

STATS 266 Handout - Data Structures & Manipulation

Qi Wang

2025-02-28

Contents

1	Introduction	2
2	Type of Variables	2
3	Data Structures	3
3.1	Atomic Vectors	3
3.2	Matrix	5
3.3	Array	7
3.4	Lists	7
3.5	Dataframe	8
4	Acknowledgement	10

1 Introduction

Welcome to **STATS 266: Introduction to R**. This handout provides an introduction about data structures and manipulation in R. By the end of this document, you should be able to:

- Identify the 5 main data types.
- Begin exploring data frames, and understand how they are related to vectors and lists.
- Be able to ask questions from R about the type, class, and structure of an object.
- Understand the information of the attributes “names”, “class”, and “dim”.

For this part, valuable materials to refer to include <http://adv-r.had.co.nz/Subsetting.html> and <https://swcarpentry.github.io/r-novice-gapminder/04-data-structures-part1.html>.

2 Type of Variables

There are five major types of atomic vectors: logical, integer, double, complex, and character. Integer and double are also known as numeric vectors.

```
v1 <- c(TRUE,FALSE)
typeof(v1)
```

```
## [1] "logical"
```

```
v2 <- c(1L,2L,3L)
typeof(v2)
```

```
## [1] "integer"
```

```
v3 <- c(1,2,3)
typeof(v3)
```

```
## [1] "double"
```

```
v4 <- c("string 1", "string 2")
typeof(v4)
```

```
## [1] "character"
```

```
v5 <- 1+1i
typeof(v5)
```

```
## [1] "complex"
```

3 Data Structures

Data in the real world will not exist as one single element. They always exist like a vector, matrix, or data frame and so on. So this section introduces the data structures in R and how we can subset them.

3.1 Atomic Vectors

3.1.1 What is an atomic vector?

An atomic vector in R is the simplest type of data structure that holds elements of the same data type. It is a one-dimensional structure where all elements belong to the same class (e.g., numeric, character, logical).

I believe most of you are now confused what does the function `c()` do here. It can be understood as a way to concatenate elements and creating a new vector.

To combine atomic vectors, the function `c()` is used. It flattens vectors, creating a new atomic vector containing all of the elements. When dealing with combining different types, we need to be careful:

```
c(1,0,TRUE,FALSE)
```

```
## [1] 1 0 1 0
```

```
c(1,0,'a','b')
```

```
## [1] "1" "0" "a" "b"
```

For missing values, R treats them as `NA` short for not applicable. Most calculations related to `NA` lead to `NA`, unless the results hold true for all possible values.

3.1.2 Subsetting an atomic vector

We use single bracket `[i]` to get the *i*-th element of the vector:

```
v1 <- c(2,4,6,8,10)
v1[1]
```

```
## [1] 2
```

```
v1[c(1,3)]
```

```
## [1] 2 6
```

```
v1[c(T,F,T,F,T)]
```

```
## [1] 2 6 10
```

More advanced cases include filtering some elements that we want. For example, if we want elements between -0.5 and 0.5 in a vector, what should we do?

```
v1 <- round(rnorm(100,0,0.5),2)
v1
```

```
## [1] 0.17 -0.30 -0.20 -0.02 -0.83 -0.48 0.22 -0.29 -0.24 -0.19 -0.88 0.86
## [13] 0.86 -0.37 -0.26 0.13 -0.12 0.57 -0.43 0.90 0.26 0.11 -0.26 -0.22
## [25] 0.74 -1.12 -0.08 -0.32 0.03 -0.27 0.23 -1.35 -0.07 0.25 -0.61 0.38
## [37] -0.10 0.72 0.10 0.56 -0.04 1.05 -0.65 0.01 0.34 0.82 -0.47 0.19
## [49] -0.87 0.14 0.31 -0.99 -0.60 0.04 -0.08 1.00 -0.38 0.43 -0.47 -0.56
## [61] 0.42 0.29 0.43 0.10 0.19 -0.29 -0.01 0.33 0.53 0.25 0.31 -0.38
## [73] 0.43 -0.17 -0.22 -0.18 -0.29 -0.28 0.45 0.76 -0.07 -0.15 0.02 0.05
## [85] -0.01 -0.04 -0.45 0.49 0.03 0.43 -0.42 0.68 0.23 -0.18 -0.76 0.95
## [97] 0.60 -0.41 -0.33 0.24
```

We want to find, which of the elements in the vector are those we want:

```
idx <- which(v1 > -0.5 & v1 < 0.5)
idx
```

```
## [1] 1 2 3 4 6 7 8 9 10 14 15 16 17 19 21 22 23 24 27
## [20] 28 29 30 31 33 34 36 37 39 41 44 45 47 48 50 51 54 55 57
## [39] 58 59 61 62 63 64 65 66 67 68 70 71 72 73 74 75 76 77 78
## [58] 79 81 82 83 84 85 86 87 88 89 90 91 93 94 98 99 100
```

These are indices! It returns to the position of the vector that we want. So we just subset the vector:

```
v1[idx]
```

```
## [1] 0.17 -0.30 -0.20 -0.02 -0.48 0.22 -0.29 -0.24 -0.19 -0.37 -0.26 0.13
## [13] -0.12 -0.43 0.26 0.11 -0.26 -0.22 -0.08 -0.32 0.03 -0.27 0.23 -0.07
## [25] 0.25 0.38 -0.10 0.10 -0.04 0.01 0.34 -0.47 0.19 0.14 0.31 0.04
## [37] -0.08 -0.38 0.43 -0.47 0.42 0.29 0.43 0.10 0.19 -0.29 -0.01 0.33
## [49] 0.25 0.31 -0.38 0.43 -0.17 -0.22 -0.18 -0.29 -0.28 0.45 -0.07 -0.15
## [61] 0.02 0.05 -0.01 -0.04 -0.45 0.49 0.03 0.43 -0.42 0.23 -0.18 -0.41
## [73] -0.33 0.24
```

Basic ideas here are to first find which positions in the series are those we want. We set filtering conditions, combine with `which` function to get the positions, then just subset the full vector.

3.2 Matrix

Note that vectors are just one dimensional, we can also create a matrix or tensor(array) in R when we have more than one dimension. For example, we can create a 4 by 4 matrix:

```
matrix(1:16, nrow = 4, ncol = 4)
```

```
##      [,1] [,2] [,3] [,4]
## [1,]    1    5    9   13
## [2,]    2    6   10   14
## [3,]    3    7   11   15
## [4,]    4    8   12   16
```

Did you notice that the elements are assigned by column? we can add an argument `byrow = TRUE` in the matrix function, to assign values by row. We can also bind two matrices by row or column:

```
m1 <- matrix(1:16, nrow = 4, ncol = 4)
m2 <- matrix(17:32, nrow = 4, ncol = 4)
cbind(m1,m2)
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## [1,]    1    5    9   13   17   21   25   29
## [2,]    2    6   10   14   18   22   26   30
## [3,]    3    7   11   15   19   23   27   31
## [4,]    4    8   12   16   20   24   28   32
```

```
rbind(m1,m2)
```

```
##      [,1] [,2] [,3] [,4]
## [1,]    1    5    9   13
## [2,]    2    6   10   14
## [3,]    3    7   11   15
## [4,]    4    8   12   16
## [5,]   17   21   25   29
## [6,]   18   22   26   30
## [7,]   19   23   27   31
## [8,]   20   24   28   32
```

In R, we have some functions to calculate the sum of the rows and columns too:

```
colSums(m1)
```

```
## [1] 10 26 42 58
```

```
rowSums(m1)
```

```
## [1] 28 32 36 40
```

More: In statistical analysis, we can do numeric matrix operations including multiplication, addition, inversion, transpose, eigen decomposition, determinant, and so on. Refer to: <http://www.philender.com/courses/multivariate/notes/matr.html> for more interesting operations.

Subsetting a Matrix To subset a matrix, there are three possible ways using bracket:

```
m1
```

```
##      [,1] [,2] [,3] [,4]
## [1,]    1    5    9   13
## [2,]    2    6   10   14
## [3,]    3    7   11   15
## [4,]    4    8   12   16
```

Subsetting a matrix by regarding it to be a vector: (**Dangerous**)

```
m1[1:5]
```

```
## [1] 1 2 3 4 5
```

Subsetting the specific row or column of the matrix:

```
m1[1,]
```

```
## [1] 1 5 9 13
```

```
m1[,1]
```

```
## [1] 1 2 3 4
```

Subsetting the specific some elements of the matrix:

```
m1[1,1]
```

```
## [1] 1
```

```
m1[1,1:2]
```

```
## [1] 1 5
```

We can also use `which()` function to subset the matrix, this is left for an exercise for you.

3.3 Array

Similar to matrix, if we have more than two dimensions of the data, we need to use an array. It's also known as tensor in deep learning literature:

```
a1 <- array(1:8, dim = c(2,2,2))
a1
```

```
## , , 1
##
##      [,1] [,2]
## [1,]    1    3
## [2,]    2    4
##
## , , 2
##
##      [,1] [,2]
## [1,]    5    7
## [2,]    6    8
```

It's like a 3D Lego now, right? We have each slices being a matrix, and put one slices above the other one. Subsetting here follows the similar rules, the only different thing from matrix is that now you need three dimensional coordinates to subset the ones we want.

3.4 Lists

A list is a collection of objects.

```
l1 <- list(1:3,
"a",
c(TRUE,FALSE,TRUE),
c(1,2)
)
l1
```

```
## [[1]]
## [1] 1 2 3
##
## [[2]]
## [1] "a"
##
## [[3]]
## [1] TRUE FALSE TRUE
##
## [[4]]
## [1] 1 2
```

```
str(l1)
```

```
## List of 4
## $ : int [1:3] 1 2 3
## $ : chr "a"
## $ : logi [1:3] TRUE FALSE TRUE
## $ : num [1:2] 1 2
```

Lists work differently from `c()`, as they can contain objects of different types. We can also assign names to the lists for each vector in the list:

```
names(l1) <- c("a", "b", "c", "d")
l1
```

```
## $a
## [1] 1 2 3
##
## $b
## [1] "a"
##
## $c
## [1] TRUE FALSE TRUE
##
## $d
## [1] 1 2
```

3.4.1 Subsetting a List

There are two ways to subset a list, by the order (just like atomic vectors), or by the name. But if we are going to subset by the order, double bracket (`[[]]`) is needed. To subset the list by name, we can use a dollar sign `$` + name.

```
l1[[1]]
```

```
## [1] 1 2 3
```

```
l1$c
```

```
## [1] TRUE FALSE TRUE
```

3.5 Dataframe

In data analysis, the lists that we frequently use are data frames and tibbles. Data frames are the lists to store the data for analysis.


```
df1 <- data.frame(x = 1:3, y = c(TRUE,FALSE,FALSE))
df1
```

```
##      x      y
## 1 1 TRUE
## 2 2 FALSE
## 3 3 FALSE
```

```
typeof(df1)
```

```
## [1] "list"
```

```
str(df1)
```

```
## 'data.frame':   3 obs. of  2 variables:
## $ x: int  1 2 3
## $ y: logi TRUE FALSE FALSE
```

```
attributes(df1)
```

```
## $names
## [1] "x" "y"
##
## $class
## [1] "data.frame"
##
## $row.names
## [1] 1 2 3
```

A dataframe can be converted to matrix, and it also works in the other direction:

```
a <- matrix(1:9, nrow = 3)
colnames(a) <- c("A", "B", "C")
a
```

```
##      A B C
## [1,] 1 4 7
## [2,] 2 5 8
## [3,] 3 6 9
```

```
is.data.frame(data.frame(a))
```

```
## [1] TRUE
```

```
is.matrix(a)
```

```
## [1] TRUE
```

3.5.1 Subsetting a dataframe

Either use the same way as subsetting the matrix, or using the dollar sign followed by the name of the column that we want:

```
df_a <- data.frame(a)
df_a$A
```

```
## [1] 1 2 3
```

```
df_a[1,]
```

```
##   A B C
## 1 1 4 7
```

We can also use the bracket combined with the column name to subset a dataframe:

```
df_a["A"]
```

```
##   A
## 1 1
## 2 2
## 3 3
```

4 Acknowledgement

This teaching material is adapted from the previous material of this course made by [Marcela Alfaro-Córdoba](#) and [Sheng Jiang](#).