

## STAT 2005 – PROGRAMMING LANGUAGES FOR STATISTICS

TUTORIAL 6 MATRIX, R PACKAGES

2020

LIU Ran

*Department of Statistics, The Chinese University of Hong Kong*

## 1 Matrix

### 1.1 Vector And Column Matrix

The difference between a vector and a column matrix:

```

1 # a vector is not a matrix
2 > a = 1:4
3 # we can use is.matrix to check whether or not the object is a matrix
4 > is.matrix(a)
5 [1] FALSE
6
7 # If we transform the vector to be a matrix, it will be a column matrix instead of a row matrix
8 > (b = as.matrix(a))
9      [,1]
10 [1,]  1
11 [2,]  2
12 [3,]  3
13 [4,]  4
14 > is.matrix(b)
15 [1] TRUE
16
17 # we can also use t transpose function to get a row matrix
18 > (c = t(a))
19      [,1] [,2] [,3] [,4]
20 [1,]  1   2   3   4
21 > is.matrix(c)
22 [1] TRUE
23
24 # t(t()) is equivalent to as.matrix() for a vector
25 > (d = t(t(a)))
26      [,1]
27 [1,]  1
28 [2,]  2
29 [3,]  3
30 [4,]  4
31 > is.matrix(d)
32 [1] TRUE

```

**Remark 1.1.** Vectors perform like column matrices but are not exactly the same. Be careful when doing matrix algebra using vectors.

### 1.2 diag

`diag` is a generic function like what I told you in the last tutorial. It will have different effects(methods) on different objects.

```

1 # input a vector
2 > diag(1:4)
3      [,1] [,2] [,3] [,4]
4 [1,]  1   0   0   0
5 [2,]  0   2   0   0
6 [3,]  0   0   3   0
7 [4,]  0   0   0   4
8
9 # input a scalar
10 > diag(4)

```

```

11      [,1] [,2] [,3] [,4]
12 [1,]    1    0    0    0
13 [2,]    0    1    0    0
14 [3,]    0    0    1    0
15 [4,]    0    0    0    1
16
17 > (a = matrix(1:12, ncol=4))
18      [,1] [,2] [,3] [,4]
19 [1,]    1    4    7   10
20 [2,]    2    5    8   11
21 [3,]    3    6    9   12
22 # input a matrix
23 # the matrix may not be an n*n square matrix
24 > diag(a)
25 [1] 1 5 9
26
27 # use the submatrix with the min dimension
28 > (b = matrix(1:12, nrow=4))
29      [,1] [,2] [,3]
30 [1,]    1    5    9
31 [2,]    2    6   10
32 [3,]    3    7   11
33 [4,]    4    8   12
34 > diag(b)
35 [1] 1 6 11
36
37 # replacement function
38 > diag(b) <- 1

```

### 1.3 Matrix Algebra

```

1 # use a vector to create a matrix
2 > a = 1:4
3 # if there are not enough elements, it will automatically replicate this vector several times to
  # match the shape of the matrix
4 > (b = matrix(a, nrow=3, ncol=4))
5      [,1] [,2] [,3] [,4]
6 [1,]    1    4    3    2
7 [2,]    2    1    4    3
8 [3,]    3    2    1    4
9
10 # each element will plus the scalar
11 > b+1
12      [,1] [,2] [,3] [,4]
13 [1,]    2    5    4    3
14 [2,]    3    2    5    4
15 [3,]    4    3    2    5
16
17 # also automatically replicate
18 > (c = matrix(2, nrow=3, ncol=4))
19      [,1] [,2] [,3] [,4]
20 [1,]    2    2    2    2
21 [2,]    2    2    2    2
22 [3,]    2    2    2    2
23
24 # same shape matrix
25 # the result will be the summation on each position
26 > b+c
27      [,1] [,2] [,3] [,4]
28 [1,]    3    6    5    4
29 [2,]    4    3    6    5
30 [3,]    5    4    3    6
31
32 # the result will be the product on each position
33 > b*c
34      [,1] [,2] [,3] [,4]
35 [1,]    2    8    6    4
36 [2,]    4    2    8    6
37 [3,]    6    4    2    8
38
39 # ncol(b) is 4; nrow(c) is 3; not match for the matrix multiplication
40 > b%*%c

```

```

41 Error in b %*% c : non-conformable arguments
42 > b%*%t(c)
43      [,1] [,2] [,3]
44 [1,]    20    20    20
45 [2,]    20    20    20
46 [3,]    20    20    20
47
48 # matrix + vector: vector will be transformed into the same shape matrix automatically first
49 > b+a # b+matrix(a, nrow=3, ncol=4)
50      [,1] [,2] [,3] [,4]
51 [1,]     2     8     6     4
52 [2,]     4     2     8     6
53 [3,]     6     4     2     8
54
55 # matrix * vector: vector will be transformed into the same shape matrix automatically first
56 > b*a # b*matrix(a, nrow=3, ncol=4)
57      [,1] [,2] [,3] [,4]
58 [1,]     1    16     9     4
59 [2,]     4     1    16     9
60 [3,]     9     4     1    16
61
62 # multiply a on the right
63 # regard the vector as a column matrix
64 # 3*4 matrix times 4*1 matrix = 3*1 matrix
65 > b%*%a
66      [,1]
67 [1,]    26
68 [2,]    28
69 [3,]    26
70
71 # multiply a on the left
72 # regard the vector as a row matrix
73 > a%*%t(b)
74      [,1] [,2] [,3]
75 [1,]    26    28    26
76
77 # create a matrix which has different dimensions
78 > (d = matrix(1:16, ncol=4))
79      [,1] [,2] [,3] [,4]
80 [1,]     1     5     9    13
81 [2,]     2     6    10    14
82 [3,]     3     7    11    15
83 [4,]     4     8    12    16
84
85 # return errors with different dimension matrices
86 > b+d
87 Error in b + d : non-conformable arrays
88 > b*d
89 Error in b * d : non-conformable arrays
90
91 # use the det function to calculate the determinant of a matrix
92 > det(d)
93 [1] 0
94 # d is a singular matrix, there is no inverse matrix for it
95 > solve(d) # try to get the inverse matrix
96 Error in solve.default(d) :
97   Lapack routine dgesv: system is exactly singular: U[3,3] = 0
98 > (e = matrix(1:4, ncol=2))
99      [,1] [,2]
100 [1,]     1     3
101 [2,]     2     4
102 > det(e)
103 [1] -2
104 > solve(e)
105      [,1] [,2]
106 [1,]    -2    1.5
107 [2,]     1   -0.5

```

**Remark 1.2.** Check dimensions before doing matrix algebra.

## 2 Mathematical calculation

### 2.1 Integration

You can use `integrate` command to get the integral of the function you input in a specific interval.

```

1 > integrate(dnorm, -1.96, 1.96)
2 0.9500042 with absolute error < 1e-11
3
4 # you can also set the infinity value in the interval
5 > integrate(dnorm, -Inf, Inf)
6 1 with absolute error < 9.4e-05
7
8 # integrate must accept vectorized function
9 # for example, a function which is not vectorized
10 > f <- function(x) 2.0 # just return the value 2.0
11 > integrate(f, 0, 1)
12 Error in integrate(f, 0, 1) :
13   evaluation of function gave a result of wrong length
14 > integrate(Vectorize(f), 0, 1) ## correct
15 2 with absolute error < 2.2e-14

```

**Remark 2.1.** Because we use numeral calculations instead of the theoretical values, it will show the absolute error here.

### 2.2 Derivatives

Compute derivatives of simple expressions:**D**

```

1 # define a expression
2 > (f = expression(x^2))
3 expression(x^2)
4 > D(f, "x") # it will show the expression for derivatives
5 2 * x
6
7 # use eval command to calculate the function value at point x=5
8 > x = 5 # assign a value to x
9 > eval(D(f, "x")) # use eval command to operate on the expression
10 [1] 10

```

**Remark 2.2.** Because here we want to get the theoretical formula of derivatives. We cannot use the `function` command to define the expression here, which means we cannot use `f<-function(x)x**2`.

**Remark 2.3.** `function` must do something(such as calculating). You must have the input, function body and the output. `expression` is just the expression. It does nothing and just show the structure of your formula.

## 3 Optimization

The function `optimize` searches the interval from lower to upper for a minimum of the function `f` with respect to its first argument.

```

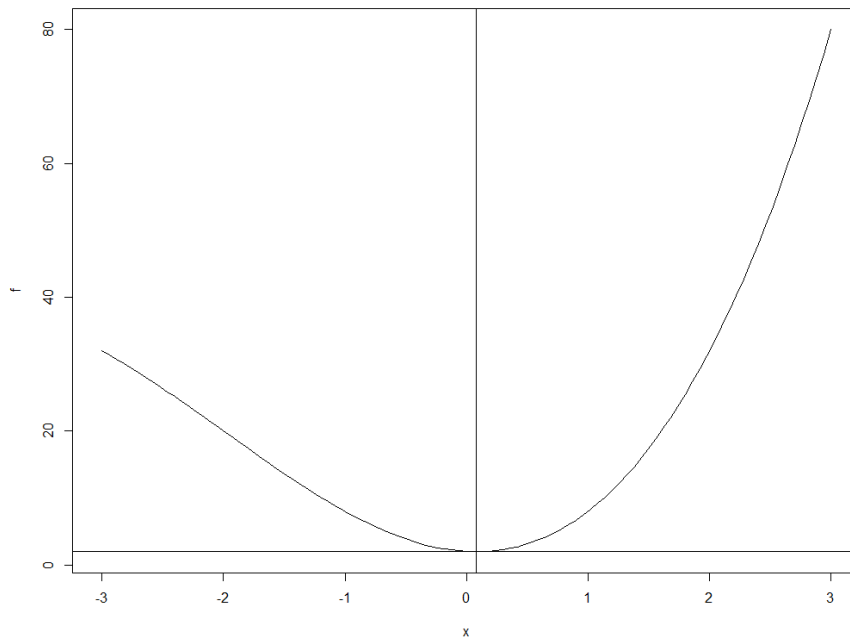
1 > f = function(x) {x^3+6*x^2-x+2}
2 > plot(f, -3, 3)
3 > (fit = optimize(f, c(-2,2), tol=0.00001))
4 $minimum # x coordinate
5 [1] 0.08166439
6
7 $objective # the function value or y coordinate
8 [1] 1.958895
9
10 # draw the point on the figure
11 > abline(v=fit$minimum)
12 > abline(h=fit$objective)

```

```

13 # we can input other arguments required for the function
14 # for example, two dimensional function f(x,y)
15 > f = function(x, y){x^3+6*x^2-x+2+y^2}
16 # optimize the f(x,y=2)
17 # y^2 = 4
18 > optimize(f, c(-2,2), tol=0.00001, y=2)
19 $minimum
20 [1] 0.08166439
21
22
23 $objective # these two objectives will differ by 4
24 [1] 5.958895

```



## 4 An exercise about the determinant

We can use the `det()` command to calculate the Determinant of a matrix.

In the case of a  $2 \times 2$  matrix the determinant may be defined as

$$|A| = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc.$$

Similarly, for a  $3 \times 3$  matrix  $A$ , its determinant is

$$\begin{aligned}
 |A| &= \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\
 &= aei + bfg + cdh - ceg - bdi - afh.
 \end{aligned}$$

Each determinant of a  $2 \times 2$  matrix in this equation is called a minor of the matrix  $A$ . This procedure can be extended to give a recursive definition for the determinant of an  $n \times n$  matrix. After the expansion, the dimension will reduce.

Therefore, we can use recursive code to calculate the determinant by ourselves.

```

1 > set.seed(2005)
2 >
3 > my_det<-function(A) {
4 +
5 +   # store the sign of one part
6 +   indicator = 1
7 +   result = 0
8 +
9 +   # stopping criterion
10 +  # if the number of the column is equal to 2, we end the recursion
11 +  if(ncol(A)==2) {
12 +    return(A[1,1]*A[2,2]-A[1,2]*A[2,1])
13 +  }
14 +
15 +  # calculate each part by the column order
16 +  for (j in 1:ncol(A)) {
17 +
18 +    # get the submatrix(minor)
19 +    temp_A = A[-1,-j]
20 +    result = result + indicator*A[1,j]*my_det(temp_A)
21 +
22 +    # change the sign of the indicator
23 +    # staggered sign
24 +    indicator = -indicator
25 +  }
26 +  return(result)
27 + }
28 >
29 > n = 4
30 > A = matrix(runif(n**2),nrow = n)
31 > my_det(A)
32 [1] 0.05846946
33 > det(A)
34 [1] 0.05846946

```

## 5 R package: dplyr(optional)

We can use `require` or `library` to load a specific R package. When you try to load a package which doesn't exist, `require` will give warning messages and return `FALSE`; `library` will return error and terminate the process.

```

1 > require('sssss')
2 Warning message:
3 In library(package, lib.loc = lib.loc, character.only = TRUE, logical.return = TRUE,  :
4 > library('sssss')
5 Error in library("sssss") :
6 > a = require('sssss')
7 > a
8 [1] FALSE

```

`dplyr` is a set of functions designed to enable dataframe manipulation in an intuitive, user-friendly way. Here are some main functions:

1. **filter**: filters out rows according to some conditions.
2. **arrange**: reorders rows according to some conditions
3. **select**: selects a subset of columns
4. **mutate**: adds a new column as a function of existing columns

`%>%` pipeline operator: All of the `dplyr` functions take a data frame as the first argument. Rather than forcing the user to either save intermediate objects or nest functions, `dplyr` provides the `%>%` operator. `x %>% f(y)` turns into `f(x, y)` so the result from one step is then “piped” into the next step. You can use the pipe to rewrite multiple operations that you can read left-to-right, top-to-bottom.

```

1 > # load the package
2 > require('dplyr')
3 > # Use the built-in data set iris
4 > str(iris)
5 'data.frame':  150 obs. of  5 variables:
6 $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
7 $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
8 $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
9 $ Petal.Width  : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
10 $ Species      : Factor w/  3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 1 ...
11 >
12 > # filter out rows: Species equal to virginica and Sepal.Length greater than 7.5
13 > iris %>%
14 +   filter(Species == 'virginica', Sepal.Length > 7.5)
15   Sepal.Length Sepal.Width Petal.Length Petal.Width  Species
16 1           7.6           3.0           6.6           2.1 virginica
17 2           7.7           3.8           6.7           2.2 virginica
18 3           7.7           2.6           6.9           2.3 virginica
19 4           7.7           2.8           6.7           2.0 virginica
20 5           7.9           3.8           6.4           2.0 virginica
21 6           7.7           3.0           6.1           2.3 virginica
22 >
23 > filter(iris, Species == 'virginica', Sepal.Length > 7.5)
24   Sepal.Length Sepal.Width Petal.Length Petal.Width  Species
25 1           7.6           3.0           6.6           2.1 virginica
26 2           7.7           3.8           6.7           2.2 virginica
27 3           7.7           2.6           6.9           2.3 virginica
28 4           7.7           2.8           6.7           2.0 virginica
29 5           7.9           3.8           6.4           2.0 virginica
30 6           7.7           3.0           6.1           2.3 virginica
31
32
33 # after filtering, sort the data by the column Sepal.Length and Sepal.Width
34 # first, compare the Sepal.Length, if they are equal, compare the Sepal.Width
35 # And finally create a new column which is the average value of two lengths
36 > iris %>%
37 +   filter(Species == 'virginica', Sepal.Length > 7.5) %>%
38 +   arrange(Sepal.Length, Sepal.Width) %>%
39 +   mutate(ave_length = (Sepal.Length+Petal.Length)/2)
40
41   Sepal.Length Sepal.Width Petal.Length Petal.Width  Species ave_length
42 1           7.6           3.0           6.6           2.1 virginica    7.10
43 2           7.7           2.6           6.9           2.3 virginica    7.30
44 3           7.7           2.8           6.7           2.0 virginica    7.20
45 4           7.7           3.0           6.1           2.3 virginica    6.90
46 5           7.7           3.8           6.7           2.2 virginica    7.20
47 6           7.9           3.8           6.4           2.0 virginica    7.15
48
49 # If you don't use the pipeline operator, it is not clear enthough.
50 > temp_data1 = filter(iris, Species == 'virginica', Sepal.Length > 7.5)
51 > temp_data2 = arrange(temp_data1, Sepal.Length, Sepal.Width)
52 > temp_data3 = mutate(temp_data2, ave_length = (Sepal.Length+Petal.Length)/2)

```

Remark 5.1. Tutorials for [dplyr](#) and [the vignettes](#).