

Contents

1	<i>Introduction to R/Python Programming</i>	9
1.1	<i>Calculator</i>	9
1.2	<i>Variable and Type</i>	11
1.3	<i>Functions</i>	12
1.4	<i>Control flows</i>	13
1.5	<i>Some built-in data structures</i>	16
1.6	<i>Miscellaneous</i>	32
2	<i>More on R/Python Programming</i>	35
2.1	<i>Write & run R/Python scripts</i>	35
2.2	<i>Debugging in R/Python</i>	36
2.3	<i>Benchmarking</i>	43
2.4	<i>Vectorization</i>	46
2.5	<i>Embarrassingly parallelism in R/Python</i>	53
2.6	<i>Scope of Variables</i>	57
2.7	<i>Miscellaneous</i>	61
3	<i>data.table and pandas</i>	63
3.1	<i>SQL</i>	63
3.2	<i>Get started with data.table and pandas</i>	66
3.3	<i>Indexing & Selecting Data</i>	67
3.4	<i>Add/Remove/Update</i>	76
3.5	<i>Group by</i>	80
3.6	<i>Join</i>	82

4	<i>Linear Regression</i>	87
4.1	<i>Basics of linear regression</i>	87
4.2	<i>Linear hypothesis testing</i>	92
4.3	<i>Ridge regression</i>	93
5	<i>Optimization in Practice</i>	99
5.1	<i>Convexity</i>	99
5.2	<i>Gradient descent</i>	100
5.3	<i>Root-finding</i>	105
5.4	<i>General purpose minimization tools in R/Python</i>	107
5.5	<i>Linear programming</i>	114
5.6	<i>Miscellaneous</i>	116
6	<i>Predictive Modeling in Practice</i>	123

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 list structure in Python is mutable. A mutable object can be changed after it is created, but an immutable object can't. However, the mutability of vector in R is a bit of complicated. If we change the value of an element in a vector without changing the type of the element, the vector is mutable. If we change an element to another type, the behavior of the vector is immutable. A variable itself is a reference or pointer to an object (usually stored in the machine's memory). To check the mutability of a variable, we can trace the memory address.

R

```

1 > x=c(1:3)
2 > tracemem(x) # print the memory
  address of x whenever the
  address changes
3 [1] "<0x7ff360c95c08>"
4 > x[1]=-x[1] # type not changed, i
  .e., from integer to integer
5 > tracemem(x)
6 [1] "<0x7ff360c95c08>"
7 > x[1]=-1.0
8 tracemem[0x7ff360c95c08 -> 0
  x7ff3604692d8]:

```

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]=1.0 # from integer to
  float
5 >>> hex(id(x))
6 '0x10592d908'

```

From the code snippet above, in Python the memory address doesn't change after we change the value of the first element because list in Python is mutable. When we try to modify the value of `x[1]` in R, the memory address of `x` doesn't change. (you probably would see different addresses on your machine). But when we change the value of `x[1]` from the integer type to a double type, the memory address got changed. It's worth noting since R 3.5 arithmetic sequences created by `1:n`, `seq_along`, and the like now use compact internal representations via the ALTREP framework ⁴. Let's the example below.

R

```

1 > x=1:3
2 > tracemem(x)
3 [1] "<0x7f828e84c110>"
4 > .Internal(inspect(x))
5 @7f828e84c110 13 INTSXP g0c0 [NAM
  (3),TR] 1 : 3 (compact)
6 > x[1]=2L
7 tracemem[0x7f828e84c110 -> 0
  x7f828fe49848]:

```

R

```

1 > x=c(1:3)
2 > tracemem(x)
3 [1] "<0x7f828fe498c8>"
4 > .Internal(inspect(x))
5 @7f828fe498c8 13 INTSXP g0c2 [NAM
  (1),TR] (len=3, t1=0) 1,2,3
6 > x[1]=2L
7 > tracemem(x)
8 [1] "<0x7f828fe498c8>"

```

⁴ <https://cran.r-project.org/doc/manuals/r-devel/NEWS.html>

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 ⁵.

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
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

```

⁵ <https://docs.python.org/3/tutorial/datastructures.html>

```

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 a 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 ⁶ 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

```

The mutability of the list structure in R is similar to the vector structure. The difference is that when we change the type of an element in a list, the memory address doesn't change in general.

R

```

1 > x=list(c(1:3), 'Hello World!')
2 > tracemem(x)
3 [1] "<0x7f828fe497c8>"
4 > x[[1]]=1.0
5 > x[[2]]=2.0
6 > x
7 [[1]]
8 [1] 1
9
10 [[2]]
11 [1] 2
12

```

⁶<http://www.numpy.org>

```

13 > tracemem(x)
14 [1] "<0x7f828fe497c8>"

```

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.

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. 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. In Python, there is no such built-in data structure since Python is a more general-purpose programming language. The solution for data.frame in Python is the pandas ⁷ module.

Before we dive into data.frame, you may be curious why we need it? In other words, why can't we just use vector, list, array/matrix and dictionary for all data manipulation tasks? I would say yes - data.frame is not a must-have feature for most of ETL (extraction, transformation and Load) operations. But data.frame provides a very intuitive way for us to understand the structured data set. A data.frame is usually flat with 2 dimensions, i.e., row and column. The row dimension is across multiple observations and the column dimension is across multiple attributes/features. If you are familiar with relational database, a data.frame can be viewed as a table.

Let's see an example of using data.frame to represent employees' information in a company.

R

```

1 > employee_df = data.frame(name=c(
    "A", "B", "C"), department=c("
    Engineering", "Operations", "
    Sales"))
2 > employee_df
3   name department
4 1    A Engineering
5 2    B Operations
6 3    C      Sales

```

Python

```

1 >>> import pandas as pd
2 >>> employee_df=pd.DataFrame({'
    name':['A','B','C'],'department
   ':['Engineering','Operations',"
    Sales"]})
3 >>> employee_df
4   name department
5 0    A Engineering
6 1    B Operations
7 2    C      Sales

```

There are quite a few ways to create data.frame. The most commonly used one is to create data.frame object from array/matrix. We may also need to convert a numeric data.frame to an array/matrix.

R

```

1 > x=array(rnorm(12),c(3,4))
2 > x
3           [,1]      [,2]      [,3]      [,4]
4 [1,] -0.8101246 -0.8594136 -2.260810  0.5727590
5 [2,] -0.9175476  0.1345982  1.067628 -0.7643533
6 [3,]  0.7865971 -1.9046711 -0.154928 -0.6807527
7 > random_df=as.data.frame(x)
8 > random_df
9           V1           V2           V3           V4
10 1 -0.8101246 -0.8594136 -2.260810  0.5727590

```

⁷<https://pandas.pydata.org/>

```

11 2 -0.9175476  0.1345982  1.067628 -0.7643533
12 3  0.7865971 -1.9046711 -0.154928 -0.6807527
13 > data.matrix(random_df)
14           V1           V2           V3           V4
15 [1,] -0.8101246 -0.8594136 -2.260810  0.5727590
16 [2,] -0.9175476  0.1345982  1.067628 -0.7643533
17 [3,]  0.7865971 -1.9046711 -0.154928 -0.6807527

```

Python

```

1 >>> import numpy as np
2 >>> import pandas as pd
3 >>> x=np.random.normal(size=(3,4))
4 >>> x
5 array([[ -0.54164878, -0.14285267, -0.39835535, -0.81522719],
6        [ 0.01540508,  0.63556266,  0.16800583,  0.17594448],
7        [-1.21598262,  0.52860817, -0.61757696,  0.18445057]])
8 >>> random_df=pd.DataFrame(x)
9 >>> random_df
10           0           1           2           3
11 0 -0.541649 -0.142853 -0.398355 -0.815227
12 1  0.015405  0.635563  0.168006  0.175944
13 2 -1.215983  0.528608 -0.617577  0.184451
14 >>> np.asarray(random_df)
15 array([[ -0.54164878, -0.14285267, -0.39835535, -0.81522719],
16        [ 0.01540508,  0.63556266,  0.16800583,  0.17594448],
17        [-1.21598262,  0.52860817, -0.61757696,  0.18445057]])

```

In general, operations on an array/matrix is much faster than that on a data frame. In R, we may use the built-in function `data.matrix` to convert a data.frame to an array/matrix. In Python, we could use the function `asarray` in `numpy` module.

Although `data.frame` is a built-in type, it is not quite efficient for many operations. I would suggest to use `data.table`⁸ whenever possible. `dplyr`⁹ is also a very popular package in R for data manipulation. 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 paradigm¹⁰. We can also do functional programming (FP) and OOP in R/Python. 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,

⁸ <https://cran.r-project.org/web/packages/data.table/index.html>

⁹ <https://dplyr.tidyverse.org/>

¹⁰ https://en.wikipedia.org/wiki/Comparison_of_programming_paradigms

i.e., S3 and S4. In addition, there is an external package R6¹¹ 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
      object of this class
5 + real = NULL,
6 + imag = NULL,
7 + # the initialize function would be called automatically when we create an
      object of the class
8 + initialize = function(real,imag){
9 +     # call functions to change real and imag values
10 +     self$set_real(real)
11 +     self$set_imag(imag)
12 + },
13 + # define a function to change the real value
14 + set_real = function(real){
15 +     self$real=real
16 + },
17 + # define a function to change the imag value
18 + set_imag = function(imag){
19 +     self$imag=imag
20 + },
21 + # override print function
22 + print = function(){
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

```

¹¹ <https://cran.r-project.org/web/packages/R6/index.html>

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 ¹².

You may be curious of the double underscore surrounding the methods, such as `__init__` and `__repr__`. These methods are well-known as magic methods ¹³. 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.

In the definition of the magic method `__repr__` in the Python code, the `format` method of `str` object ¹⁴ is used.

¹² [https://en.wikipedia.org/wiki/Inheritance_\(object-oriented_programming\)](https://en.wikipedia.org/wiki/Inheritance_(object-oriented_programming))

¹³ <https://rszalski.github.io/magicmethods/>

¹⁴ <https://docs.python.org/3.7/library/string.html>

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 ¹⁵.

```

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

¹⁵ https://docs.python.org/3.3/reference/lexical_analysis.html

use `.` in the variable names ¹⁶. 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` ¹⁷ 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.

<- vs. =

If you have known R before, you probably heard of the advice ¹⁸ 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.

¹⁶ <https://cran.r-project.org/doc/manuals/r-release/R-lang.html#Identifiers>

¹⁷ <https://packaging.python.org/tutorials/installing-packages/>

¹⁸ <https://google.github.io/styleguide/Rguide.xml>

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
```

2

More on R/Python Programming

We have learned quite a few of basic R/Python programming in the previous chapter. I hope this chapter could be used as an intermediate level R/Python programming tutorial. There are a few topics to cover, including debugging, vectorization and some other useful features of R/Python.

2.1 Write & run R/Python scripts

In [chapter 1](#) we are coding within the interactive mode of R/Python. When we are working on the real world projects, using an Integrated development environment (IDE) is a more pragmatic choice. There are not many choices for R IDE, and among them RStudio ¹⁹ is the best one I have used so far. As for Python, I would recommend either Visual Studio Code ²⁰ or PyCharm ²¹. But of course, you could use any text editor to write R/Python scripts.

Let's write a simple script to print Hello World! in R/Python. I have made a directory `chapter2` on my disk, the R script is saved as `hello_world.R` and the Python script is saved as `hello_world.py`, inside the directory.

R

```
chapter2/hello_world.R
1 print("Hello World!")
```

Python

```
chapter2/hello_world.py
1 print("Hello World!")
```

There are a few ways to run the R script. For example, we can run the script from the console with the `r -f filename` command. Also, we can open the R interactive session and use the `source()` function. I would recommend the second approach with `source()` function. As for the Python script, we can run it from the console.

¹⁹ <https://www.rstudio.com/products/rstudio/>

²⁰ <https://www.jetbrains.com/pycharm/>

²¹ <https://www.jetbrains.com/pycharm/>

R

```

1 chapter2 $ls
2 hello_world.R hello_world.py
3 chapter2 $r -f hello_world.R
4 > print("Hello World!")
5 [1] "Hello World!"
6
7 chapter2 $r
8 > source('hello_world.R')
9 [1] "Hello World!"

```

Python

```

1 chapter2 $ls
2 hello_world.R hello_world.py
3
4 chapter2 $python3.7 hello_world.py
5 Hello World!

```

2.2 Debugging in R/Python

Debugging is one of the most important aspects of programming. What is debugging in programming? The programs we write might include errors/bugs and debugging is a step-by-step process to find and remove the errors/bugs in order to get the desired results.

If you are smart enough or the bugs are evident enough then you can debug the program on your mind without using a computer at all. But in general we need some tools/techniques to help us with debugging.

PRINT

Most of programming languages provide the functionality of printing, which is a natural way of debugging. By trying to place print statements at different positions we may finally catch the bugs. When I use `print` to debug, it's feeling like playing the game of minesweeper. In Python, there is a module called `logging`²² which could be used for debugging like the `print` function, but in a more elegant fashion.

BROWSER IN R AND PDB IN PYTHON

In R, there is a function `browser()` which interrupts the execution and allows the inspection of the current environment. Similarly, there is a module `pdb` in Python that provides more debugging features. We would only focus on the basic usages of `browser()` and the `set_trace()` function in `pdb` module. The essential difference between debugging using `print()` and `browser()` and `set_trace()` is that the latter functions allows us to debug in an interactive mode.

Let's write a function which takes a sorted vector/list `v` and a target value `x` as input and returns the leftmost index `pos` of the sorted vector/list so that `v[pos]>=x`. Since `v` is already sorted, we may simply loop through it from left to right to find `pos`.

²² <https://docs.python.org/3/library/logging.html>

R

chapter2/find_pos.R

```

1 find_pos=function(v,x){
2   for (i in 1:length(v)){
3     if (v[i]>=x){
4       return(i)
5     }
6   }
7 }
8
9 v=c(1,2,5,10)
10 print(find_pos(v,-1))
11 print(find_pos(v,4))
12 print(find_pos(v,11))

```

Python

chapter2/find_pos.py

```

1 def find_pos(v,x):
2   for i in range(len(v)):
3     if v[i]>=x:
4       return i
5
6 v=[1,2,5,10]
7 print(find_pos(v,-1))
8 print(find_pos(v,4))
9 print(find_pos(v,11))

```

Now let's run these two scripts.

R

```

1 chapter2 $r
2 > source('find_pos.R')
3 [1] 1
4 [1] 3
5 NULL

```

Python

```

1 chapter2 $python3.7 find_pos.py
2 0
3 2
4 None

```

When $x=11$, the function returns `NULL` in R and `None` in Python because there is no such element in v larger than x . The implementation above is trivial, but not efficient. If you have some background in data structures and algorithms, you probably know this question can be solved by binary search. The essential idea of binary search comes from Divide-and-conquer²³. Since v is already sorted, we may divide it into two partitions by cutting it from the middle, and then we get the left partition and the right partition. v is sorted implies that both the left partition and the right partition are also sorted. If the target value x is larger than the rightmost element in the left partition, we can just discard the left partition and search x within the right partition. Otherwise, we can discard the right partition and search x within the left partition. Once we have determined which partition to search, we may apply the idea recursively so that in each step we reduce the size of v by half. If the length of v is denoted as n , in terms of big O notation²⁴, the run time complexity of binary search is $\mathcal{O}(\log n)$, compared with $\mathcal{O}(n)$ of the for-loop implementation.

The code below implements the binary search solution to our question (It is more intuitive to do it with recursion but here I write it with iteration since tail recursion optimization²⁵ in R/Python is not supported).

²³ https://en.wikipedia.org/wiki/Divide-and-conquer_algorithm

²⁴ https://en.wikipedia.org/wiki/Big_O_notation

²⁵ https://en.wikipedia.org/wiki/Tail_call

R

chapter2/find_binary_search_buggy.R

```

1 binary_search_buggy=function(v,x){
2   start = 1
3   end = length(v)
4   while (start<end){
5     mid = (start+end) %% 2 # %% is the floor division operator
6     if (v[mid]>=x){
7       end = mid
8     }else{
9       start = mid+1
10    }
11  }
12  return(start)
13 }
14 v=c(1,2,5,10)
15 print(binary_search_buggy(v,-1))
16 print(binary_search_buggy(v,5))
17 print(binary_search_buggy(v,11))

```

Python

chapter2/find_binary_search_buggy.py

```

1 def binary_search_buggy(v,x):
2   start,end = 0,len(v)-1
3   while start<end:
4     mid = (start+end)//2 # // is the floor division operator
5     if v[mid]>=x:
6       end = mid
7     else:
8       start = mid+1
9   return start
10
11 v=[1,2,5,10]
12 print(binary_search_buggy(v,-1))
13 print(binary_search_buggy(v,5))
14 print(binary_search_buggy(v,11))

```

Now let's run these two `binary_search` scripts.

R

```
1 chapter2 $r
2 > source('binary_search_buggy.R')
3 [1] 1
4 [1] 3
5 [1] 4
```

Python

```
1 chapter2 $python3.7
   binary_search_buggy.py
2 0
3 2
4 3
```

The binary search solutions don't work as expected when `x=11`. We write two new scripts.

R

chapter2/find_binary_search_buggy_debug.R

```
1 binary_search_buggy=function(v,x){
2   browser()
3   start = 1
4   end = length(v)
5   while (start<end){
6     mid = (start+end)
7     if (v[mid]>=x){
8       end = mid
9     }else{
10      start = mid+1
11    }
12  }
13  return(start)
14 }
15 v=c(1,2,5,10)
16 print(binary_search_buggy(v,11))
```

Python

chapter2/find_binary_search_buggy_debug.py

```
1 from pdb import set_trace
2 def binary_search_buggy(v,x):
3   set_trace()
4   start,end = 0,len(v)-1
5   while start<end:
6     mid = (start+end)//2
7     if v[mid]>=x:
8       end = mid
```

```

9     else:
10         start = mid+1
11     return start
12
13 v=[1,2,5,10]
14 print(binary_search_buggy(v, 11))

```

Let's try to debug the programs with the help of `browser()` and `set_trace()`.

R

```

1 > source('binary_search_buggy_debug.R')
2 Called from: binary_search_buggy(v, 11)
3 Browse[1]> ls()
4 [1] "v" "x"
5 Browse[1]> n
6 debug at binary_search_buggy_debug.R#3: start = 1
7 Browse[2]> n
8 debug at binary_search_buggy_debug.R#4: end = length(v)
9 Browse[2]> n
10 debug at binary_search_buggy_debug.R#5: while (start < end) {
11     mid = (start + end)%/%2
12     if (v[mid] >= x) {
13         end = mid
14     }
15     else {
16         start = mid + 1
17     }
18 }
19 Browse[2]> n
20 debug at binary_search_buggy_debug.R#6: mid = (start + end)%/%2
21 Browse[2]> n
22 debug at binary_search_buggy_debug.R#7: if (v[mid] >= x) {
23     end = mid
24 } else {
25     start = mid + 1
26 }
27 Browse[2]> n
28 debug at binary_search_buggy_debug.R#10: start = mid + 1
29 Browse[2]> n
30 debug at binary_search_buggy_debug.R#5: (while) start < end
31 Browse[2]> n
32 debug at binary_search_buggy_debug.R#6: mid = (start + end)%/%2
33 Browse[2]> n

```

```

34 debug at binary_search_buggy_debug.R#7: if (v[mid] >= x) {
35     end = mid
36 } else {
37     start = mid + 1
38 }
39 Browse[2]> n
40 debug at binary_search_buggy_debug.R#10: start = mid + 1
41 Browse[2]> n
42 debug at binary_search_buggy_debug.R#5: (while) start < end
43 Browse[2]> start
44 [1] 4
45 Browse[2]> n
46 debug at binary_search_buggy_debug.R#13: return(start)
47 Browse[2]> n
48 [1] 4

```

In the R code snippet above, we placed the `browser()` function on the top of the function `binary_search_buggy`. Then when we call the function we enter into the debugging environment. By calling `ls()` we see all variables in the current debugging scope, i.e., `v`, `x`. Typing `n` will evaluate the next statement. After typing `n` a few times, we finally exit from the while loop because `start` = 4 such that `start` < `end` is FALSE. As a result, the function just returns the value of `start`, i.e., 4. To exit from the debugging environment, we can type `Q`; to continue the execution we can type `c`.

The root cause is that we didn't deal with the corner case when the target value `x` is larger than the last/largest element in `v` correctly.

Let's debug the Python function using `pdb` module.

Python

```

1 chapter2 $python3.7 binary_search_buggy_debug.py
2 > chapter2/binary_search_buggy_debug.py(4)binary_search_buggy()
3 -> start,end = 0,len(v)-1
4 (Pdb) n
5 > chapter2/binary_search_buggy_debug.py(5)binary_search_buggy()
6 -> while start<end:
7 (Pdb) l
8 1 from pdb import set_trace
9 2 def binary_search_buggy(v,x):
10 3     set_trace()
11 4     start,end = 0,len(v)-1
12 5 -> while start<end:
13 6         mid = (start+end)//2
14 7         if v[mid]>=x:
15 8             end = mid
16 9         else:

```

```

17 10          start = mid+1
18 11      return start
19 (Pdb) b 7
20 Breakpoint 1 at chapter2/binary_search_buggy_debug.py:7
21 (Pdb) c
22 > chapter2/binary_search_buggy_debug.py(7)binary_search_buggy()
23 -> if v[mid]>=x:
24 (Pdb) c
25 > chapter2/binary_search_buggy_debug.py(7)binary_search_buggy()
26 -> if v[mid]>=x:
27 (Pdb) mid
28 2
29 (Pdb) n
30 > chapter2/binary_search_buggy_debug.py(10)binary_search_buggy()
31 -> start = mid+1
32 (Pdb) n
33 > chapter2/binary_search_buggy_debug.py(5)binary_search_buggy()
34 -> while start<end:
35 (Pdb) start
36 3
37 (Pdb) n
38 > chapter2/binary_search_buggy_debug.py(11)binary_search_buggy()
39 -> return start

```

Similar to R, command `n` would evaluate the next statement in `pdb`. Typing command `l` would show the current line of current execution. Command `b line_number` would set the corresponding line as a breakpoint; and `c` would continue the execution until the next breakpoint (if exists).

In R, besides the `browser()` function there are a pair of functions `debug()` and `undebug()` which are also very handy when we try to debug a function; especially when the function is wrapped in a package. More specifically, the `debug` function would invoke the debugging environment whenever we call the function to debug. See the example below how we invoke the debugging environment for the `sd` function (standard deviation calculation).

R

```

1 > x=c(1,1,2)
2 > debug(sd)
3 > sd(x)
4 debugging in: sd(x)
5 debug: sqrt(var(if (is.vector(x) || is.factor(x)) x else as.double(x),
6     na.rm = na.rm))
7 Browse[2]> ls()
8 [1] "na.rm" "x"
9 Browse[2]> Q

```

```

10 > undebug(sd)
11 > sd(x)
12 [1] 0.5773503

```

The binary_search solutions are fixed below.

R

chapter2/find_binary_search.py

```

1 binary_search=function(v,x){
2   if (x>v[length(v)]) {return(NULL)}
3   start = 1
4   end = length(v)
5   while (start<end){
6     mid = (start+end)
7     if (v[mid]>=x){
8       end = mid
9     }else{
10      start = mid+1
11    }
12  }
13  return(start)
14 }

```

Python

chapter2/find_binary_search.py

```

1 def binary_search(v,x):
2   if x>v[-1]: return
3   start,end = 0,len(v)-1
4   while start<end:
5     mid = (start+end)//2
6     if v[mid]>=x:
7       end = mid
8     else:
9       start = mid+1
10  return start

```

2.3 Benchmarking

By benchmarking, I mean measuring the entire operation time of a piece of program. There is another term called profiling which is related to benchmarking. But profiling is more complex since it commonly aims at

understanding the behavior of the program and optimizing the program in terms of time elapsed during the operation.

In R, I like using the `microbenchmark` package. And in Python, `timeit` module is a good tool to use when we want to benchmark a small bits of Python code.

As mentioned before, the run time complexity of binary search is better than that of a for-loop search. We can do benchmarking to compare the two algorithms.

R

chapter2/benchmark.R

```

1 library(microbenchmark)
2 source('binary_search.R')
3 source('find_pos.R')
4
5 v=1:10000
6
7 # call each function 1000 times;
8 # each time we randomly select an integer as the target value
9
10 # for-loop solution
11 set.seed(2019)
12 print(microbenchmark(find_pos(v, sample(10000,1)), times=1000))
13 # binary-search solution
14 set.seed(2019)
15 print(microbenchmark(binary_search(v, sample(10000,1)), times=1000))

```

In the R code above, `times=1000` means we want to call the function 1000 times in the benchmarking process. The `sample()` function is used to draw samples from a set of elements. Specifically, we pass the argument 1 to `sample()` to draw a single element. It's the first time we use `set.seed()` function in this book. In R/Python, we draw random numbers based on the pseudorandom number generator (PRNG) algorithm²⁶. The sequence of numbers generated by PRNG is completely determined by an initial value, i.e., the seed. Whenever a program involves the usage of PRNG, it is better to set the seed in order to get replicable results (see the example below).

R

```

1 > set.seed(2019)
2 > rnorm(1)
3 [1] 0.7385227
4 > rnorm(1)
5 [1] -0.5147605

```

R

```

1 > set.seed(2019)
2 > rnorm(1)
3 [1] 0.7385227
4 > set.seed(2019)
5 > rnorm(1)
6 [1] 0.7385227

```

²⁶ https://en.wikipedia.org/wiki/Pseudorandom_number_generator

Now let's run the R script to see the benchmarking result. llr

```

1 > source('benchmark.R')
2 Unit: microseconds
3           expr    min       lq      mean   median       uq
4 find_pos(v, sample(10000, 1)) 3.96 109.5385 207.6627 207.5565 307.8875
5   536.171
6 neval
7 1000
8 Unit: microseconds
9           expr    min       lq      mean   median       uq      max
10 binary_search(v, sample(10000, 1)) 5.898 6.3325 14.2159 6.6115 7.3635 6435.57
11 neval
12 1000

```

The binary_search solution is much more efficient based on the benchmarking result. Doing the same benchmarking in Python is a bit of complicated.

Python

chapter2/benchmark.py

```

1 from binary_search import binary_search
2 from find_pos import find_pos
3 import timeit
4 import random
5
6 v=list(range(1,10001))
7
8 def test_for_loop(n):
9     random.seed(2019)
10    for _ in range(n):
11        find_pos(v,random.randint(1,10000))
12
13 def test_bs(n):
14     random.seed(2019)
15     for _ in range(n):
16         binary_search(v,random.randint(1,10000))
17
18 # for-loop solution
19 print(timeit.timeit('test_for_loop(1000)',setup='from __main__ import
20     test_for_loop',number=1))
21 # binary_search solution
22 print(timeit.timeit('test_bs(1000)',setup='from __main__ import test_bs',
23     number=1))

```

The most interesting part of the Python code above is `from __main__ import`. Let's ignore it for now, and we would revisit it later.

Below is the benchmarking result in Python (the unit is second).

Python

```
1 chapter2 $python3 benchmark.py
2 0.284618441
3 0.003966589000000002
```

2.4 Vectorization

In parallel computing, automatic vectorization ²⁷ means a program in a scalar implementation is converted to a vector implementation which process multiple pairs of operands simultaneously by compilers that feature auto-vectorization. For example, let's calculate the element-wise sum of two arrays `x` and `y` of the same length in C programming language.

```
1 int x[4] = {1,2,3,4};
2 int y[4] = {0,1,2,3};
3 int z[4];
4 for (int i=0;i<4;i++){
5     z[i]=x[i]+y[i];
6 }
```

The C code above might be vectorized by the compiler so that the actual number of iterations performed could be less than 4. If 4 pairs of operands are processed at once, there would be only 1 iteration. Automatic vectorization may make the program runs much faster in some languages like C. However, when we talk about vectorization in R/Python, it is different from automatic vectorization. Vectorization in R/Python usually refers to the human effort paid to avoid for-loops. First, let's see some examples of how for-loops may slow your programs in R/Python.

R

chapter2/vectorization_1.R

```
1 library(microbenchmark)
2
3 # generate n standard normal r.v
4 rnorm_loop = function(n){
5   x=rep(0,n)
6   for (i in 1:n) {x[i]=rnorm(1)}
7 }
8
9 rnorm_vec = function(n){
```

²⁷ https://en.wikipedia.org/wiki/Automatic_vectorization


```

10 x=rnorm(n)
11 }
12
13 n=100
14 # for loop
15 print(microbenchmark(rnorm_loop(n),times=1000))
16 # vectorize
17 print(microbenchmark(rnorm_vec(n),times=1000))

```

Running the R code results in the following result on my local machine.

```

1 > source('vectorization_1.R')
2 Unit: microseconds
3      expr      min       lq      mean   median      uq      max  neval
4  rnorm_loop(n) 131.622 142.699 248.7603 145.3995 270.212 16355.6  1000
5 Unit: microseconds
6      expr      min       lq      mean  median      uq      max  neval
7  rnorm_vec(n)  6.696   7.128  10.87463   7.515   8.291 2422.338  1000

```

Python

```

1 import timeit
2 import numpy as np
3
4 def rnorm_for_loop(n):
5     x=[0]*n # create a list with n 0s
6     np.random.seed(2019)
7     for _ in range(n):
8         np.random.normal(0,1,1)
9
10 def rnorm_vec(n):
11     np.random.seed(2019)
12     x = np.random.normal(0,1,n)
13
14 print("for loop")
15 print(f'{timeit.timeit("rnorm_for_loop(100)",setup="from __main__ import rnorm_for_loop",number=1000):.6f} seconds')
16 print("vectorized")
17 print(f'{timeit.timeit("rnorm_vec(100)",setup="from __main__ import rnorm_vec",number=1000):.6f} seconds')

```

Please note that in this Python example we are using the random submodule of numpy module instead of the built-in random module since random module doesn't provide the vectorized version of random number generation function. Running the Python code results in the following result on my local machine.

Python

```
1 chapter2 $python3.7 vectorization_1.py
2 for loop
3 0.258466 seconds
4 vectorized
5 0.008213 seconds
```

In either R or Python, the vectorized version of random normal random variable (r.v.) is significantly faster than the scalar version. It is worth noting the usage of the `print(f"")` statement in the Python code, which is different from the way how we print the object of `Complex` class in [chapter 1](#). In the code above, we use the `f-string`²⁸ which is a literal string prefixed with `'f'` containing expressions inside `{}` which would be replaced with their values. `f-string` was a feature introduced since Python 3.6. If you are familiar with Scala, you may find that this feature is quite similar with the string interpolation mechanism introduced since Scala 2.10.

It's also worth noting that lots of built-in functions in R are already vectorized, such as the basic arithmetic operators, comparison operators, `ifelse()`, element-wise logical operators `&`, `|`. But the logical operators `&&`, `||` are not vectorized.

In addition to vectorization, there are also some built-in functions which may help to avoid the usages of for-loops. For example, in R we might be able use the `apply` family of functions to replace for-loops; and in Python the `map()` function can also be useful. In the Python pandas module, there are also many usages of `map/apply` methods. But in general the usage of `apply/map` functions has little or nothing to do with performance improvement. However, appropriate usages of such functions may help with the readability of the program. Compared with the `apply` family of functions in R, I think the `do.call()` function is more useful in practice. We would spend some time in `do.call()` later.

²⁸ <https://www.python.org/dev/peps/pep-0498/>

Application: Biham–Middleton–Levine (BML) traffic model

Considering the importance of vectorization in scientific programming, let's try to get more familiar with vectorization through the Biham–Middleton–Levine (BML) traffic model²⁹. The BML model is very important in modern studies of traffic flow since it exhibits a sharp phase transition from free flowing status to a fully jammed status. A simplified BML model could be characterized as follows:

- Initialized on a 2-D lattice, each site of which is either empty or occupied by a colored particle (blue or red);
- Particles are distributed randomly through the initialization according to a uniform distribution; the two colors of particles are equally distributed.
- On even time steps, all blue particles attempt to move one site up and an attempt fails if the site to occupy is not empty;
- On Odd time steps, all red particles attempt to move one site right and an attempt fails if the site to occupy is not empty;
- The lattice is assumed periodic which means when a particle moves out of the lattice, it would move into the lattice from the opposite side.

The BML model specified above is implemented in both R/Python as follows to illustrate the usage of vectorization.

R

chapter2/BML.R

```

1 library(R6)
2 BML = R6Class(
3   "BML",
4   public = list(
5     # alpha is the parameter of the uniform distribution to control particle
      # distribution's density
6     # m*n is the dimension of the lattice
7     alpha = NULL,
8     m = NULL,
9     n = NULL,
10    lattice = NULL,
11    initialize = function(alpha, m, n) {
12      self$alpha = alpha
13      self$m = m
14      self$n = n
15      self$initialize_lattice()
16    },
17    initialize_lattice = function() {
18      # 0 -> empty site

```

```

19     # 1 -> blue particle
20     # 2 -> red particle
21     u = runif(self$m * self$n)
22     # the usage of L is to make sure the elements in particles are of type
        integer;
23     # otherwise they would be created as double
24     particles = rep(0L, self$m * self$n)
25     # doing inverse transform sampling
26     particles[(u > self$alpha) &
27               (u <= (self$alpha + 1.0) / 2)] = 1L
28     particles[u > (self$alpha + 1.0) / 2] = 2L
29     self$lattice = array(particles, c(self$m, self$n))
30 },
31 odd_step = function() {
32     blue.index = which(self$lattice == 1L, arr.ind = TRUE)
33     # make a copy of the index
34     blue.up.index = blue.index
35     # blue particles move 1 site up
36     blue.up.index[, 1] = blue.index[, 1] - 1L
37     # periodic boundary condition
38     blue.up.index[blue.up.index[, 1] == 0L, 1] = self$m
39     # find which moves are feasible
40     blue.movable = self$lattice[blue.up.index] == 0L
41     # move blue particles one site up
42     # drop=FALSE prevents the 2D array degenerates to 1D array
43     self$lattice[blue.up.index[blue.movable, , drop = FALSE]] = 1L
44     self$lattice[blue.index[blue.movable, , drop = FALSE]] = 0L
45 },
46 even_step = function() {
47     red.index = which(self$lattice == 2L, arr.ind = TRUE)
48     # make a copy of the index
49     red.right.index = red.index
50     # red particles move 1 site right
51     red.right.index[, 2] = red.index[, 2] + 1L
52     # periodic boundary condition
53     red.right.index[red.right.index[, 2] == (self$n + 1L), 2] = 1
54     # find which moves are feasible
55     red.movable = self$lattice[red.right.index] == 0L
56     # move red particles one site right
57     self$lattice[red.right.index[red.movable, , drop = FALSE]] = 2L
58     self$lattice[red.index[red.movable, , drop = FALSE]] = 0L
59 }
60 )
61 )

```

Now we can create a simple BML system on a 5×5 lattice using the R code above.

R

```

1 > source('BML.R')
2 > set.seed(2019)
3 > bml=BML$new(0.4,5,5)
4 > bml$lattice
5      [,1] [,2] [,3] [,4] [,5]
6 [1,]    2    0    2    1    1
7 [2,]    2    2    1    0    1
8 [3,]    0    0    0    2    2
9 [4,]    1    0    0    0    0
10 [5,]    0    1    1    1    0
11 > bml$odd_step()
12 > bml$lattice
13      [,1] [,2] [,3] [,4] [,5]
14 [1,]    2    0    2    1    0
15 [2,]    2    2    1    0    1
16 [3,]    1    0    0    2    2
17 [4,]    0    1    1    1    0
18 [5,]    0    0    0    0    1
19 > bml$even_step()
20 > bml$lattice
21      [,1] [,2] [,3] [,4] [,5]
22 [1,]    0    2    2    1    0
23 [2,]    2    2    1    0    1
24 [3,]    1    0    0    2    2
25 [4,]    0    1    1    1    0
26 [5,]    0    0    0    0    1

```

In the initialization step, we used the inverse transform sampling approach³⁰ to generate the status of each site. Inverse transform sampling method is basic but powerful approach to generate r.v. from any probability distribution given its cumulative distribution function (CDF). Reading the wiki page is enough to master this sampling method.

Python

```

1 import numpy as np
2
3 class BML:

```

³⁰ https://en.wikipedia.org/wiki/Inverse_transform_sampling

```

4     def __init__(self, alpha, m, n):
5         self.alpha = alpha
6         self.shape = (m, n)
7         self.initialize_lattice()
8
9     def initialize_lattice(self):
10        u = np.random.uniform(0.0, 1.0, self.shape)
11        # instead of using default list, we use np.array to create the lattice
12        self.lattice = np.zeros_like(u, dtype=int)
13        # the parentheses below can't be ignored
14        self.lattice[(u > self.alpha) & (u <= (1.0+self.alpha)/2.0)] = 1
15        self.lattice[u > (self.alpha+1.0)/2.0] = 2
16
17    def odd_step(self):
18        # please note that np.where returns a tuple which is immutable
19        blue_index = np.where(self.lattice == 1)
20        blue_index_i = blue_index[0] - 1
21        blue_index_i[blue_index_i < 0] = self.shape[0]-1
22        blue_movable = self.lattice[(blue_index_i, blue_index[1])] == 0
23        self.lattice[(blue_index_i[blue_movable],
24                      blue_index[1][blue_movable])] = 1
25        self.lattice[(blue_index[0][blue_movable],
26                      blue_index[1][blue_movable])] = 0
27
28    def even_step(self):
29        red_index = np.where(self.lattice == 2)
30        red_index_j = red_index[1] + 1
31        red_index_j[red_index_j == self.shape[1]] = 0
32        red_movable = self.lattice[(red_index[0], red_index_j)] == 0
33        self.lattice[(red_index[0][red_movable],
34                      red_index_j[red_movable])] = 2
35        self.lattice[(red_index[0][red_movable],
36                      red_index[1][red_movable])] = 0

```

The Python implementation is also given.

R

```

1 >>> import numpy as np
2 >>> np.random.seed(2019)
3 >>> from BML import BML
4 >>> bml=BML(0.4,5,5)
5 >>> bml.lattice
6 array([[2, 0, 1, 1, 2],

```

```

7         [0, 2, 2, 2, 1],
8         [1, 0, 0, 2, 0],
9         [2, 0, 1, 0, 2],
10        [1, 1, 0, 2, 1]])
11 >>> bml.odd_step()
12 >>> bml.lattice
13 array([[2, 0, 0, 1, 2],
14        [1, 2, 2, 2, 1],
15        [0, 0, 1, 2, 0],
16        [2, 1, 0, 0, 2],
17        [1, 0, 1, 2, 1]])
18 >>> bml.even_step()
19 >>> bml.lattice
20 array([[0, 2, 0, 1, 2],
21        [1, 2, 2, 2, 1],
22        [0, 0, 1, 0, 2],
23        [2, 1, 0, 0, 2],
24        [1, 0, 1, 2, 1]])

```

Please note that although we have imported `numpy` in `BML.py`, we import it again in the code above in order to set the random seed. If we change the line to `from BML import *`, we don't need to import `numpy` again. But it is not recommended to `import *` from a module.

2.5 Embarrassingly parallelism in R/Python

According to the explanation of wikipedia ³¹, single-threading is the processing of one command at a time, and its opposite is multithreading. A process is the instance of a computer program executed by one or many threads ³². Multithreading is not the same as parallelism. In a single processor environment, the processor can switch execution between threads, which results in concurrent execution. However, it is possible a process with multithreads runs on on a multiprocessor environment and each of the threads on a separate processor, which results in parallel execution.

Both R and Python are single-threaded. In Python, there is a `threading` package ³³, which support multithreading on a single core. It may suit some specific tasks. For example, in web scraping multithreading on a single core may be able to speed up the program if the download time is in excess of the CPU processing time.

Now let's talk about embarrassingly parallelism by multi-processing. Embarrassingly parallel problem is one where little or no effort is needed to separate the problem into a number of parallel tasks ³⁴. In R there are various packages supporting multi-processing on multiple CPU cores, for example, the `parallel` package, which is my favorite one. In Python, there are also some available modules, such as `multiprocessing`, `joblib` and `concurrent.futures`. Let's see an application of embarrassingly parallelism to calculate π

³¹ [https://en.wikipedia.org/wiki/Thread_\(computing\)](https://en.wikipedia.org/wiki/Thread_(computing))

³² [https://en.wikipedia.org/wiki/Process_\(computing\)](https://en.wikipedia.org/wiki/Process_(computing))

³³ <https://docs.python.org/3.7/library/threading.html>

³⁴ https://en.wikipedia.org/wiki/Embarrassingly_parallel

using Monte Carlo simulation ³⁵.

Application: Monte Carlo simulation to estimate π via parallelization

Monte Carlo simulation provides a simple and straightforward way to estimate π . We know the area of a circle with radius 1 is just π . Thus, we can convert the original problem of π calculation to a new problem, i.e., how to calculate the area of a circle with radius 1. We also know the area of a square with side length 2 is equal to 4. Thus, π can be calculated as $4r_{c/s}$ where $r_{c/s}$ denotes the ratio of areas of a circle with radius 1 and a square with side length 2. Now the problem is how to calculate the ratio $r_{c/s}$? When we randomly throw n points into the square and m of these points fall into the inscribed circle, then we can estimate the ratio as m/n . As a result, a natural estimate of π is $4m/n$. This problem is an embarrassingly parallel problem by its nature. Let's see how we implement the idea in R/Python.

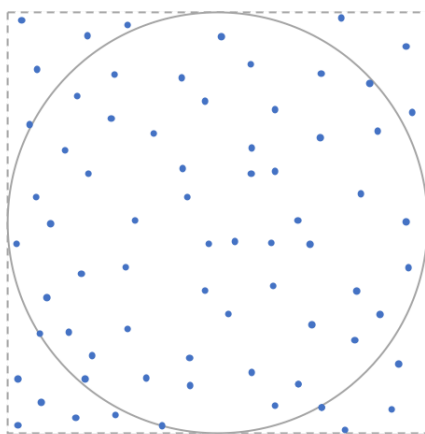


Figure 2.1: Generate points within a square and count how many times these points fall into the inscribed circle

R

chapter2/pi.R

```
1 library(parallel)
2 count_inside_point = function(n){
3   m = 0
4   for (i in 1:n){
5     p_x = runif(1, -1, 1)
6     p_y = runif(1, -1, 1)
7     if (p_x^2 + p_y^2 <=1){
8       m = m+1
9     }
10  }
11  m
}
```

³⁵ https://en.wikipedia.org/wiki/Monte_Carlo_method


```

12 }
13
14 # now let's use the mclapply for parallelization
15 generate_points_parallel = function(n){
16   # detectCores() returns the number of cores available
17   # we assign the task to each core
18   unlist(mclapply(X = rep(n %% detectCores(), detectCores()), FUN=count_inside_
19     _point))
20 }
21
22 # now let's use vectorization
23 generate_points_vectorized = function(n){
24   p = array(runif(n*2, -1, 1), c(n, 2))
25   sum((p[, 1]^2 + p[, 2]^2) <= 1)
26 }
27
28 pi_naive = function(n) cat('naive: pi -', 4*count_inside_point(n)/n, '\n')
29
30 pi_parallel = function(n) cat('parallel: pi -', 4*sum(generate_points_parallel
31   (n))/n, '\n')
32
33 pi_vectorized = function(n) cat('vectorized: pi -', 4*sum(generate_points_
34   vectorized(n))/n, '\n')

```

In the above R code snippet, we use the function `mclapply` which is not currently available on some operation systems³⁶. When it is not available, we may consider to use `parLapply` instead.

Python

chapter2/pi.py

```

1 # now let's try the parallel approach
2 # each process uses the same seed, which is not desired
3 def generate_points_parallel(n):
4     pool = mp.Pool()
5     # we ask each process to generate n//mp.cpu_count() points
6     return pool.map(count_inside_point, [n//mp.cpu_count()]*mp.cpu_count())
7
8 # set seed for each process
9 # first, let's define a helper function
10 def helper(args):
11     n, s = args
12     seed(s)
13     return count_inside_point(n)

```

³⁶ <https://stat.ethz.ch/R-manual/R-devel/library/parallel/doc/parallel.pdf>

```

14
15 def generate_points_parallel_set_seed(n):
16     pool = mp.Pool() # we can also specify the number of processes by Pool(
        number)
17     # we ask each process to generate n//mp.cpu_count() points
18     return pool.map(helper, list(zip([n//mp.cpu_count()]*mp.cpu_count(), range(
        mp.cpu_count()))))
19
20 # another optimization via vectorization
21 def generate_points_vectorized(n):
22     p = uniform(-1, 1, size=(n,2))
23     return np.sum(np.linalg.norm(p, axis=1) <= 1)
24
25 def pi_naive(n):
26     print(f'pi: {count_inside_point(n)/n*4:.6f}')
27
28 def pi_parallel(n):
29     print(f'pi: {sum(generate_points_parallel_set_seed(n))/n*4:.6f}')
30
31 def pi_vectorized(n):
32     print(f'pi: {generate_points_vectorized(n)/n*4:.6f}')

```

In the above Python code snippet, we defined a function `generate_points_parallel` which returns a list of number of inside points. But the problem of this function is that each process in the pool shares the same random state. As a result, we will obtain the a list of duplicated numbers by calling this function. To overcome the issue, we defined another function `generate_points_parallel_set_seed`. Let's see the results of our estimation for π running on a laptop.

R

```

1 > source('pi.R')
2 > system.time(pi_naive(1e6))
3 naive: pi - 3.144592
4   user  system elapsed
5   8.073   1.180   9.333
6 > system.time(pi_parallel(1e6))
7 parallel: pi - 3.1415
8   user  system elapsed
9   4.107   0.560   4.833
10 > system.time(pi_vectorized(1e6))
11 vectorized: pi - 3.141224
12   user  system elapsed
13   0.180   0.031   0.214

```

Python

```

1 >>> import timeit
2 >>> from pi import pi_naive, pi_parallel, pi_vectorized
3 >>> print(f'naive - {timeit.timeit("pi_naive(1000000)",setup="from __main__
      import pi_naive",number=1):.6f} seconds')
4 pi: 3.141056
5 naive - 4.363822 seconds
6 >>> print(f'parallel - {timeit.timeit("pi_parallel(1000000)",setup="from
      __main__ import pi_parallel",number=1):.6f} seconds')
7 pi: 3.141032
8 parallel - 2.204697 seconds
9 >>> print(f'vectorized - {timeit.timeit("pi_vectorized(1000000)",setup="from
      __main__ import pi_vectorized",number=1):.6f} seconds')
10 pi: 3.139936
11 vectorized - 0.148950 seconds

```

We see the winner in this example is vectorization, and the parallel solution is better than the naive solution. However, when the problem cannot be vectorized we may use parallelization to achieve better performance.

2.6 Scope of Variables

We have seen how to define variables in R/Python in [chapter 1](#), and we have known that a variable is an identifier to a location in memory. What is the scope of a variable and why does it matter? Let's first have a look at the code snippets below.

R

```

1 > x=1
2 > var_func_1 = function(){print(x)}
3 > var_func_1()
4 [1] 1
5 > var_func_2 = function(){x=x+1;
      print(x)}
6 > var_func_2()
7 [1] 2
8 > x
9 [1] 1

```

Python

```

1 >>> x=1
2 >>> def var_func_1():print(x)
3 >>> var_func_1()
4 1
5 >>> def var_func_2():x+=1
6 ...
7 >>> var_func_2()
8 Traceback (most recent call last):
9   File "<stdin>", line 1, in <
      module>
10   File "<stdin>", line 1, in
      var_func_2
11 UnboundLocalError: local variable
    'x' referenced before
    assignment

```

The results of the code above seem strange before knowing the concept of variable scope. Inside a function, a variable may refer to a function argument/parameter or it could be formally declared inside the function which is called a local variable. But in the code above, `x` is neither a function argument nor a local variable. How does the `print()` function know where the identifier `x` points to?

The scope of a variable determines where the variable is available/accessible (can be referenced). Both R and Python apply lexical/static scoping for variables, which set the scope of a variable based on the structure of the program. In static scoping, when an 'unknown' variable referenced, the function will try to find it from the most closely enclosing block. That explains how the `print()` function could find the variable `x`.

In the R code above, `x=x+1` the first `x` is a local variable created by the `=` operator; the second `x` is referenced inside the function so the static scoping rule applies. As a result, a local variable `x` which is equal to 2 is created, which is independent with the `x` outside of the function `var_func_2()`. However, in Python when a variable is assigned a value in a statement the variable would be treated as a local variable and that explains the `UnboundLocalError`.

Is it possible to change a variable inside a function which is declared outside the function? Based on the static scoping rule only, it's impossible. But there are workarounds in both R/Python. In R, we need the help of environment; and in Python we can use the keyword `global`.

So what is an environment in R? An environment is a place where objects are stored. When we invoke the interactive R session, an environment named as `.GlobalEnv` is created automatically. We can also use the function `environment()` to get the present environment. The `ls()` function can take an environment as the argument to list all objects inside the environment.

R

```

1 $r
2 > typeof(.GlobalEnv)
3 [1] "environment"
4 > environment()
5 <environment: R_GlobalEnv>
6 > x=1
7 > ls(environment())
8 [1] "x"
9 > env_func_1=function(x){
10 +   y=x+1
11 +   print(environment())
12 +   ls(environment())
13 + }
14 > env_func_1(2)
15 <environment: 0x7fc59d165a20>
16 [1] "x" "y"
17 > env_func_2=function(){print(environment())}
18 > env_func_2()
19 <environment: 0x7fc59d16f520>

```

The above code shows that each function has its own environment containing all function arguments and

local variables declared inside the function. In order to change a variable declared outside of a function, we need the access of the environment enclosing the variable to change. There is a function `parent_env(e)` that returns the parent environment of the given environment `e` in R. Using this function, we are able to change the value of `x` declared in `.GlobalEnv` inside a function which is also declared in `.GlobalEnv`. The `global` keyword in Python works in a totally different way, which is simple but less flexible.

R

```

1 > x=1
2 > env_func_3=function(){
3 +   cur_env=environment()
4 +   par_env=parent.env(cur_env)
5 +   par_env$x=2
6 + }
7 > env_func_3()
8 > x
9 [1] 2

```

Python

```

1 >>> def env_func_3():
2 ...     global x
3 ...     x = 2
4 ...
5 >>> x=1
6 >>> env_func_3()
7 >>> x
8 2

```

I seldomly use the `global` keyword in Python, if any. But the environment in R could be very handy in some occasions. In R, environment could be used as a purely mutable version of the `list` data structure.

R

```

1 # list is not purely mutable
2 > x=list(1)
3 > tracemem(x)
4 [1] "<0x7f829183f6f8>"
5 > x$a=2
6 > tracemem(x)
7 [1] "<0x7f828f4d05c8>"

```

R

```

1 # environment is purely mutable
2 > x=new.env()
3 > x
4 <environment: 0x7f8290aee7e8>
5 > x$a=2
6 > x
7 <environment: 0x7f8290aee7e8>

```

Actually, the object of an R6 class type is also an environment.

R

```

1 > # load the Complex class that we defined in chapter 1
2 > x = Complex$new(1,2)
3 > typeof(x)
4 [1] "environment"

```

In Python, we can assign values to multiple variables in one line.

Python

```

1 >>> x,y = 1,2
2 >>> x
3 1
4 >>> y
5 2

```

Python

```

1 >>> x,y=(1,2)
2 >>> print(x,y)
3 1 2
4 >>> (x,y)=(1,2)
5 >>> print(x,y)
6 1 2
7 >>> [x,y]=(1,2)
8 >>> print(x,y)
9 1 2

```

Even though in the left snippet above there aren't parentheses embracing 1,2 after the = operator, a tuple is created first and then the tuple is unpacked and assigned to x, y. Such mechanism doesn't exist in R, but we can define our own multiple assignment operator with the help of environment.

R

chapter2/multi_assignment.R

```

1 `%= `% = function(left, right) {
2   # we require the RHS to be a list strictly
3   stopifnot(is.list(right))
4   # dest_env is the destination environment enclosing the variables on LHS
5   dest_env = parent.env(environment())
6   left = substitute(left)
7
8   recursive_assign = function(left, right, dest_env) {
9     if (length(left) == 1) {
10       assign(x = deparse(left),
11             value = right,
12             envir = dest_env)
13       return()
14     }
15     if (length(left) != length(right) + 1) {
16       stop("LHS and RHS must have the same shapes")
17     }
18
19     for (i in 2:length(left)) {
20       recursive_assign(left[[i]], right[[i - 1]], dest_env)
21     }
22   }
23
24   recursive_assign(left, right, dest_env)

```

25 }

Before going into the script deeper, first let's see the usage of the multiple assignment operator we defined.

R

```

1 > source('multi_assignment.R')
2 > c(x,y,z) %=% list(1,"Hello World!",c(2,3))
3 > x
4 [1] 1
5 > y
6 [1] "Hello World!"
7 > z
8 [1] 2 3
9 > list(a,b) %=% list(1,as.Date('2019-01-01'))
10 > a
11 [1] 1
12 > b
13 [1] "2019-01-01"

```

In the `%=%` operator defined above, we used two functions `substitute`, `deparse` which are very powerful but less known by R novices. To better understand these functions as well as some other less known R functions, the Rchaeology ³⁷ tutorial is worth reading.

It is also interesting to see that we defined the function `recursive_assign` inside the `%=%` function. Both R and Python support the concept of first class functions. More specifically, a function in R/Python is an object, which can be

1. stored as a variable;
2. passed as a function argument;
3. returned from a function.

The essential idea behind the `recursive_assign` function is a depth-first search (DFS), which is a fundamental graph traversing algorithm ³⁸. In the context of the `recursive_assign` function, we use DFS to traverse the parse tree of the left argument created by calling `substitute(left)`.

2.7 Miscellaneous

We have introduced the basics of R/Python programming so far. There are much more to learn to become an advanced user of R/Python. For example, the appropriate usages of `iterator`, `generator`, `decorator` could improve both the conciseness and readability of your Python code. The `generator` ³⁹ is commonly seen in machine learning programs to prepare training/testing samples. `decorator` is a kind of syntactic sugar to allow the modification of a function's behavior in a simple way. In R there are no built-in `iterator`, `generator`, `decorator`, but you may find some third-party libraries to mimic these features; or you may

³⁷ <https://cran.r-project.org/web/packages/rockchalk/vignettes/Rchaeology.pdf>

³⁸ https://en.wikipedia.org/wiki/Depth-first_search

³⁹ <https://docs.python.org/3/howto/functional.html>

try to implement your own.

One advantage of Python over R is that there are some built-in modules containing high-performance data structures or commonly-used algorithms implemented efficiently. For example, I enjoy using the deque structure in the Python collections module ⁴⁰, but there is no built-in counterpart in R. We have written our own binary search algorithm earlier in this Chapter, which can also be replaced by the functions in the built-in module bisect ⁴¹ in Python.

Another important aspect of programming is testing. Unit testing is a typical software testing method that is commonly adopted in practice. In R there are two third-party packages testthat and RUnit. In Python, the built-in unittest is quite powerful. Also, the third-party module pytest ⁴² is very popular.

⁴⁰ <https://docs.python.org/3/library/collections.html>

⁴¹ <https://docs.python.org/3.7/library/bisect.html>

⁴² <https://docs.pytest.org/en/latest/>

4

Linear Regression

After finishing the first two Chapters, I was thinking of the topic of Chapter 3. Finally I chose linear regression to write, because of the importance of these models in machine learning and data science.

There are numerous books/online courses available on the theory of linear regression, among which my favorite one is - The Elements of Statistical Learning [1]. So what is the purpose to write a chapter about these models in another book of data science? Many audience would be interested in how to implement their own regression models rather than using the off-the-shelf software packages. By the end of this chapter, we will build up our own linear regression models in R/Python. We would also see how we would reuse some functions in different regressions, such as linear regression, and these regressions with L2 penalties.

4.1 Basics of linear regression

We start from the matrix form of linear regression.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (4.1)$$

where the column vector \mathbf{y} contains n observations on the dependent variable, \mathbf{X} is a $n \times (p + 1)$ matrix ($n > p$) of independent variables with constant vector $\mathbf{1}$ in the first column, $\boldsymbol{\beta}$ is a column vector of unknown population parameters to estimate based on the data, and $\boldsymbol{\epsilon}$ is the error term (or noise). For the sake of illustration, (4.1) can be extended as in

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{21} & \cdots & X_{p1} \\ 1 & X_{12} & X_{22} & \cdots & X_{p2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{1n} & X_{2n} & \cdots & X_{pn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \quad (4.2)$$

We apply ordinary least squares (OLS)⁴⁹ approach to estimate the model parameter $\boldsymbol{\beta}$ since it requires fewer assumptions than other estimation methods such as maximum likelihood estimation⁵⁰. Suppose the estimated model parameter is denoted as $\hat{\boldsymbol{\beta}}$, we define the residual vector of the system as $\mathbf{e} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$. The

⁴⁹ https://en.wikipedia.org/wiki/Ordinary_least_squares

⁵⁰ https://en.wikipedia.org/wiki/Maximum_likelihood_estimation

idea of OLS is to find $\hat{\beta}$ which can minimize the sum of squared residuals (SSR), i.e.,

$$\min_{\hat{\beta}} e'e \quad (4.3)$$

Now the question is how to solve the optimization problem (4.2). First, let's expand the SSR.

$$\begin{aligned} e'e &= (y - X\hat{\beta})'(y - X\hat{\beta}) \\ &= y'y - 2\hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta} \end{aligned} \quad (4.4)$$

The first and second order derivatives are calculated as follows

$$\begin{cases} \frac{\partial e'e}{\partial \hat{\beta}} = -2X'y + 2X'X\hat{\beta} \\ \frac{\partial^2 e'e}{\partial \hat{\beta} \partial \hat{\beta}'} = 2X'X. \end{cases} \quad (4.5)$$

We see that the second order derivative is positive semidefinite which implies the SSR in OLS is a convex function (see section 3.1.4 in [2]) and for an unconstrained convex optimization problem, the necessary as well as sufficient condition for optimality is that the first order derivative equals 0 (see section 4.2.3 in [2]). Optimization of convex function is very important in machine learning. Actually, the parameter estimations of many machine learning models are convex optimization problems.

Based on the analysis above, the solution of (4.3) is given in (4.6).

$$\hat{\beta} = (X'X)^{-1}X'y \quad (4.6)$$

Now it seems we are ready to write our own linear regression model in R/Python. The solution in (4.6) involves matrix transportation, multiplication and inversion, all of which are supported in both R and Python. In Python, we can use the `numpy` module for the matrix operations.

However, in practice we don't solve linear regression with (4.6) directly. Why? Let's see an example with

$$x = \begin{bmatrix} 1e+6 & -1 \\ -1 & 1e-6 \end{bmatrix}.$$

R

```
1 > x=array(c(10^6,-1,-1,10^-6),c
2       (2,2))
3 > solve(t(x) %*% x) # solve()
   would return the inverse matrix
Error in solve.default(t(x) %*% x)
:
4 system is computationally
   singular: reciprocal
   condition number = 2.22044e
   -28
```

Python

```
1 >>> import numpy as np
2 >>> x=np.array([[1e6,-1],[-1,1e
3       -6]])
4 >>> np.linalg.inv(np.dot(x.
5       transpose(),x))
array([[4.50359963e+03, 4.50359963
6       e+09],
7       [4.50359963e+09, 4.50359963
8       e+15]])
```

The R code above throws an error because of the singularity of $bmX'X$. It's interesting that the corresponding Python code doesn't behave in the same way as R, which has been reported as an issue on github

⁵¹.

When the matrix $bmX'X$ is singular, how to solve the OLS problem? In this book, we would focus on the QR decomposition based solution. Singular value decomposition (SVD) can also be used to solve OLS, which would not be covered in this book.

In linear algebra, a QR decomposition ⁵² of matrix X would factorize X into a product, i.e., $X = QR$ where Q are orthogonal matrices and R is an upper triangular matrix. Since the matrix Q is orthogonal ($Q' = Q^{-1}$), we have

$$\begin{aligned}\hat{\beta} &= (X'X)^{-1}X'y \\ &= (R'Q'QR)^{-1}R'Q'y \\ &= (R'R)^{-1}R'Q'y \\ &= R^{-1}Q'y\end{aligned}\tag{4.7}$$

Now we are ready to write our simple R/Python functions for linear regression with the help of QR decomposition according to (4.7).

R

chapter4/qr_solver.R

```
1 qr_solver=function(x,y){
2   qr.coef(qr(x),y)
3 }
```

Python

chapter4/qr_solver.py

```
1 import numpy as np
2
3 def qr_solver(x,y):
4   q,r=np.linalg.qr(x)
5   p = np.dot(q.T,y)
6   return np.dot(np.linalg.inv(r),p
   )
```

Of course, we don't need to implement our own OLS solvers in a production environment; but if you do, still you may find some well-written and well-tested functions such as `np.linalg.lstsq` to save your time and effort from doing it from scratch.

Ok, now we have finished the training part of a linear regression model in both R and Python. After we train a model we want to use it, i.e., to make predictions based on the model. For most of machine learning models, training is much more complex than prediction (Exceptions include Lazy-learning models such as KNN). Let's continue developing our own linear regression model by including the prediction function and enclose everything in an object.

R

chapter4/linear_regression.R

⁵¹ <https://github.com/numpy/numpy/issues/10471>

⁵² https://en.wikipedia.org/wiki/QR_decomposition

```

1
2 library(R6)
3 LR = R6Class(
4   "LR",
5   public = list(
6     coef = NULL,
7     initialize = function() {
8
9     },
10    fit = function(x, y) {
11      self$qr_solver(cbind(1, x), y)
12    },
13    qr_solver = function(x, y) {
14      self$coef = qr.coef(qr(x), y)
15    },
16    predict = function(new_x) {
17      cbind(1, new_x) %*% self$coef
18    }
19  )
20 )

```

Python

chapter4/linear_regression.py

```

1 import numpy as np
2
3 class LR:
4     def __init__(self):
5         self.coef = None
6
7     def qr_solver(self, x, y):
8         q, r = np.linalg.qr(x)
9         p = np.dot(q.T, y)
10        return np.dot(np.linalg.inv(r), p)
11
12    def fit(self, x, y):
13        self.coef = self.qr_solver(np.hstack((np.ones((x.shape[0], 1)), x)), y)
14
15    def predict(self, x):
16        return np.dot(np.hstack((np.ones((x.shape[0], 1)), x)), self.coef)

```

⁵³ <https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html>

Now let's try to use our linear regression model to solve a real regression problem with the Boston dataset⁵³, and check the results.

R

```

1 > source('linear_regression.R')
2 >
3 > library(MASS) # load Boston data from this package
4 >
5 > lr = LR$new()
6 > # -i means excluding the ith column from the data.frame
7 > lr$fit(data.matrix(Boston[,-ncol(Boston)]), Boston$medv)
8 > print(lr$coef)
9
10      crim      zn      indus      chas
11 3.645949e+01 -1.080114e-01 4.642046e-02 2.055863e-02 2.686734e+00
12      nox      rm      age      dis      rad
13 -1.776661e+01 3.809865e+00 6.922246e-04 -1.475567e+00 3.060495e-01
14      tax      ptratio      black      lstat
15 -1.233459e-02 -9.527472e-01 9.311683e-03 -5.247584e-01
16 > # let's make prediction on the same data
17 > pred=lr$predict(data.matrix(Boston[,-ncol(Boston)]))
18 > print(pred[1:5])
19 [1] 30.00384 25.02556 30.56760 28.60704 27.94352
20 >
21 > # compare it with the R built-in linear regression model
22 > rlm = lm(medv ~ ., data=Boston)
23 > print(rlm$coef)
24
25      (Intercept)      crim      zn      indus      chas
26 3.645949e+01 -1.080114e-01 4.642046e-02 2.055863e-02 2.686734e+00
27      nox      rm      age      dis      rad
28 -1.776661e+01 3.809865e+00 6.922246e-04 -1.475567e+00 3.060495e-01
29      tax      ptratio      black      lstat
30 -1.233459e-02 -9.527472e-01 9.311683e-03 -5.247584e-01
31 > print(rlm$fitted[1:5])
32
33      1      2      3      4      5
34 30.00384 25.02556 30.56760 28.60704 27.94352

```

Python

```

1 >>> from sklearn.datasets import load_boston
2 >>> from linear_regression import LR
3 >>> boston = load_boston()
4 >>> X, y = boston.data, boston.target

```

```

5 >>> # first, let's run our own linear regression
6 ... lr = LR()
7 >>> lr.fit(X, y)
8 >>> print(lr.coef)
9 [ 3.64594884e+01 -1.08011358e-01  4.64204584e-02  2.05586264e-02
10    2.68673382e+00 -1.77666112e+01  3.80986521e+00  6.92224640e-04
11   -1.47556685e+00  3.06049479e-01 -1.23345939e-02 -9.52747232e-01
12    9.31168327e-03 -5.24758378e-01]
13 >>> print(lr.predict(X)[:5])
14 [30.00384338 25.02556238 30.56759672 28.60703649 27.94352423]
15 >>>
16 >>> # now, use sklearn's linear regression model
17 ... from sklearn.linear_model import LinearRegression
18 >>> reg = LinearRegression().fit(X, y)
19 >>> print(reg.coef_)
20 [-1.08011358e-01  4.64204584e-02  2.05586264e-02  2.68673382e+00
21   -1.77666112e+01  3.80986521e+00  6.92224640e-04 -1.47556685e+00
22    3.06049479e-01 -1.23345939e-02 -9.52747232e-01  9.31168327e-03
23   -5.24758378e-01]
24 >>> print(reg.predict(X)[:5])
25 [30.00384338 25.02556238 30.56759672 28.60703649 27.94352423]

```

The results from our own linear regression models are almost identical to the results from `lm()` function or the `sklearn.linear_model` module, which means we have done a great job so far.

4.2 Linear hypothesis testing

I'm not a big fan of applying hypothesis testing in data science or machine learning. But sometimes it is irreplaceable, or at least useful, for example in the Design of Experiment (DOE).

Most of applications of hypothesis testing in machine learning models are about feature/variable selections. There are lots of debates on whether hypothesis-testings-based feature selections are good or bad. The major criticism of such approach is that the entire process is built on the data used for model training and the model performance on testing data is not considered at all.

I think it is still worth giving a brief introduction of hypothesis testing in linear regression, as it is still popular among data scientists with statistician's mindset. I would assume the readers already have basic ideas of hypothesis-testing, p-value, significance level.

If you have done linear regressions using a computer software (R, Stata, SPSS, Minitab etc.), you may have noticed that the outputs of these softwares contain the p-values⁵⁴ and t-statistics of the coefficient of each variable. If the p-value is less than a pre-determined significance level (usually 0.1 or 0.05 are used in practice), the null hypothesis (always denoted as H_0) should be rejected. The hypothesis against the null hypothesis is called alternative hypothesis (denoted as H_1). An example of H_0 and H_1 regarding model (4.2)

⁵⁴ https://www.statsdirect.com/help/basics/p_values.htm

could be stated as:

$$\begin{aligned} H_0 : \beta_1 &= 0 \\ H_1 : \beta_1 &\neq 0. \end{aligned} \tag{4.8}$$

If the p-value of β_1 suggests not to reject H_0 , we may exclude the corresponding feature and re-fit the linear model.

The theory behind the hypothesis testing in (4.8) is not complex. But we have to make an additional assumption for our linear regression model. More specifically, we assume the error term ϵ follows a multivariate normal distribution, i.e., $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$. Based on this assumption, we could conclude that $\hat{\beta}$ also follows a multivariate normal distribution. Furthermore, we can construct a test statistic⁵⁵ containing $\hat{\beta}_j$ which follows a t-distribution⁵⁶. The additional assumption we made for hypothesis testing is not required for least square estimation. It is also not one of assumptions for the Gauss-Markov theorem⁵⁷.

The example given in (4.8) focuses on a single coefficient. Actually, it is a special case of a more general hypothesis testing, which is called linear hypothesis testing. In linear hypothesis testing, people are interested in a linear combination of the coefficients, i.e.,

$$\begin{aligned} H_0 : \mathbf{A}\beta + \mathbf{b} &= 0 \\ H_1 : \mathbf{A}\beta + \mathbf{b} &\neq 0. \end{aligned} \tag{4.9}$$

The linear hypothesis testing is usually conducted by constructing a F-distributed test statistic. The general hypothesis testing is very powerful. For example, we may start from a full model with a larger list of variables to train a linear regression model; we may also exclude some variables from the full list to train a reduced model. By setting proper values of \mathbf{A} and \mathbf{b} in (4.9) we can conduct a linear hypothesis testing to accept or reject the reduced model. For example, if the full model involves three variables and in the reduced model we only keep the first variable, we will set

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

4.3 Ridge regression

What is the problem of solving linear regression model specified by (4.3)? There is nothing wrong at all with that approach. But I would like to cite my favorite quote - "Essentially, all models are wrong, but some are useful"⁵⁸. (4.3) provides one solution to model (4.1). There are some alternative solutions, such as lasso regression and ridge regression. In this section, let's focus on ridge regression.

What is ridge regression? Ridge regression doesn't change the model itself. Instead, it changes the way to estimate model parameters. By naive OLS, we minimize the SSR directly. In ridge regression, we change the objective function (which is commonly called loss function in machine learning models) by adding an

⁵⁵ https://en.wikipedia.org/wiki/Test_statistic

⁵⁶ https://en.wikipedia.org/wiki/Student%27s_t-distribution

⁵⁷ https://en.wikipedia.org/wiki/Gauss-Markov_theorem

⁵⁸ https://en.wikipedia.org/wiki/All_models_are_wrong

additional penalty

$$\min_{\hat{\beta}} e'e + \lambda \beta' \beta. \quad (4.10)$$

The optimization problem (4.10) is still an unconstrained optimization problem and it is convex. It can also be formulated as a constrained convex optimization problem as

$$\begin{aligned} \min_{\hat{\beta}} \quad & e'e \\ \text{subject to} \quad & \beta' \beta \leq t. \end{aligned} \quad (4.11)$$

The theory behind ridge regression can be found from The Elements of Statistical Learning [1]. Let's turn our attention to the implementation of ridge regression. The solution to (4.11) can be obtained in the same way as the solution to (4.3), i.e.,

$$\hat{\beta} = (X'X + \lambda I)^{-1} X'y. \quad (4.12)$$

Again, in practice we don't use (4.12) to implement ridge regression for the same reasons that we don't use (4.6) to solve linear regression without penalty.

Actually, we don't need new techniques to solve (4.10). Let's make some transformation on the objective function in (4.10):

$$e'e + \lambda \beta' \beta = \sum_{i=1}^n (y_i - x_i' \beta)^2 + \sum_{i=1}^p (0 - \sqrt{\lambda} \beta_i)^2 \quad (4.13)$$

, where $x_i = [1, X_{1i}, X_{2i}, \dots, X_{pi}]'$.

Let's define an augmented data set:

$$X_{\lambda} = \begin{bmatrix} 1 & X_{11} & X_{21} & \cdots & X_{p1} \\ 1 & X_{12} & X_{22} & \cdots & X_{p2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & \cdots & X_{pn} \\ \sqrt{\lambda} & 0 & 0 & \cdots & 0 \\ 0 & \sqrt{\lambda} & 0 & \cdots & 0 \\ 0 & 0 & \sqrt{\lambda} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \sqrt{\lambda} \end{bmatrix}, \quad y_{\lambda} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (4.14)$$

If we regress y_{λ} on X_{λ} , the OLS solution is just what we are looking after. However, usually the penalty

is not applied on the intercept. Thus, we modify y_λ and X_λ to

$$X_\lambda = \begin{bmatrix} 1 & X_{11} & X_{21} & \cdots & X_{p1} \\ 1 & X_{12} & X_{22} & \cdots & X_{p2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & \cdots & X_{pn} \\ 0 & \sqrt{\lambda} & 0 & \cdots & 0 \\ 0 & 0 & \sqrt{\lambda} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \sqrt{\lambda} \end{bmatrix}, \quad y_\lambda = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (4.15)$$

Now we are ready to implement our own ridge regression model based on the description above. It is also common to normalize⁵⁹ the independent variables before applying ridge regression to make different variables in the same order of magnitude.

R

chapter4/linear_regression_ridge.R

```

1 library(R6)
2 LR_Ridge = R6Class(
3   "LR_Ridge",
4   public = list(
5     coef = NULL,
6     mu = NULL,
7     sd = NULL,
8     lambda = NULL,
9     initialize = function(lambda) {
10       self$lambda = lambda
11     },
12     scale = function(x) {
13       self$mu = apply(x, 2, mean)
14       self$sd = apply(x, 2, function(e) {
15         sqrt((length(e) - 1) / length(e)) * sd(e)
16       })
17     },
18     transform = function(x) {
19       return(t((t(x) - self$mu) / self$sd))
20     },
21     fit = function(x, y) {
22       self$scale(x)
23       x_transformed = self$transform(x)
24       x_lambda = rbind(x_transformed, diag(rep(sqrt(self$lambda), ncol(x))))
25       y_lambda = c(y, rep(0, ncol(x)))

```

⁵⁹ [https://en.wikipedia.org/wiki/Normalization_\(statistics\)](https://en.wikipedia.org/wiki/Normalization_(statistics))

```

26     self$qr_solver(cbind(c(rep(1, nrow(
27         x
28     )), rep(0, ncol(
29         x
30     ))), x_lambda), y_lambda)
31 },
32 qr_solver = function(x, y) {
33     self$coef = qr.coef(qr(x), y)
34 },
35 predict = function(new_x) {
36     new_x_transformed = self$transform(new_x)
37     cbind(rep(1, nrow(new_x)), new_x_transformed) %*% self$coef
38 }
39 )
40 )

```

Python

chapter4/linear_regression_ridge.py

```

1  import numpy as np
2  from sklearn.preprocessing import StandardScaler
3
4
5  class LR_Ridge:
6      def __init__(self, l):
7          self.l = l
8          self.coef = None
9          self.scaler = StandardScaler()
10
11      def qr_solver(self, x, y):
12          q, r = np.linalg.qr(x)
13          p = np.dot(q.T, y)
14          return np.dot(np.linalg.inv(r), p)
15
16      def fit(self, x, y):
17          x_transformed = self.scaler.fit_transform(x)
18          x_lambda = np.vstack(
19              (x_transformed, np.diag([self.l**0.5]*x.shape[1])))
20          x_lambda = np.hstack(
21              (np.vstack((np.ones((x.shape[0], 1)), np.zeros((x.shape[1], 1)))),
22               x_lambda))
23          y_lambda = np.hstack((y, np.zeros((x.shape[1],))))
24          self.coef = self.qr_solver(x_lambda, y_lambda)

```

```

24
25     def predict(self, x):
26         new_x_transformed = self.scaler.transform(x)
27         new_x_transformed = np.hstack(
28             (np.ones((x.shape[0],1)), new_x_transformed)
29         )
30         return np.dot(new_x_transformed, self.coef)

```

in R, we implement our own scaler; but in the python implementation, we use StandardScaler from `sklearn.preprocessing` module, which is very handy. Please pay attention to how we calculate the standard deviation ⁶⁰ in R - we are using the formula of population standard deviation rather than the formula of sample standard deviation. Actually using which formula to calculate the standard deviation doesn't matter at all. I made the choice to use population standard deviation formula in order to generate consistent result of the StandardScaler since in StandardScaler uses the formula of population standard deviation.

The selection of best λ requires solving the OLS problem repeatedly with different values of λ , which implies the QR decomposition procedure would be called multiple times. Using SVD decomposition could be more efficient in terms of selection of best λ . But it wouldn't be covered in the current version of this book.

We are ready to run our own ridge regression on the Boston dataset.

R

```

1 > source('linear_regression_ridge.R')
2 >
3 > library(MASS) # load Boston data from this package
4 >
5 > # let's try lambda = 0.5
6 > ridge = LR_Ridge$new(0.5)
7 > ridge$fit(data.matrix(Boston[, -ncol(Boston)]), Boston$medv)
8 > print(ridge$coef)
9
10      crim      zn      indus      chas      nox
11 22.53280632 -0.92396151 1.07393055 0.12895159 0.68346136 -2.04275750
12      rm      age      dis      rad      tax      ptratio
13 2.67854971 0.01627328 -3.09063352 2.62636926 -2.04312573 -2.05646414
14      black      lstat
15 0.84905910 -3.73711409
16 > # let's make prediction on the same data
17 > pred=ridge$predict(data.matrix(Boston[, -ncol(Boston)]))
18 > print(pred[1:5])
[1] 30.01652 25.02429 30.56839 28.61521 27.95385

```

⁶⁰ https://en.wikipedia.org/wiki/Standard_deviation

Python

```

1 >>> from sklearn.datasets import load_boston
2 >>> from linear_regression_ridge import LR_Ridge
3 >>>
4 >>> boston = load_boston()
5 >>> X, y = boston.data, boston.target
6 >>> # first, let's run our own linear regression
7 ...
8 >>> ridge = LR_Ridge(0.5)
9 >>> ridge.fit(X, y)
10 >>> print(ridge.coef)
11 [ 2.25328063e+01 -9.23961511e-01  1.07393055e+00  1.28951591e-01
12    6.83461360e-01 -2.04275750e+00  2.67854971e+00  1.62732755e-02
13   -3.09063352e+00  2.62636926e+00 -2.04312573e+00 -2.05646414e+00
14    8.49059103e-01 -3.73711409e+00]
15 >>> print(ridge.predict(X)[:5])
16 [30.01652397 25.02429359 30.56839459 28.61520864 27.95385422]

```

It's exciting to see the outputs from R and Python are quite consistent.

Linear regression is not as simple as it seems. To learn more about it, I recommend to read [\[3\]](#).