

Introduction to Data Science

DSA1101

Semester 1, 2018/2019

Week 1, 14 August

Introduction to the R programming language



R is a programming language and software framework for statistical analysis and graphics distributed under the GNU general public license.

R can be downloaded from <http://cran.r-project.org/>
Linux/Mac OS X/Windows versions available

Basic commands in R

- R uses *functions* to perform operations

To run a function called `funcname`, we type `funcname(input1, input2, ...)` into the R console

- Example: the function `c()` concatenates the numbers inside the parentheses to form a vector

```
1 > x <- c(1, 3, 2, 5)
2 > x
3 [1] 1 3 2 5
```

Basic commands in R

- **Help file** is available by **typing ?funcname** into the console, for example

```
1 >? c
```

brings up the HTML page

```
c {base}
```

R Documentation

Combine Values into a Vector or List

Description

This is a generic function which combines its arguments.

The default method combines its arguments to form a vector. All arguments are coerced to a common type which is the type of the returned value, and all attributes except names are removed.

Usage

```
## S3 Generic function
c(...)

## Default S3 method:
c(..., recursive = FALSE, use.names = TRUE)
```

Arguments

...
objects to be concatenated.

recursive

Basic commands in R

- Addition, subtraction, multiplication and division are performed by `+`, `-`, `*` and `/` respectively
- Each element of a vector will be changed by operation with a constant, for example

```
1 > a <- c(1,2,3,4)
2 > a
3 [1] 1 2 3 4
4 > a + 5
5 [1] 6 7 8 9
6 > a - 10
7 [1] -9 -8 -7 -6
8 > a*4
9 [1] 4 8 12 16
10 > a/5
11 [1] 0.2 0.4 0.6 0.8
```

Basic commands in R

- Element-wise operations are possible for two vectors of the same length
- Length of a vector can be checked by the function `length()`

```
1 > x = c(1,6,2)
2 > x
3 [1] 1 6 2
4 > y = c(1,4,3)
5 > y
6 [1] 1 4 3
7 > length(x)
8 [1] 3
9 > length(y)
10 [1] 3
11 > x+y
12 [1] 2 10 5
```

Basic commands in R

- The function `sqrt()` returns the square root of each element of a vector
- The command `x^2` raises each element of `x` to the power 2; any powers are possible, including fractional or negative powers.

```
1 > x=c(2,3,9,16)
2 > sqrt(x)
3 [1] 1.414214 1.732051 3.000000 4.000000
4 > x^2
5 [1] 4 9 81 256
6 > x^3
7 [1] 8 27 729 4096
8 > x^-1
9 [1] 0.5000000 0.3333333 0.1111111 0.0625000
```

Basic commands in R

- We can quickly get a vector of numbers in sequence, $x = c(a_1, a_1 + 1, a_1 + 2, \dots, b)$, using $x = a : b$
- We can quickly get a vector of numbers in sequence with a specified step, $x = c(a_1, a_1 + d, a_1 + 2d, \dots, b)$, using $x = seq(a_1, b, d)$

```
1 > x = 1:5
2 > x
3 [1] 1 2 3 4 5
4 > x = 10:15
5 > x
6 [1] 10 11 12 13 14 15
7 > y = seq(1,10,2)
8 > y
9 [1] 1 3 5 7 9
```


Basic commands in R

- The function `ls()` provides a list of all the objects created so far
- The function `rm()` can delete any unwanted objects `rm` = remove

```
1 > ls()
2 [1] "a"  "x"  "y"
3 > rm(x,y)
4 > ls()
5 [1] "a"
6 > rm(a)
7 character (0)
```

- To remove all objects at once:

```
1 > rm(list=ls())
```

Basic statistics with R

Measures of central tendency



- Empirically speaking data shows a tendency to agglomerate around a central value. This is called central tendency. This central value can be used to summarize data. Suppose $x = c(x_1, x_2, \dots, x_N)$ is a given vector of numbers.
- The **mean** is computed as

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- The function **mean()** computes this value in R:

```
1 > x=c(1,3,5,7,10)
2 > mean(x)
3 [1] 5.2
4 > sum(x)/length(x)
5 [1] 5.2
```

Measures of central tendency

- The sample **median** is computed as

$$\text{median}(x) = \begin{cases} x_{(N+1)/2}, & \text{if } N \text{ is odd} \\ \frac{x_{N/2} + x_{N/2+1}}{2}, & \text{if } N \text{ is even} \end{cases} \quad (1)$$

- The function **median()** computes this value in R:

```
1 > x=c(1,3,5,7,10)
2 > median(x)
3 [1] 5
4 > x=c(1,3,5,7,10,11)
5 > median(x)
6 [1] 6
```

Measures of central tendency

- The sample **mode** is the most frequent value in a data set
- Note that the function **mode()** in R returns the type of object, not the sample mode
- R allows users to write their own functions:

```
1 #function to get mode
2 > getmode <- function(v) {
3 +   uniqv <- unique(v)
4 +   uniqv[which.max(tabulate(match(v, uniqv)))]
5 + }
6 > x=c(1,3,5,7,10,10,11)
7 > getmode(x)
8 [1] 10
```

Measures of dispersion

- The tendency for data points to spread around a central value is called dispersion.
- The **sample variance** is computed as

$$var(x) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

- The function **var()** computes this value in R:

```
1 > x=c(1,3,5,7,10,10,11)
2 > var(x)
3 [1] 14.90476
```

Measures of dispersion

- The **sample standard deviation** computed by function `sd()` is the square root of the variance

$$sd(x) = \sqrt{var(x)}$$

```
1 > x=c(1,3,5,7,10,10,11)
2 > var(x)
3 [1] 14.90476
4 > sqrt(var(x))
5 [1] 3.860669
6 > sd(x)
7 [1] 3.860669
```

- The **range** is the difference between the largest and smallest values

```
1 > x=c(1,3,5,7,10,10,11)
2 > max(x)-min(x)
3 [1] 10
```

Measures of association

- For two data vectors $x = c(x_1, x_2, \dots, x_N)$ and $y = c(y_1, y_2, \dots, y_N)$, their sample covariance can be computed by the function `cov()` and is given by

$$\text{cov}(x, y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

```
1 > x=c(2,4,7,1,10)
2 > y=c(1,6,3,10,9)
3 > cov(x,y)
4 [1] 1.95
5 #compute sample covariance from previous functions
6 > sum((x-mean(x))*(y-mean(y)))/(length(x)-1)
7 [1] 1.95
```


Measures of association

- The **sample correlation coefficient** for two data vectors can be computed by the function **cor()** and is given by

e.g height (cm) vs weight (kg)

$$r_{xy} = \frac{\text{cov}(x, y)}{\text{sd}(x)\text{sd}(y)}$$

- Unlike the sample covariance, the relation $-1 \leq r_{xy} \leq 1$ holds
- One can **compare the degree of association using r_{xy}** , which is invariant to scale and location changes

```
1 > x=c(2,4,7,1,10)
2 > y=c(1,6,3,10,9)
3 > cor(x,y)
4 [1] 0.1374092
5 > cov(x,y)/(sd(x)*sd(y))
6 [1] 0.1374092
7 #correlation invariant to scaling
8 > cor(100*x,100*y)
9 [1] 0.1374092
10 > cor(100*x,20*y)
11 [1] 0.1374092
```

Effects of scale and location changes

- Location changes: for two data vectors x and y , suppose we create two new data vectors $u = x + a$ and $v = y + b$, where a and b are constants
- Then the mean, median and mode for u will be different from those for x (similarly for v and y)
- Sample variance, standard deviation and covariance are invariant to changes in location

```
1 > x=c(2,2,7,1,10)
2 > y=c(1,6,3,10,10)
3 > u=x+10
4 > v=y-5
5 > stats_x=c(mean(x),median(x),getmode(x))
6 > stats_x
7 [1] 4.4 2.0 2.0
8 > stats_u=c(mean(u),median(u),getmode(u))
9 > stats_u
10 [1] 14.4 12.0 12.0
```

Effects of scale and location changes

```
1 > x_stats= c(var(x),sd(x),cov(x,y),cor(x,y))
2 > x_stats
3 [1] 15.3000000  3.9115214  3.2500000  0.2045482
4 > u_stats= c(var(u),sd(u),cov(u,v),cor(u,v))
5 [1] 15.3000000  3.9115214  3.2500000  0.2045482
```

- r_{xy} is invariant to location changes

Effects of scale and location changes

- Scale changes: for two data vectors x and y , suppose we create two new data vectors $u = a * x$ and $v = b * y$, where a and b are constants
- Then the mean, median and mode for u will be different from those for x (similarly for v and y)
- Sample variance, standard deviation and covariance between the two data vectors will also be different

```
1 > x=c(2,2,7,1,10)
2 > y=c(1,6,3,10,10)
3 > u=-5*x
4 > v=-10*y
5 > x_stats= c(mean(x),median(x),getmode(x))
6 > x_stats
7 [1] 4.4 2.0 2.0
8 > u_stats= c(mean(u),median(u),getmode(u))
9 > u_stats
10 [1] -22 -10 -10
```

Effects of scale and location changes

```
1 > x_stats= c(var(x),sd(x),cov(x,y),cor(x,y))
2 > x_stats
3 [1] 15.3000000  3.9115214  3.2500000  0.2045482
4 > u_stats= c(var(u),sd(u),cov(u,v),cor(u,v))
5 > u_stats
6 [1] 382.5000000  19.5576072 162.5000000
   0.2045482
```

- Notice that sample variance changes multiplicatively by the square of the scale, $var(u) = (-5)^2 \times var(x)$
- While sample standard deviation changes according to the absolute value of the scale, $sd(u) = |-5| \times sd(x)$
- r_{xy} is invariant to scale changes (provided the constants a and b are of the same sign)

Data manipulation with R

Data manipulation in R

- Besides the console, data from external sources can be read directly into R
- Example: a comma-separated-value (CSV) file containing the annual sales in U.S. dollars for 10,000 retail customers

Code **comments** are preceded by **'#'**

You can either **specify the complete path** to file in function **read.csv()**:

```
1 # import a CSV file of the total annual sales
2 sales <- read.csv("c:/data/yearly_sales.csv")
```

Or **set the working directory to the one containing the file** **yearly_sales.csv**

```
1 setwd("c:/data")
2 sales <- read.csv("yearly_sales.csv")
```

Data manipulation in R

- R offers other functions to read in data with different default values for column separator and decimal symbol

Function	Headers	Separator	Decimal Point
<code>read.table()</code>	FALSE	" "	"."
<code>read.csv()</code>	TRUE	" "	"."
<code>read.csv2()</code>	TRUE	";"	","
<code>read.delim()</code>	TRUE	"\t"	"."
<code>read.delim2()</code>	TRUE	"\t"	","

True: skip first line intro; number separator, ./,

Data manipulation in R

- Examine the imported data

The function `head()` displays the first few records in the dataset

```
1 # examine the imported dataset
2 head(sales)
3   cust_id  sales_total  num_of_orders  gender
4 1  100001         800.64             3      F
5 2  100002         217.53             3      F
6 3  100003          74.58             2      M
7 4  100004         498.60             3      M
8 5  100005         723.11             4      F
9 6  100006          69.43             2      F
```

Data manipulation in R

- Descriptive statistics

The function `summary()` provides summary statistics such as mean, median and quantiles

```
1 # examine the imported dataset
2 summary(sales)
3   cust_id      sales_total  num_of_orders  gender
4   Min.   :100001   Min.    : 30.02   Min.    : 1.000   F:5035
5   1st Qu.:102501   1st Qu. : 80.29   1st Qu. : 2.000   M:4965
6   Median :105001   Median  : 151.65   Median  : 2.000
7   Mean   :105001   Mean    : 249.46   Mean    : 2.428
8   3rd Qu.:107500   3rd Qu. : 295.50   3rd Qu. : 3.000
9   Max.   :110000   Max.    :7606.09   Max.    :22.000
```

Data manipulation in R

- Preliminary observations on data characteristics

The function `dim()` tells us that the data in `sales` has 10000 observations or rows, and 4 variables or columns

```
1 > dim(sales)
2 [1] 10000 4
```

The function `names()` tells us the variable names in the data set

```
1 > names(sales)
2 [1] "cust_id"      "sales_total"  "num_of_orders"
   "gender"
```

Data manipulation in R

- Preliminary observations on data characteristics

The function `dim()` tells us that the data in `sales` has 10000 observations or rows, and 4 variables or columns

```
1 > dim(sales)
2 [1] 10000      4
```

The function `names()` tells us the variable names in the data set

```
1 > names(sales)
2 [1] "cust_id"      "sales_total"  "num_of_orders"
   "gender"
```