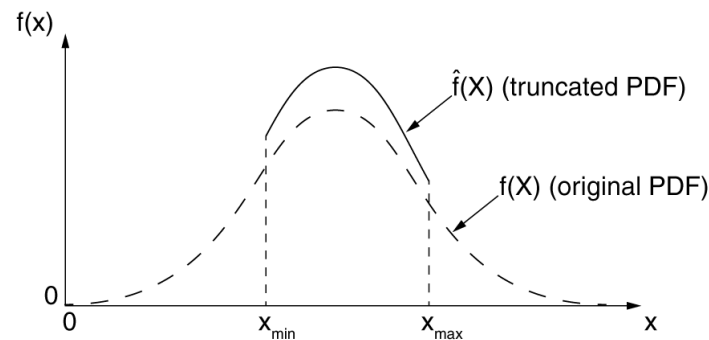
Remark 4.1. Please remember to set the stopping criterion every time when using the loops. Unstoppable loops may cause very serious problems.

Remark 4.2. If you only know the stopping criterion instead of the ending times, use `repeat` or `while` instead of `for`.

Remark 4.3. Although using loop can simplify the problem a lot, it is much slower than vectorized operations in general. Therefore, avoid using loops inside loop unless there is no other choice.

5 An example

We want to get samples from a truncated normal distribution: the truncated normal distribution is the probability distribution derived from that of a normally distributed random variable by bounding the random variable from either below or above (or both).



```

1 # get samples from a truncated normal distribution
2 rnorm_trunc<-function(lower, upper, count, mean, sd) {
3   sample <- rnorm(count, mean = mean, sd = sd)
4   looped <- 0
5
6   # exclude samples outside the interval
7   # re-sample the same size of bad samples from the original distribution
8
9   # any() will return False as long as one element is False
10  while (any(sample < lower | sample > upper)) {
11
12    # get the number of bad samples
13    bad <- (sample < lower | sample > upper) # logical indexes
14    bad_values <- sample[bad]
15    count_bad <- length(bad_values)
16
17    # re-sample from the original distribution
18    new_vals <- rnorm(count_bad, mean = mean, sd = sd)
19
20    # fix the samples
21    sample[bad] <- new_vals
22
23    # add a stopping criterion to avoid some extreme cases
24    # if after 10000*count times, we still cannot find a proper sample vector.
25    if(looped > 10000*count){
26      return(NA)
27    }
28    looped <- looped + 1
29  }
30  return(sample)
31 }
32
33 > rnorm_trunc(lower = -1, upper = 0, count = 5, mean = 0, sd = 1)
34 [1] -0.2781313 -0.8784863 -0.8221632 -0.4314080 -0.4344664

```

Similarly, if we want to get the truncated samples from any distribution:

Here, we'll use the ellipsis arguments `...` to collect an arbitrary set of parameters and pass them on to internal function calls. When defining a function to take `...`, it is usually specified at last. So, we'll write a function called `sample_trunc()` that takes five parameters:

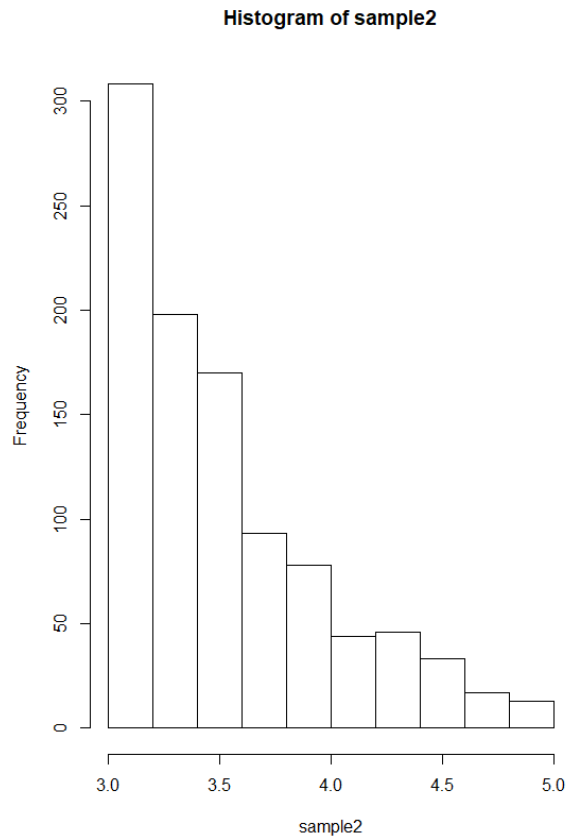
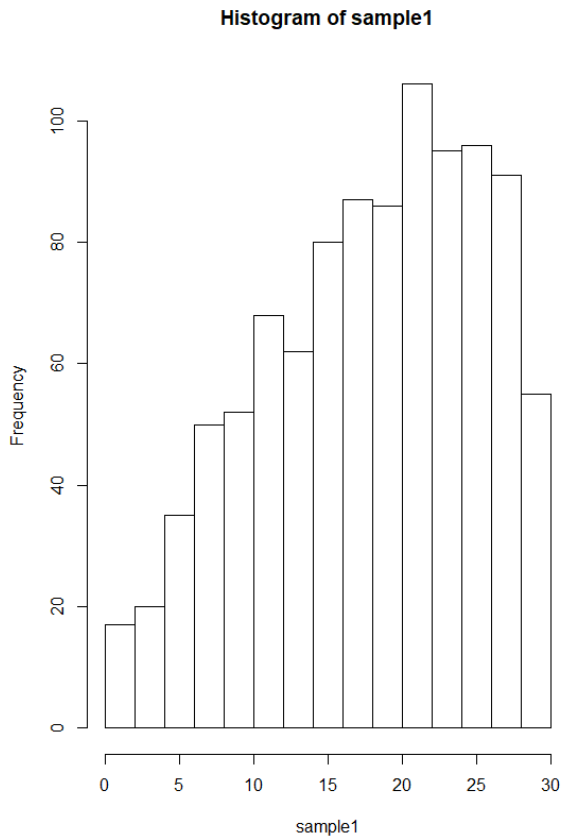
1. The lower limit, `lower`.

2. The upper limit, upper.
3. The sample size to generate, count.
4. The function to call to generate samples, `sample_func()`.
5. Additional parameters to pass on to `sample_func()`

```

1 sample_trunc<-function(lower, upper, count, sample_func, ...){
2   sample <- sample_func(count, ...)
3   looped <- 0
4
5   while (any(sample < lower | sample > upper)) {
6     # fix the sample
7     bad <- (sample < lower | sample > upper) # logical
8     bad_values <- sample[bad]
9     count_bad <- length(bad_values)
10
11    new_vals <- sample_func(count_bad, ...)
12    sample[bad] <- new_vals
13
14    if(looped > 10000*count){
15      return(NA)
16    }
17    looped <- looped + 1
18  }
19  return(sample)
20 }
21
22
23 sample1 <- sample_trunc(0, 30, 1000, rnorm, mean = 20, sd = 10)
24 sample2 <- sample_trunc(3, 5, 1000, rexp, rate = 1.5)
25
26
27 par(mfrow=c(1,2))
28 hist(sample1)
29 hist(sample2)

```



6 Attributes(optional)

This is the prior knowledge of the object-oriented programming(oop), which is optional for you in this course, but I think it is very important for you to learn how to code in an oop way.

Most packages are programmed in this way, if you are familiar with this oop, you will have a deeper understanding of them. And then, you may write your own package.

Both the names and the dimensions of matrices and arrays are stored in R as attributes of the object. These attributes can be seen as labeled values you can attach to any object.

You can define and extract an attribute by the function `attr()`.

```

1 > set.seed(2005)
2 > samples <- rnorm(10, mean = 0, sd = 10)
3
4 # define an attribute and assign a value to it
5 > attr(samples, 'disttype') <- 'norm'
6 > print(samples)
7 [1]  9.640380 12.689196  4.731471  5.144074 -4.108097  0.694464 -2.861926  2.512387
8 [9]  9.600874 -14.079819
9 attr(,"disttype")
10 [1] "norm"
11
12 # one way to nicely print the structure of a list (or other data object) is 'str'
13 # emphasize on the structure
14 > str(samples)
15 num [1:10] 9.64 12.69 4.73 5.14 -4.11 ...
16 - attr(*, "disttype")= chr "norm"
17
18 # extract the attribute
19 > attr(samples, "disttype")
20 [1] "norm"
21
22 # another way, attributes() will return a list which contains all the attributes
23 > attributes(samples)$disttype
24 [1] "norm"

```

Many R functions return sorts of complex attribute-laden lists. Consider the `t.test()` function, which compares the means of two vectors for statistical equality(Hypothesis test):

```

1 set.seed(2005)
2 samp1 <- rnorm(100,mean = 10,sd = 5)
3 samp2 <- rnorm(100,mean = 8,sd = 5)
4 tresult <- t.test(samp1, samp2)
5
6 # it belongs to 'htest' class
7 > class(tresult)
8 [1] "htest"
9
10 # When printed, the result is a nicely formatted, human-readable result.
11 > print(tresult)
12
13      Welch Two Sample t-test
14
15 data:  samp1 and samp2
16 t = 1.0739, df = 197.79, p-value = 0.2842
17 alternative hypothesis: true difference in means is not equal to 0
18 95 percent confidence interval:
19  -0.5789316  1.9634773
20 sample estimates:
21 mean of x mean of y
22  9.319114  8.626841

```

If we run `str(tresult)`, however, we will find the true nature of result: it's a list.

```

1 > str(tresult)
2 List of 10
3 $ statistic : Named num 1.07
4 ..- attr(*, "names")= chr "t"
5 $ parameter : Named num 198

```

```

6  ..- attr(*, "names")= chr "df"
7  $ p.value      : num 0.284
8  $ conf.int    : num [1:2] -0.579 1.963
9  ..- attr(*, "conf.level")= num 0.95
10 $ estimate     : Named num [1:2] 9.32 8.63
11 ..- attr(*, "names")= chr [1:2] "mean of x" "mean of y"
12 $ null.value  : Named num 0
13 ..- attr(*, "names")= chr "difference in means"
14 $ stderr     : num 0.645
15 $ alternative: chr "two.sided"
16 $ method      : chr "Welch Two Sample t-test"
17 $ data.name   : chr "samp1 and samp2"
18 - attr(*, "class")= chr "htest"

```

Given knowledge of this structure, we can easily extract specific elements, such as the confidence level:

```

1 > attr(tresult$conf.int, 'conf.level')
2 [1] 0.95

```

Actually, the nice formatted result for `print` is because the class of `tresult` is `'htest'`.

```

1 > class(tresult) # use the 'class' function
2 [1] "htest"
3 > attr(tresult, 'class') # get the attribute's value
4 [1] "htest"
5
6 > attr(tresult, 'class')<-NULL # delete this attribute
7 > class(tresult) # back to the original data object
8 [1] "list"
9
10 > print(tresult) # print as a list instead of a nice formatted result
11 $statistic
12      t
13 1.073929
14
15 $parameter
16      df
17 197.7907
18
19 $p.value
20 [1] 0.284163
21
22 $conf.int
23 [1] -0.5789316  1.9634773
24 attr(,"conf.level")
25 [1] 0.95
26
27 $estimate
28 mean of x mean of y
29  9.319114  8.626841
30
31 $null.value
32 difference in means
33                0
34
35 $stderr
36 [1] 0.6446171
37
38 $alternative
39 [1] "two.sided"
40
41 $method
42 [1] "Welch Two Sample t-test"
43
44 $data.name
45 [1] "samp1 and samp2"

```

The interpreter notices that the "class" attribute is set to "htest", and searches for another function with a different name to actually run: `print.htest()`. You can check it in `methods(print)`.

I will give more detail in the next tutorial.