

# ASSIGNMENT-4

2451-18-733-001

M. Ramani Priya

1. Write a program to find the factorial of a number using R.

```
factorial ← function(n) {  
  if (n == 0) return(1)  
  else  
    return (n * factorial(n-1))  
}
```

```
factorial  
n = 5  
answer ← factorial(n)  
print(answer).
```

2. Describe the data structures in R programming language.

Data structures in R,

1. atomic vector.
2. list
3. matrix
4. data frame.
5. factors.

1. Vector: It is the most common & basic data structure in R.

Vectors can be of 2 types:

(i) atomic vector.

(ii) list

A vector is a collection of elements that are most commonly of mode character, logical, integer or numeric.

→ `vector()` → creates an empty vector of default mode-logical.

→ Common approach to creating vectors is to use `character()`, `numeric()` etc.

2457-18-733-001  
→ vectors can be created directly by specifying their content. R will then guess appropriate mode of storage of vector.

$x \leftarrow c(1, 2, 3) \Rightarrow \text{mode-numeric}$

→  $x1 \leftarrow c(1L, 2L, 3L)$  - integer mode.

→  $z \leftarrow c("Sarah", "Tracy", "John")$  mode-character.

### Examining vectors.

the functions - `typeof()`, `length()`, `class()` & `str()` provide useful information about your vectors & R-objects.

### Adding elements.

the function `c()` (for combine) can also be used to add elements to vectors.

→ Vectors can also be created as a sequence of vectors

`uSeries <- 1:10`

`seq(10)`

→ 1 2 3 4 5 6 7 8 9 10

→ R supports missing data in vectors. They are represented by NA - not available.

The function `isna()` indicates the elements of vectors that represent missing data, & the function `anyNA()` returns TRUE if vector contains any missing values.

→ Inf is infinity.

## 2. Lists

A list is a generic object consisting of an ordered collection of objects. Lists are heterogeneous data structures.

These are also one-dimensional data structures.

A list can be a list of vectors, list of matrices, a list of characters or a list of functions & so on.

Eg: `empld = c(1, 2, 3, 4)`

`empName = c("D", "C", "B", "A")`

`numberOfEmp = 4`

`emplist = list(empld, empName, numberOfEmp)`

⇒ `emplist` ⇒

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

[[2]]
[1] "D" "C" "B" "A"

[[3]]
[1] 4
```

## 3. Dataframes

Dataframes are generic data objects of R which are used to store the tabular data. Dataframes are the foremost popular data objects in R programming because we are comfortable in seeing the data within the tabular form.

They are 2 dimensional, heterogeneous data structures. These are lists of vectors of equal length.

Dataframes have the following constraints placed on them:

1. A data frame must have column name & every row should have a unique name.

2. Each column must have the identical number of items.
3. Each item in a single column must be of same data type.
4. Different columns <sup>may</sup> have different data types

To create a dataframe we use `data.frame()` function.

Eq: `a = c("b", "c", "d")`

`x = c("p", "q", "r")`

`y = c(22, 25, 28)`

`df = data.frame(a, x, y)`

df  $\Rightarrow$

|   | a | x | y  |
|---|---|---|----|
| 1 | b | p | 22 |
| 2 | c | q | 25 |
| 3 | d | r | 28 |

### 3. Matrix

A matrix is a rectangular arrangement of numbers in rows & columns. In a matrix, the rows are ones that run horizontally & columns run vertically. Matrices are 2 dimensional, homogenous data structures.

Matrices can be created using `matrix()` function.

The arguments of `matrix()` are set of elements in the vector.

By default, matrices are in column wise order.

Eq: `A = matrix(c(1, 2, 3, 4, 5, 6, 7, 8, 9), nrow=3, ncol=3, byrow=TRUE)`

A  $\Rightarrow$

|      | [,1] | [,2] | [,3] |
|------|------|------|------|
| [1,] | 1    | 2    | 3    |
| [2,] | 4    | 5    | 6    |
| [3,] | 7    | 8    | 9    |



## 5. Arrays

Arrays are the R objects which store the data in more than 2 dimensions. Arrays are n dimensional data structures.

For example, if we create an array of dimensions (2, 3, 3) then it creates 3 rectangular matrices each with 2 rows & 3 columns. They are homogeneous data structures.

Arrays can be created using `array()` function.

The arguments to `array()` are the set of elements in vectors & you have to pass a vector containing dimensions of array.

Eg: `a = array (`  
`c(1, 2, 3, 4, 5, 6, 7, 8),`  
`dim = c(2, 2, 2)`  
`)`

, , 1

|       | [, 1] | [, 2] |       | [, 1] | [, 2] |
|-------|-------|-------|-------|-------|-------|
| [1, ] | 1     | 3     | [1, ] | 5     | 7     |
| [2, ] | 2     | 4     | [2, ] | 6     | 8     |

## 6. Factors

Factors are the data objects which are used to categorize the data & store it as levels. They are useful for storing categorical data. They can store both strings & integers. They are useful in data analysis for statistical modeling.

Factors can be created using a function - `factor()`.

The argument to `factor()` is a vector.

Factors are useful to categorize unique values in columns like "TRUE" or "FALSE", ~~"MALE"~~ "MALE" or "FEMALE". etc.

Eg: `fac = factor (c("Male", "Female", "Male", "Male", "Female"))`

`fac`  $\Rightarrow$  [1] Male Female Male Male Female

— Levels: Female Male

3. Define objects. List the methods for measuring distance between objects.

Objects: Objects are the instances of <sup>the</sup> class. Also, everything in R is an object & they can have their attributes like class, attributes, idimnames, names etc.

In R programming,

→ `stats::dist()` is a <sup>base</sup> method used to calcul get distance between objects.

however this method can compute only following distances, euclidean, maximum, manhattan, canberra, binary & minkowski.

→ `distance()` function is implemented using same logic as `stats::dist()` & takes as input matrix or data frame.

The corresponding matrix & data frame should store probability density functions (as rows) for which distance computations should be performed.

→ When defining matrix the probability vectors should be combined as rows using `rbind()`.

→ `distance()` function allows to choose from 46 distance/similarity measures.

→ `getDistMethods()` can give information about which methods are implemented in `distance()`.

→ `distance()` returns a symmetric matrix whereas `stats::dist()` returns only a part of matrix.

The arguments of `distance()` are matrix, method, use.rows.names (boolean), as.dist.obj (boolean).

Eg:  $P \leftarrow 1:10 / \text{sum}(1:10)$

$Q \leftarrow 20:29 / \text{sum}(20:29)$

$x \leftarrow \text{rbind}(P, Q)$  // combine P & Q as matrix object.

$\text{distance}(x, \text{method} = "euclidean", \text{use.row.names} = \text{TRUE},$   
 $\text{as.dist.obj} = \text{TRUE})$

$\text{stats} :: \text{dist}(x, \text{method} = "euclidean")$ .

Methods for measuring distances:

### 1. Euclidean Distance

The most common distance is Euclidean distance. The euclidean distance between two vectors  $x$  &  $y$  is defined as

$$\text{distance}(x, y) \leftarrow \text{sqr}t((x[1] - y[1])^2 + (x[2] - y[2])^2 + \dots)$$

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

### 2. Manhattan distance (city block)

Manhattan distance measure distance in the number of horizontal & vertical units it takes to get from one (real valued) point to another (no diagonal moves).

$$\text{distance}(x, y) \leftarrow \text{sum}(\text{abs}(x[1] - y[1]) + \text{abs}(x[2] - y[2]) + \dots)$$

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

### 3. Cosine similarity

is a common similarity metric in text analysis.

It measures the smallest angle between 2 vectors & is assumed to be between  $0$  &  $90^\circ$ .

Two perpendicular vectors ( $\theta = 90^\circ$ ) are the most dissimilar,



the cosine of  $90^\circ$  is 0.

Two parallel vectors are similar,  $\cos(0) = 1$

$$\text{dot}(x, y) \leftarrow \text{sum}(x[1] * y[1] + x[2] * y[2] + \dots)$$

$$\text{cosSim}(x, y) \leftarrow \text{dot}(x, y) / (\sqrt{\text{dot}(x, x) * \text{dot}(y, y)})$$

$$\cos \theta = \frac{A \cdot B}{|A| |B|}$$

4. List out the various control structures supported by R programming language.

There are 8 control statements supported by R programming language. Control statements are expressions used to control the execution & flow of program based on the conditions provided in statements. The statements are used to make a decision after assessing the variable.

1. if condition

This control statement structure contains checks the expression provided in parenthesis is true or not. If true, the execution of statements in braces {} continues.

Syntax: `if(expression){`

`Statements`

`....`  
`....`

`}`

2. if-else condition

It is similar to if condition but when the test expression in if condition fails, then statements in else condition are executed.



Syntax: if (expression) {  
 Statements  
 ....  
 ....  
 }  
 else {  
 statements  
 ....  
 ....  
 }

### 3. for loop.

It is a type of loop or sequence of statements executed repeatedly until exit condition is reached.

Syntax: for (value in vector) {  
 Statements  
 ....  
 ....  
 }

### 4. Nested loops

Nested loops are similar to simple loops. Nested means loops inside loop. Moreover, nested loops are used to manipulate the matrix.

Eg:  $m \leftarrow \text{matrix}(2:15, 2)$   
 for (r in seq(nrow(m))) {  
 for (c in seq(ncol(m))) {  
 print(m[r, c])  
 }  
 }

### 5. while loop

while loop is another kind of loop iterated until a condition is satisfied. The testing expression is checked first before executing the body of loop.

Syntax: while (expression) {  
 statement  
 ....  
 ....  
 }

## 6. repeat loop &amp; break statement

repeat is a loop which can be iterated many number of times but there is no exit condition to come out from the loop. So, break statement is used to exit from loop. break statement can be used in any type of loop to exit from loop.

Syntax: repeat {

statement

....

....

if (expression) {

break

}

}

## 7. return statement.

return statement is used to return the result of an executed function & returns control to the calling function.

Syntax: return (expression)

## 8. next statement.

next statement is used to skip the current iteration without executing further statements & continues the next iteration cycle without terminating the loop.

Eg:  $x \leftarrow 1:10$

for (i in x) {

if (i % 2 != 0) {

next # jumps to next loop

}

print(i)

}

5) Same as 2nd.

2451-18-733-001  
6. Define list & Data frame in R & explain various operations on lists & dataframes with suitable examples.

### Lists:

Lists are R objects which contain elements of different types like - numbers, strings, vectors & another list inside it. A list can contain a matrix or a function as its elements.

- List is created using `list()` function.
- The list elements can be given names & can be accessed using these names using `names()` function.
- The elements of list can be accessed by index of element in list. In case of named list it can be also be accessed using names.
- We can add, delete, update list elements.  
We can add & delete only at beginning or end of list.  
But we can update any element of list.
- We can merge many lists into one list by placing all the lists in one `list()` function.
- A list can be converted to vector so that the elements of the vector can be used for further manipulation. All the arithmetic operations on vectors can be applied after the list is converted into vector. To do this conversion, we use `unlist()` function.

Eg: `list_data ← list(c("a", "b", "c"), matrix(c(3, 9, 5, -1, 2, 8),  
row = 2),  
list("g", 2, 3))`

`names(list_data) ← c("1st quarter", "Matrix", "Inner list")`



```
> print(list_data[1])
```

```
1st Quarter
```

```
"a" "b" "c"
```

```
> list_data[4] <- "New element"
```

```
> # Merging Lists
```

```
list1 <- list(1, 2, 3)
```

```
list2 <- list("A", "B", "C")
```

```
merged_list <- c(list1, list2)
```

```
> # Converting list to vector.
```

```
list1 <- list(1:5)
```

```
list2 <- list(10:14)
```

```
v1 <- unlist(list1)
```

```
v2 <- unlist(list2)
```

```
result <- v1 + v2
```

## Data frames

Data frames are generic data objects of R which are used to store the tabular data. Data frames are considered to be the most popular data objects in R objects because it is more comfortable to analyze the data in the tabular forms.

Operations that can be performed using Dataframes are,

### 1. Creating Dataframe

(1) Creating data frame from vectors:

We can use data.frame() function.

Eg: name = c("a", "b", "c")

language = c("R", "P", "J")

age = c(22, 25, 28)

df = data.frame(name, language, age)

(ii) Creating dataframe using data from a file:

Dataframes can also be created by importing the data from a file.

newDF = read.table(path = "Path of file")

newDF = read.csv(path = "Path of file")

2. Accessing rows and columns.

df[val1, val2]

df = dataframe object

val1 = rows of a dataframe

val2 = columns of dataframe

accessing rows.,

df[1:2, ]

accessing columns

df[, 1:2]

3. Selecting subset of Dataframe.

A subset can also be created based on certain conditions

newDF = subset(df, col conditions)

df = Original dataframe

conditions = certain conditions.

4. Editing Dataframes.

Dataframes can be edited in 2 ways,

(i) Editing dataframes by direct assignments: Much like the list in R you can edit the data frames by a direct assignment.

(ii) Editing dataframes using edit() command.

Step 1: create an empty instance of dataframe.

Step 2: use edit function to launch viewer.

`myTable = data.frame()`

`myTable = edit(myTable)`

Step 3: enter the values in table & close the editor

Step 4: Check resulting table.

`> myTable`

5. Adding rows & columns to dataframe.

(i) Adding rows:

function: `rbind(df, the entries of row)`

`newDF = rbind(df, data.frame(name = "d",  
language = "C",  
age = 40))`

(ii) Adding extra columns.

function: `cbind(df, the entries of column)`

`newDF = cbind(df, Rank = c(3, 5, 4, 1))`

6. Adding new variables to dataframes.

We can add variables to dataframe based on existing ones.

To do that we have to first call the `dplyr` library.

using command `library()`. Then call `mutate()` function.

will add extra variable columns based on existing ones

`library(dplyr)`

`newDF = mutate(df, new-var = [existing-var])`

`df` = original dataframe.

`new-var` = Name of new variable.

`existing-var` = The modify action you are taking (e.g. `log` value, multiply by 10).



7. Deleting rows & columns from a data frame. 2451-18-733-001

To delete a row or column you need to access that row or column & then insert a negative sign before that row or column. It indicates that you need to delete that row or column.

$\text{newDF} = \text{df}[-\text{rowNo}, -\text{colNo}]$