

1

Introduction to R/Python Programming

THERE has been considerable debate over choosing R vs. Python for Data Science. I started to learn Python when I was an undergraduate in 2006. At that time I never heard of Data Science. Five years later I read an R script for the first time. In my opinion, both R and Python are great languages and are worth learning; so why not learn them together?

In this Chapter, I would give an introduction on general R and Python programming, in a parallel fashion.

1.1 Calculator

R and Python are general-purpose programming languages that can be used for writing softwares in a variety of domains. But for now, let us start from using them as basic calculators. The first thing is to have them installed. R ¹ and Python ² can be downloaded from their official website. In this book, I would keep using R 3.5 and Python 3.7.

To use R/Python as basic calculators, let's get familiar with the interactive mode. After the installation, we can type R or Python (it is case insensitive so we can also type r/python) to invoke the interactive mode. Since Python 2 is installed by default on many machines, in order to avoid invoking Python 2 we type python3.7 instead.

R

```
1
2 ~ $R
3
4 R version 3.5.1 (2018-07-02) — "Feather Spray"
5 Copyright (C) 2018 The R Foundation for Statistical Computing
6 Platform: x86_64-apple-darwin15.6.0 (64-bit)
7
8 R is free software and comes with ABSOLUTELY NO WARRANTY.
9 You are welcome to redistribute it under certain conditions.
10 Type 'license()' or 'licence()' for distribution details.
```

¹ <https://www.r-project.org>

² <https://www.python.org>

```

11
12 Natural language support but running in an English locale
13
14 R is a collaborative project with many contributors.
15 Type 'contributors()' for more information and
16 'citation()' on how to cite R or R packages in publications.
17
18 Type 'demo()' for some demos, 'help()' for on-line help, or
19 'help.start()' for an HTML browser interface to help.
20 Type 'q()' to quit R.
21
22 >

```

Python

```

1 ~ $python3.7
2 Python 3.7.1 (default, Nov 6 2018, 18:45:35)
3 [Clang 10.0.0 (clang-1000.11.45.5)] on darwin
4 Type "help", "copyright", "credits" or "license" for more information.
5 >>>

```

The messages displayed by invoking the interactive mode depend on both the version of R/Python installed and the machine. Thus, you may see different messages on your local machine. As the messages said, to quit R we can type `q()`. There are 3 options prompted by asking the user if the workspace should be saved or not. Since we just want to use R as a basic calculator, we quit without saving workspace.

To quit Python, we can simply type `exit()`.

R

```

1 > q()
2 Save workspace image? [y/n/c]: n
3 ~ $

```

Once we are inside the interactive mode, we can use R/Python as a calculator.

R

```

1 > 1+1
2 [1] 2
3 > 2*3+5
4 [1] 11
5 > log(2)
6 [1] 0.6931472
7 > exp(0)

```

```
8 [1] 1
```

Python

```
1 >>> 1+1
2 2
3 >>> 2*3+5
4 11
5 >>> log(2)
6 Traceback (most recent call last):
7   File "<stdin>", line 1, in <module>
8   NameError: name 'log' is not defined
9 >>> exp(0)
10 Traceback (most recent call last):
11   File "<stdin>", line 1, in <module>
12   NameError: name 'exp' is not defined
```

From the code snippet above, R is working as a calculator perfectly. However, errors are raised when we call `log(2)` and `exp(2)` in Python. The error messages are self-explanatory - `log` function and `exp` function don't exist in the current Python environment. In fact, `log` function and `exp` function are defined in the `math` module in Python. A module³ is a file consisting of Python code. When we invoke the interactive mode of Python, a few built-in modules are loaded into the current environment by default. But the `math` module is not included in these built-in modules. That explains why we got the `NameError` when we try to use the functions defined in the `math` module. To resolve the issue, we should first load the functions to use by using the `import` statement as follows.

Python

```
1 >>> from math import log,exp
2 >>> log(2)
3 0.6931471805599453
4 >>> exp(0)
5 1.0
```

1.2 Variable and Type

In the previous section we have seen how to use R/Python as calculators. Now, let's see how to write real programs. First, let's define some variables.

³ <https://docs.python.org/3/tutorial/modules.html>

R

```

1 > a=2
2 > b=5.0
3 > x='hello world'
4 > a
5 [1] 2
6 > b
7 [1] 5
8 > x
9 [1] "hello world"
10 > e=a*2+b
11 > e
12 [1] 9

```

Python

```

1 >>> a=2
2 >>> b=5.0
3 >>> x='hello world'
4 >>> a
5 2
6 >>> b
7 5.0
8 >>> x
9 'hello world'
10 >>> e=a*2+b
11 >>> e
12 9.0

```

Here, we defined 4 different variables `a`, `b`, `x`, `e`. To get the type of each variable, we can utilize the function `typeof()` in R and `type()` in Python, respectively.

R

```

1 > typeof(x)
2 [1] "character"
3 > typeof(e)
4 [1] "double"

```

Python

```

1 >>> type(x)
2 <class 'str'>
3 >>> type(e)
4 <class 'float'>

```

The type of `x` in R is called character, and in Python is called str.

1.3 Functions

We have seen two functions `log` and `exp` when we use R/Python as calculators. A function is a block of code which performs a specific task. A major purpose of wrapping a block of code into a function is to reuse the code.

It is simple to define functions in R/Python.

R

```

1 > fun1=function(x){return(x*x)}
2 > fun1
3 function(x){return(x*x)}
4 > fun1(2)
5 [1] 4

```

Python

```

1 >>> def fun1(x):
2 ...     return x*x # note the
3 ...     indentation
4 >>> fun1(2)
5 4

```

Here, we defined a function `fun1` in R/Python. This function takes `x` as input and returns the square of `x`. When we call a function, we simply type the function name followed by the input argument inside a pair of parentheses. It is worth noting that input or output are not required to define a function. For example, we can define a function `fun2` to print `Hello World!` without input and output.

One major difference between R and Python codes is that Python codes are structured with indentation. Each logical line of R/Python code belongs to a certain group. In R, we use `{}` to determine the grouping of statements. However, in Python we use leading whitespace (spaces and tabs) at the beginning of a logical line to compute the indentation level of the line, which is used to determine the statements' grouping. Let's see what happens if we remove the leading whitespace in the Python function above.

Python

```
1 >>> def fun1(x):
2 ...   return x*x # note the indentation
3 File "<stdin>", line 2
4     return x*x # note the indentation
5         ^
6 IndentationError: expected an indented block
```

We got an `IndentationError` because of missing indentation.

R

```
1 > fun2=function(){print('Hello
   World!')}
2 > fun2()
3 [1] "Hello World!"
```

Python

```
1 >>> def fun2(): print('Hello World
   !')
2 ...
3 >>> fun2()
4 Hello World!
```

Let's go back to `fun1` and have a closer look at the `return`. In Python, if we want to return something we have to use the keyword `return` explicitly. `return` in R is a function but it is not a function in Python and that is why no parenthesis follows `return` in Python. In R, `return` is not required even though we need to return something from the function. Instead, we can just put the variables to return in the last line of the function defined in R. That being said, we can define `fun1` as follows.

R

```
1 > fun1=function(x){x*x}
```

1.4 Control flows

To implement a complex logic in R/Python, we may need control flows.

IF/ELSE

Let's define a function to return the absolute value of input.

R

```

1 > fun3=function(x){
2 +   if (x>=0){
3 +     return(x)}
4 +   else{
5 +     return(-x)}
6 + }
7 > fun3(2.5)
8 [1] 2.5
9 > fun3(-2.5)
10 [1] 2.5

```

Python

```

1 >>> def fun3(x):
2 ...   if x>=0:
3 ...     return x
4 ...   else:
5 ...     return -x
6 ...
7 >>> fun3(2.5)
8 2.5
9 >>> fun3(-2.5)
10 2.5

```

The code snippet above shows how to use **if/else** in R/Python. The subtle difference between R and Python is that the condition after **if** must be embraced by parenthesis in R but it is optional in Python.

We can also put **if** after **else**. But in Python, we use **elif** as a shortcut.

R

```

1 > fun4=function(x){
2 +   if (x==0){
3 +     print('zero')}
4 +   else if (x>0){
5 +     print('positive')}
6 +   else{
7 +     print('negative')}
8 + }
9 > fun4(0)
10 [1] "zero"
11 > fun4(1)
12 [1] "positive"
13 > fun4(-1)
14 [1] "negative"

```

Python

```

1 >>> def fun4(x):
2 ...   if x==0:
3 ...     print('zero')
4 ...   elif x>0:
5 ...     print('positive')
6 ...   else:
7 ...     print('negative')
8 ...
9 >>> fun4(0)
10 zero
11 >>> fun4(1)
12 positive
13 >>> fun4(-1)
14 negative

```

FOR LOOP

Similar to the usage of **if** in R, we also have to use parenthesis after the keyword **for** in R. But in Python there should be no parenthesis after **for**.

R

```

1 > for (i in 1:3){print(i)}
2 [1] 1
3 [1] 2
4 [1] 3

```

Python

```

1 >>> for i in range(1,4):print(i)
2 ...
3 1
4 2
5 3

```

There is something more interesting than the `for` loop itself in the snippets above. In the R code, the expression `1:3` creates a vector with elements 1,2 and 3. In the Python code, we use the `range()` function for the first time. Let's have a look at the type of them.

R

```

1 > typeof(1:3)
2 [1] "integer"

```

Python

```

1 >>> type(range(1,4))
2 <class 'range'>

```

`range()` function returns a `range` type object, which represents an immutable sequence of numbers. `range()` function can take three arguments, i.e., `range(start, stop, step)`. However, `start` and `step` are both optional. It's critical to keep in mind that the `stop` argument that defines the upper limit of the sequence is exclusive. And that is why in order to loop through 1 to 3 we have to pass 4 as the `stop` argument to `range()` function. The `step` argument specifies how much to increase from one number to the next. The default values of `start` and `step` are 0 and 1, respectively.

WHILE LOOP

R

```

1 > i=1
2 > while (i<=3){
3 +   print(i)
4 +   i=i+1
5 + }
6 [1] 1
7 [1] 2
8 [1] 3

```

Python

```

1 >>> i=1
2 >>> while i<=3:
3 ...   print(i)
4 ...   i+=1
5 ...
6 1
7 2
8 3

```

You may have noticed that in Python we can do `i+=1` to add 1 to `i`, which is not feasible in R by default. Both `for` loop and `while` loop can be nested.

BREAK/CONTINUE

`Break/continue` helps if we want to break the `for/while` loop earlier, or to skip a specific iteration. In R,

the keyword for continue is called `next`, in contrast to `continue` in Python. The difference between `break` and `continue` is that calling `break` would exit the innermost loop (when there are nested loops, only the innermost loop is affected); while calling `continue` would just skip the current iteration and continue the loop if not finished.

R

```

1 > for (i in 1:3){
2 +   print(i)
3 +   if (i==1) break
4 + }
5 [1] 1
6 > for (i in 1:3){
7 +   if (i==2){next}
8 +   print(i)
9 + }
10 [1] 1
11 [1] 3

```

Python

```

1 >>> for i in range(1,4):
2 ...   print(i)
3 ...   if i==1: break
4 ...
5 1
6 >>> for i in range(1,4):
7 ...   if i==2: continue
8 ...   print(i)
9 ...
10 1
11 3

```

1.5 Some built-in data structures

In the previous sections, we haven't seen much difference between R and Python. However, regarding the built-in data structures, there are some significant differences we would see in this section.

VECTOR IN R AND LIST IN PYTHON

In R, we can use function `c()` to create a Vector; A vector is a sequence of elements with the same type. In Python, we can use `[]` to create a list, which is also a sequence of elements. But the elements in a list don't need to have the same type. To get the number of elements in a Vector in R, we use the function `length()`; and to get the number of elements in a list in Python, we use the function `len()`.

R

```

1 > x=c(1,2,5,6)
2 > y=c('hello', 'world', '!')
3 > x
4 [1] 1 2 5 6
5 > y
6 [1] "hello" "world" "!"
7 > length(x)
8 [1] 4
9 > z=c(1, 'hello')
10 > z
11 [1] "1"      "hello"

```

Python

```

1 >>> x=[1,2,5,6]
2 >>> y=['hello', 'world', '!']
3 >>> x
4 [1, 2, 5, 6]
5 >>> y
6 ['hello', 'world', '!']
7 >>> len(x)
8 4
9 >>> z=[1, 'hello']
10 >>> z
11 [1, 'hello']

```


In the code snippet above, the first element in the variable `z` in R is coerced from `1` (numeric) to `"1"` (character) since the elements must have the same type.

To access a specific element from a Vector or list, we could use `[]`. In R, sequence types are indexed beginning with the one subscript; In contrast, sequence types in Python are indexed beginning with the zero subscript.

R

```
1 > x=c(1,2,5,6)
2 > x[1]
3 [1] 1
```

Python

```
1 >>> x=[1,2,5,6]
2 >>> x[1]
3 2
4 >>> x[0]
5 1
```

What if the index to access is out of boundary?

R

```
1 > x=c(1,2,5,6)
2 > x[-1]
3 [1] 2 5 6
4 > x[0]
5 numeric(0)
6 > x[length(x)+1]
7 [1] NA
8 > length(numeric(0))
9 [1] 0
10 > length(NA)
11 [1] 1
```

Python

```
1 >>> x=[1,2,5,6]
2 >>> x[-1]
3 6
4 >>> x[len(x)+1]
5 Traceback (most recent call last):
6   File "<stdin>", line 1, in <
    module>
7 IndexError: list index out of
    range
```

In Python, negative index number means indexing from the end of the list. Thus, `x[-1]` points to the last element and `x[-2]` points to the second-last element of the list. But R doesn't support indexing with negative number in the same way as Python. Specifically, in R `x[-index]` returns a new Vector with `x[index]` excluded.

When we try to access with an index out of boundary, Python would throw an `IndexError`. The behavior of R when indexing out of boundary is more interesting. First, when we try to access `x[0]` in R we get a `numeric(0)` whose length is also 0. Since its length is 0, `numeric(0)` can be interpreted as an empty numeric vector. When we try to access `x[length(x)+1]` we get a `NA`. In R, there are also `NaN` and `NULL`.

`NaN` means "Not A Number" and it can be verified by checking its type - "double". `0/0` would result in a `NaN` in R. `NA` in R generally represents missing values. And `NULL` represents a NULL (empty) object. To check if a value is `NA`, `NaN` or `NULL`, we can use `is.na()`, `is.nan()` or `is.null()`, respectively.

R

```

1 > typeof(NA)
2 [1] "logical"
3 > typeof(NaN)
4 [1] "double"
5 > typeof(NULL)
6 [1] "NULL"
7 > is.na(NA)
8 [1] TRUE
9 > is.null(NULL)
10 [1] TRUE
11 > is.nan(NaN)

```

Python

```

1 >>> type(None)
2 <class 'NoneType'>
3 >>> None is None
4 True
5 >>> 1 == None
6 False

```

In Python, there is no built-in NA or NaN. The counterpart of **NULL** in Python is **None**. In Python, we can use the **is** keyword or **==** to check if a value is equal to **None**.

From the code snippet above, we also notice that in R the boolean type value is written as "TRUE/FALSE", compared with "True/False" in Python. Although in R "TRUE/FALSE" can also be abbreviated as "T/F", I don't recommend to use the abbreviation.

There is one interesting fact that we can't add a **NULL** to a Vector in R, but it is feasible to add a **None** to a list in Python.

R

```

1 > x=c(1, NA, NaN, NULL)
2 > x
3 [1] 1 NA NaN
4 > length(x)
5 [1] 3

```

Python

```

1 >>> x=[1, None]
2 >>> x
3 [1, None]
4 >>> len(x)
5 2

```

Beside accessing a specific element from a Vector/list, we may also need to do slicing, i.e., to select a subset of the Vector/list. There are two basic approaches of slicing:

- Integer-based

R

```

1 > x=c(1,2,3,4,5,6)
2 > x[2:4]
3 [1] 2 3 4
4 > x[c(1,2,5)] # a Vector of indices
5 [1] 1 2 5

```

```

6 > x[seq(1,5,2)] # seq creates a Vector to be used as indices
7 [1] 1 3 5

```

Python

```

1 >>> x=[1,2,3,4,5,6]
2 >>> x[1:4] # x[start:end] start is inclusive but end is exclusive
3 [2, 3, 4]
4 >>> x[0:5:2] # x[start:end:step]
5 [1, 3, 5]

```

The code snippet above uses hash character `#` for comments in both R and Python. Everything after `#` on the same line would be treated as comment (not executable). In the R code, we also used the function `seq()` to create a Vector. When I see a function that I haven't seen before, I might either google it or use the builtin helper mechanism. Specifically, in R use `?` and in Python use `help()`.

R

```

1 > ?seq

```

Python

```

1 >>> help(print)

```

- Condition-based

Condition-based slicing means to select a subset of the elements which satisfy certain conditions. In R, it is quite straightforward by using a boolean Vector whose length is the same as the Vector to slice.

R

```

1 > x=c(1,2,5,5,6,6)
2 > x[x %% 2==1] # %% is the modulo operator in R; we select the odd elements
3 [1] 1 5 5
4 > x %% 2==1 # results in a boolean Vector with the same length as x
5 [1] TRUE FALSE TRUE TRUE FALSE FALSE

```

The condition-based slicing in Python is quite different from that in R. The prerequisite is list comprehension which provides a concise way to create new lists in Python. For example, let's create a list of squares of another list.

Python

```

1 >>> x=[1,2,5,5,6,6]
2 >>> [e**2 for e in x] # ** is the exponent operator, i.e., x**y means x to
   the power of y
3 [1, 4, 25, 25, 36, 36]

```

We can also use `if` statement with list comprehension to filter a list to achieve list slicing.

Python

```

1 >>> x=[1,2,5,5,6,6]
2 >>> [e for e in x if e%2==1] # % is the modulo operator in Python
3 [1, 5, 5]

```

It is also common to use `if/else` with list comprehension to achieve more complex operations. For example, given a list `x`, let's create a new list `y` so that the non-negative elements in `x` are squared and the negative elements are replaced by 0s.

Python

```

1 >>> x=[1,-1,0,2,5,-3]
2 >>> [e**2 if e>=0 else 0 for e in x]
3 [1, 0, 0, 4, 25, 0]

```

The example above shows the power of list comprehension. To use `if` with list comprehension, the `if` statement should be placed in the end after the `for` loop statement; but to use `if/else` with list comprehension, the `if/else` statement should be placed before the `for` loop statement.

We can also modify the value of an element in a Vector/list variable.

R

```

1 > x=c(1,2,3)
2 > x[1]=-1
3 > x
4 [1] -1  2  3

```

Python

```

1 >>> x=[1,2,3]
2 >>> x[0]=-1
3 >>> x
4 [-1, 2, 3]

```

Although the Vector structure in R and the list structure in Python looks similar regarding their usages and purposes, the implementation of these two structures are essentially different. The Vector structure is actually immutable, compared with the mutable list structure in Python. A mutable object can be changed after it is created, but an immutable object can't.

You may have the question - If a Vector object can't be changed, why it is feasible to change the value of a Vector in R? To answer the question, first we have to understand the relation between variable and object in R/Python. A variable itself is a reference or pointer to an object (usually stored in the machine's memory). In R, when we try to modify a Vector, the original Vector is not modified at all. Instead, the variable to that Vector is linked to a new Vector object to reflect the modification. Let's verify it by checking the memory address.

R

```

1 > x=1:1000
2 > tracemem(x) # print the memory
   address of x whenever the
   internal code copies the object
3 [1] "<0x7f840927a240>"
4 > x[1]=-x[1]
5 tracemem[0x7f840927a240 -> 0
   x7f8408b2ae00]:

```

Python

```

1 >>> x=list(range(1,1001)) # list()
   convert a range object to a
   list
2 >>> hex(id(x)) # print the memory
   address of x
3 '0x10592d908'
4 >>> x[0]=-x[0]
5 >>> hex(id(x))
6 '0x10592d908'

```

From the code snippet above, when we try to modify the value of `x[1]` in R, `x` points to a new object at address `0x7f8408b2ae00` (you probably would see different addresses on your machine). But in Python, the memory address doesn't change after we change the value of the first element because list in Python is mutable. Two or multiple Vectors/lists can be concatenated easily.

R

```

1 > x=c(1,2)
2 > y=c(3,4)
3 > z=c(5,6,7,8)
4 > c(x,y,z)
5 [1] 1 2 3 4 5 6 7 8

```

Python

```

1 >>> x=[1,2]
2 >>> y=[3,4]
3 >>> z=[5,6,7,8]
4 >>> x+y+z
5 [1, 2, 3, 4, 5, 6, 7, 8]

```

As the list structure in Python is mutable, there are many things we can do with list.

Python

```

1 >>> x=[1,2,3]
2 >>> x.append(4) # append a single value to the list x
3 >>> x
4 [1, 2, 3, 4]
5 >>> y=[5,6]
6 >>> x.extend(y) # extend list y to x
7 >>> x
8 [1, 2, 3, 4, 5, 6]
9 >>> last=x.pop() # pop the last element from x
10 >>> last
11 6
12 >>> x
13 [1, 2, 3, 4, 5]

```

Is there any immutable data structure in Python? Yes, for example tuple is immutable, which contains a sequence of elements. The element accessing and subset slicing of tuple is following the same rules of list in Python.

Python

```

1 >>> x=(1,2,3,) # use () to create a tuple in Python, it is better to always
    put a comma in the end
2 >>> type(x)
3 <class 'tuple'>
4 >>> len(x)
5 3
6 >>> x[0]
7 1
8 >>> x[0]=-1
9 Traceback (most recent call last):
10   File "<stdin>", line 1, in <module>
11 TypeError: 'tuple' object does not support item assignment

```

I like the list structure in Python much more than the Vector structure in R. list in Python has a lot more useful features which can be found from the python official documentation ⁴.

ARRAY

Array is one of the most important data structures in scientific programming. In R, there is also an object type "matrix", but according to my own experience, we can almost ignore its existence and use array instead. We can definitely use list as array in Python, but lots of linear algebra operations are not supported for the list type. Fortunately, there is a Python package *numpy* off the shelf.

R

```

1 > x=1:12
2 > array1=array(x,c(4,3)) # convert Vector x to a 4 rows * 3 cols array
3 > array1
4      [,1] [,2] [,3]
5 [1,]    1    5    9
6 [2,]    2    6   10
7 [3,]    3    7   11
8 [4,]    4    8   12
9 > y=1:6
10 > array2=array(y,c(3,2)) # convert Vector y to a 3 rows * 2 cols array
11 > array2
12      [,1] [,2]
13 [1,]    1    4

```

⁴ <https://docs.python.org/3/tutorial/datastructures.html>

```

14 [2,]    2    5
15 [3,]    3    6
16 > array3 = array1 %*% array2 # %*% is the matrix multiplication operator
17 > array3
18      [,1] [,2]
19 [1,]   38   83
20 [2,]   44   98
21 [3,]   50  113
22 [4,]   56  128
23 > dim(array3) # get the dimension of array3
24 [1] 4 2

```

Python

```

1 >>> import numpy as np # we import the numpy module and alias it as np
2 >>> array1=np.reshape(list(range(1,13)),(4,3)) # convert a list to a 2d np.
      array
3 >>> array1
4 array([[ 1,  2,  3],
5        [ 4,  5,  6],
6        [ 7,  8,  9],
7        [10, 11, 12]])
8 >>> type(array1)
9 <class 'numpy.ndarray'>
10 >>> array2=np.reshape(list(range(1,7)),(3,2))
11 >>> array2
12 array([[1, 2],
13        [3, 4],
14        [5, 6]])
15 >>> array3=np.dot(array1,array2) # matrix multiplication using np.dot()
16 >>> array3
17 array([[ 22,  28],
18        [ 49,  64],
19        [ 76, 100],
20        [103, 136]])
21 >>> array3.shape # get the shape(dimension) of array3
22 (4, 2)

```

You may have noticed that the results of the R code snippet and Python code snippet are different. The reason is that in R the conversion from a Vector to an array is by-column; but in numpy the reshape from a list to an 2D numpy.array is by-row. There are two ways to reshape a list to a 2D numpy.array by column.

Python

```

1 >>> array1=np.reshape(list(range(1,13)),(4,3),order='F') # use order='F'
2 >>> array1
3 array([[ 1,  5,  9],
4        [ 2,  6, 10],
5        [ 3,  7, 11],
6        [ 4,  8, 12]])
7 >>> array2=np.reshape(list(range(1,7)),(2,3)).T # use .T to transpose an array
8 >>> array2
9 array([[1, 4],
10        [2, 5],
11        [3, 6]])
12 >>> np.dot(array1,array2) # now we get the same result as using R
13 array([[ 38,  83],
14        [ 44,  98],
15        [ 50, 113],
16        [ 56, 128]])

```

To learn more about numpy, the official website ⁵. has great documentation/tutorials.

LIST IN R AND DICTIONARY IN PYTHON

Yes, in R there is also an object type called list. The major difference between a Vector and a list in R is that a list could contain different types of elements. list in R supports integer-based accessing using [[]] (compared to [] for Vector).

R

```

1 > x=list(1, 'hello world!')
2 > x
3 [[1]]
4 [1] 1
5
6 [[2]]
7 [1] "hello world!"
8
9 > x[[1]]
10 [1] 1
11 > x[[2]]
12 [1] "hello world!"
13 > length(x)
14 [1] 2

```

list in R is also immutable, like Vector. list in R could be named and support accessing by name via either [[]] or \$ operator. But Vector in R can also be named and support accessing by name.

⁵<http://www.numpy.org>

R

```

1 > x=c('a'=1, 'b'=2)
2 > names(x)
3 [1] "a" "b"
4 > x['b']
5 b
6 2
7 > l=list('a'=1, 'b'=2)
8 > l[['b']]
9 [1] 2
10 > l$b
11 [1] 2
12 > names(l)
13 [1] "a" "b"

```

However, elements in list in Python can't be named as R. If we need the feature of accessing by name in Python, we can use the dictionary structure. If you used Java before, you may consider dictionary in Python as the counterpart of HashMap in Java. Essentially, a dictionary in Python is a collection of key:value pairs.

Python

```

1 >>> x={'a':1, 'b':2} # {key:value} pairs
2 >>> x
3 {'a': 1, 'b': 2}
4 >>> x['a']
5 1
6 >>> x['b']
7 2
8 >>> len(x) # number of key:value pairs
9 2
10 >>> x.pop('a') # remove the key 'a' and we get its value 1
11 1
12 >>> x
13 {'b': 2}

```

Unlike dictionary in Python, list in R doesn't support the pop() operation (also because list in R is immutable). Thus, in order to modify a list in R, a new one would be created explicitly or implicitly.

DATA.FRAME, DATA.TABLE AND PANDAS

data.frame is a built-in type in R for data manipulation. data.table⁶ is an R package and pandas⁷ is a Python module. To use data.table and pandas, we have to install the them. Although data.frame is a built-in type, it is not quite efficient for many operations. I would suggest to use data.table whenever possible. dplyr

⁶ <https://cran.r-project.org/web/packages/data.table/index.html>

⁷ <https://pandas.pydata.org/>

⁸ is also a very popular package in R for data manipulation. But I favor `data.table` over `dplyr`. Many good online resources are available online to learn `data.table` and `pandas`. Thus, I would not cover the usage of these tools for now.

OBJECT-ORIENTED PROGRAMMING (OOP) IN R/PYTHON

All the codes we wrote above follow the procedural programming paradigm ⁹. We can also do functional programming (FP) and OOP in R/Python. R is more like a functional programming language than Python as most of the built-in data structures in R are immutable. In this section, let's focus on OOP in R/Python.

Class is the key concept in OOP. In R there are two commonly used built-in systems to define classes, i.e., S3 and S4. In addition, there is an external package R6 ¹⁰ which defines R6 classes. S3 is a light-weight system but its style is quite different from OOP in many other programming languages. S4 system follows the principles of modern object oriented programming much better than S3. However, the usage of S4 classes is quite tedious. I would ignore S3/S4 and introduce R6, which is more close to the class in Python.

Let's build a class in R/Python to represent complex numbers.

R

```

1 > library(R6) # load the R6 package
2 >
3 > Complex = R6Class("Complex",
4 + public = list( # only elements declared in this list are accessible by the
5 +               # object of this class
6 +               real = NULL,
7 +               imag = NULL,
8 +               # the initialize function would be called automatically when we create an
9 +               # object of the class
10 +               initialize = function(real,imag){
11 +                 # call functions to change real and imag values
12 +                 self$set_real(real)
13 +                 self$set_imag(imag)
14 +               },
15 +               # define a function to change the real value
16 +               set_real = function(real){
17 +                 self$real=real
18 +               },
19 +               # define a function to change the imag value
20 +               set_imag = function(imag){
21 +                 self$imag=imag
22 +               },
23 +               # override print function
24 +               print = function() {
```

⁸ <https://dplyr.tidyverse.org/>

⁹ https://en.wikipedia.org/wiki/Comparison_of_programming_paradigms

¹⁰ <https://cran.r-project.org/web/packages/R6/index.html>

```

23 +   cat(paste0(as.character(self$real), '+', as.character(self$imag), 'j'), '\n')
24 + }
25 + )
26 + )
27 > # let's create a complex number object based on the Complex class we defined
    above using the new function
28 > x = Complex$new(1,2)
29 > x
30 1+2j
31 > x$real # the public attributes of x could be accessed by $ operator
32 [1] 1

```

Python

```

1 >>> class Complex:
2 ...     # the __init__ function would be called automatically when we create an
    object of the class
3 ...     def __init__(self, real, imag):
4 ...         self.real = None
5 ...         self.imag = None
6 ...         self.set_real(real)
7 ...         self.set_imag(imag)
8 ...     # define a function to change the real value
9 ...     def set_real(self, real):
10 ...         self.real=real
11 ...     # define a function to change the imag value
12 ...     def set_imag(self, imag):
13 ...         self.imag=imag
14 ...     def __repr__(self):
15 ...         return "{0}+{1}j".format(self.real, self.imag)
16 ...
17 >>> x = Complex(1,2)
18 >>> x
19 1+2j
20 >>> x.real # different from the $ operator in R, here we use . to access the
    attribute of an object
21 1

```

By overriding the print function in the R6 class, we can have the object printed in the format of real+imag j. To achieve the same effect in Python, we override the method `__repr__`. In Python, we call the functions defined in classes as methods. And overriding a method means changing the implementation of a method provided by one of its ancestors. To understand the concept of ancestors in OOP, one needs to understand

the concept of inheritance ¹¹.

You may be curious of the double underscore surrounding the methods, such as `__init__` and `__repr__`. These methods are well-known as magic methods ¹². Magic methods could be very handy if we use them in the suitable cases. For example, we can use the magic method `__add__` to implement the `+` operator for the `Complex` class we defined above.

Python

```

1 >>> class Complex:
2 ...     def __init__(self, real, imag):
3 ...         self.real = None
4 ...         self.imag = None
5 ...         self.set_real(real)
6 ...         self.set_imag(imag)
7 ...     def set_real(self, real):
8 ...         self.real=real
9 ...     def set_imag(self, imag):
10 ...         self.imag=imag
11 ...     def __repr__(self):
12 ...         return "{0}+{1}j".format(self.real, self.imag)
13 ...     def __add__(self, another):
14 ...         return Complex(self.real+another.real, self.imag+another.imag)
15 ...
16 >>> x=Complex(1,2)
17 >>> y=Complex(2,4)
18 >>> x+y # + operator works now
19 3+6j

```

We can also implement the `+` operator for `Complex` class in R like what we have done for Python.

R

```

1 > `+.Complex` = function(x,y){
2 +   Complex$new(x$real+y$real, x$imag+y$imag)
3 + }
4 > x=Complex$new(1,2)
5 > y=Complex$new(2,4)
6 > x+y
7 3+6j

```

The most interesting part of the code above is ``+.Complex``. First, why do we use `` to quote the function name? Before getting into this question, let's have a look at the Python 3's variable naming rules ¹³.

¹¹ [https://en.wikipedia.org/wiki/Inheritance_\(object-oriented_programming\)](https://en.wikipedia.org/wiki/Inheritance_(object-oriented_programming))

¹² <https://rszalski.github.io/magicmethods/>

¹³ https://docs.python.org/3.3/reference/lexical_analysis.html

¹ Within the ASCII `range` (U+0001..U+007F), the valid characters `for` identifiers (also referred to `as` names) are the same `as in` Python 2.x: the uppercase `and` lowercase letters A through Z, the underscore `_` `and`, `except for` the first character, the digits 0 through 9.

According to the rule, we can't declare a variable with name `2x`. Compared with Python, in R we can also use `.` in the variable names ¹⁴. However, there is a workaround to use invalid variable names in R with the help of ```.

R

```

1 > 2x = 5
2 Error: unexpected symbol in "2x"
3 > .x = 3
4 > .x
5 [1] 3
6 > `+2x%` = 0
7 > `+2x%`
8 [1] 0

```

Python

```

1 >>> 2x = 5
2     File "<stdin>", line 1
3         2x = 5
4         ^
5 SyntaxError: invalid syntax
6 >>> .x = 3
7     File "<stdin>", line 1
8         .x = 3
9         ^
10 SyntaxError: invalid syntax

```

Now it is clear the usage of ``` in ``+.Complex`` is to define a function with invalid name. Placing `.Complex` after `+` is related to S3 method dispatching which would not be discussed here.

1.6 Miscellaneous

There are some items that I haven't discussed so far, which are also important in order to master R/Python.

PACKAGE/MODULE INSTALLATION

- Use `install.packages()` function in R
- Use R IDE to install packages
- Use `pip` ¹⁵ to install modules in Python

VIRTUAL ENVIRONMENT

Virtual environment is a tool to manage dependencies in Python. There are different ways to create virtual environments in Python. But I suggest to use the `venv` module shipped with Python 3. Unfortunately, there is nothing like a real virtual environment in R as far as I know although there quite a few of packages management tools/packages.

¹⁴ <https://cran.r-project.org/doc/manuals/r-release/R-lang.html#Identifiers>

¹⁵ <https://packaging.python.org/tutorials/installing-packages/>

<- vs. =

If you have known R before, you probably heard of the advice ¹⁶ to use `<-` to rather than `=` for value assignment. I always use `=` for value assignment. Let's see an example when `<-` makes a difference when we do value assignment.

R

```

1 > x=1
2 > a=list(x <- 2)
3 > a
4 [[1]]
5 [1] 2
6
7 > x
8 [1] 2

```

R

```

1 > x=1
2 > a=list(x = 2)
3 > a
4 $x
5 [1] 2
6
7 > x
8 [1] 1

```

¹⁶ <https://google.github.io/styleguide/Rguide.xml>

2

More on R/Python Programming