Introduction to Data Science (11372 & G 11516)

Semester 1 2021



INTRODUCTION TO DATA SCIENCE

Lecture 4

Dr. Ibrahim Radwan

DISTINCTIVE BY DESIGN

OUTLINE



- Data Structures
 - Factors
 - Arrays and Matrices
 - Lists
 - Data Frames

DATA STRUCTURES



Similar type of items

Homogeneous data structures

Atomic vector 1-D

Matrix 2-D

Array n-D

Dissimilar type of items

Heterogeneous data structures

List 1-D

Data frame 2-D

FACTOR



- Special case of vector, used to store nominal/categorical data
- A factor can be thought as an integer vector, where each integer has a label
 - It is more efficient to have a variable with "Female" and "Male" values than a variable with 1 and 2 values
 - Integers are more efficient when comparing values or when searching for a value than using characters
 - However, integers are not self describing and characters are
 - Factor is the solution that integrates both
- Creating a factor

FACTOR (2)

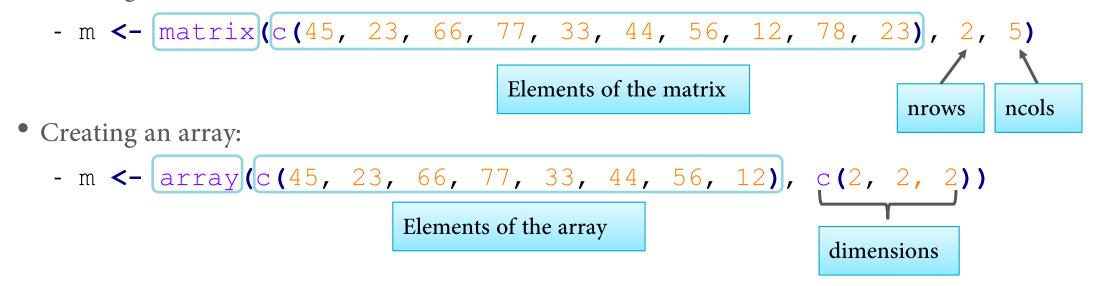


```
# Factor example
          # Create character vector for types of the properties
          houses.types <- c("House", "Unit", "House", "Unit")
          houses.types
          str(houses.types)
          # Create integral vector for types of the properties
          houses.types \leftarrow c(1L, 2L, 1L, 2L)
          houses.types
          str(houses.types)
      11
      12 # create a factor
Code Example
          houses.types <- factor(c("House", "Unit", "House", "Unit"))
          houses.types
      15 # underlying representation of factor
          unclass(houses.types)
          str(houses.types)
      18 # Searching using characters
          houses.types[houses.types == "House"] # not effecient
      20 # Searching using integer
          houses.types[as.numeric(houses.types) == 1] # effecient
      22
      23
          employees.gender <- rep(c("female", "male", "male", "female", "male"), 10)</pre>
          employees.gender
          employees <- factor(employees.gender)</pre>
          str(employees)
      28
          # order of levels
          x <- factor(c("yes", "no", "yes", "no"))</pre>
      31
          y <- factor(c("yes", "no", "yes", "no"), levels= c("yes", "no"))</pre>
      33
```

ARRAYS AND MATRICES



- Arrays store data elements in several dimensions. Matrices are special cases of arrays with only two dimensions. Arrays and Matrices in R are nothing more than vectors with a attribute that is the dimension.
- Creating a matrix:



ARRAYS AND MATRICES (2)

arrays

my.array

my.matrix



```
# Q: how are the items ordered?

# A: It is ordered column-wise

# change the dimensions of a vector in R

# my.vector <- 1:24

# dim(my.vector) <- c(3,4,2)

# check that the contents of two objects are equal using identical function identical(my.array, my.vector)

# matrices
```

[my.matrix <- matrix(1:24, nrow= 4, ncol= 6)]

my.array \leftarrow array(1:24, dim=c(3,4,2))

16

example for sub-setting

ARRAYS AND MATRICES (3)



```
# Subsetting arrays and matrices
         m \leftarrow matrix(c(45, 23, 66, 77, 33, 44, 56, 12, 78, 23, 38, 17), nrow= 3, ncol=4)
         # extract the element of the second row and third column
      5 m[2, 3]
        # extract all elements of the second row
      7 m[2,]
      8 # extract all elements of the third column
      9 m[,3]
     10 # extract all element of the second row except the the third element of this row
     11 m[2, -3]
     12 # extract the second and fourth elements of the second row
     13 m[2, c(2, 4)]
        # extract all elements of the second row except the second and fourth elements
     15 m[2, -c(2, 4)]
     16
     17 # all of the above are either column or row, so the results were vectors
     18 # The results can be a sub-matrix
        a \leftarrow m[c(1,2), c(2, 4)]
Code
     21 class(a)
     22 # If you want the result of subsetting to be a matrix even for row and column
         b <- m[1, ,drop= FALSE] # result in matrix
     24
     25 class(b)
```

ARRAYS AND MATRICES (4)



• You can use *cbind()* and *rbind()* functions to join two or more vectors together or matrices, by columns or by rows, respectively. The following examples illustrates this:

```
# rbind and cbind example
m1 <- matrix(c(45, 23, 66, 77, 33, 44, 56, 12, 78, 23, 38, 17), nrow=3, ncol=4)
m1
res <- cbind(c(4, 76, 12), m1[, 4])
res
res
res <- rbind(c(4, 76, 12, 14), m1[3, ])
res
res
m2
res <- matrix(rep(10, 16), nrow= 4, ncol= 4)
m2
m3 <- rbind(m1[1, ], m2[3, ])
m3
```

ARRAYS AND MATRICES (5)



- You can also give names to the columns and rows of matrices, using the functions
- *colnames()* and *rownames()*. This facilitates memorizing the data positions.

```
# using rownames() and colnames() functions
        results \leftarrow matrix(c(10, 30, 40, 50, 43, 56, 21, 30), 2, 4, byrow= TRUE)
       # names of the columns
       colnames(results) <- c("1st_qrt", "2nd_qrt", "3rd_qrt", "4th_qrt")
Sode example
                                                                     You can also use
       # names of the rows
       rownames(results) <- c("store1", "store2")
                                                                  dim(matrix) to get the
       # show the matrix
                                                                   dimensions of arrays
   11
       results
   12
                                                                       or matrices
   13
       # indexing using the names of the columns/rows
   14
       results["store1", ]
       results["store2", c("1st_qrt", "4th_qrt")]
```

LISTS



- A list is a data structure that holds collection of possibly unrelated (or heterogeneous) objects in order under one name.
 - Think of a vector, when you would like to combine the name, age, salary of an employee in one vector, is it possible?
 - No, most likely all of the numeric values (e.g. age and salary) will be converted to character to fulfill the rule that the items of a vector have to be from the same data type
 - However, with the lists, this is possible:
 - > employee <- list(name='Jack', age=25L, salary=5210.43)</pre>

Different type of objects are concatenated together

• Combining two lists will results in a list

LISTS (2)



```
# vectors
        # Character
        houses.addresses <- c("7 George st", "18/5 Irwan cresent", "8 Morad close", "1/2 London Circuit")
        # Numeric
        houses.area <- c(420.5, 220.15, 750.4, 120.5)
        # Integer
        houses.bedrooms <- c(4L, 3L, 5L, 2L)
        # Logical
        houses.has.garden <- c(TRUE, FALSE, TRUE, FALSE)
     10
Example
     11
        # Combine the features of each proprty
        # using vector
     13
         houses.sample <- c(address= houses.addresses[1], area=houses.area[1], bedrooms= houses.bedrooms[1],
    14
                            has_garden=houses.has.garden[1])
    15
        print(houses.sample)
        mode(houses.sample)
    16
Code
        houses.sample[1]
     18
        houses.sample[2]
     19
        houses.sample['bedrooms']
     20
        str(houses.sample)
     21
        # using list
         houses.sample <- list(address= houses.addresses[1], area=houses.area[1], bedrooms= houses.bedrooms[1],
     23
                               has_garden=houses.has.garden[1])
     24
         print(houses.sample)
     25
        mode(houses.sample)
        houses.sample[1]
     26
        houses.sample[2]
        houses.sample['bedrooms']
     28
     29
        houses.sample$bedrooms
        mode(houses.sample$bedrooms)
     31
        str(houses.sample)
```

DATA FRAMES



- A data frame is a series of records represented by rows (observations), where each row contains values in several fields/columns (variables).
- They are like matrices in structure as they are also bi-dimensional.
 - All the operations, we have used for matrices can be applied on data frames as well
 - such as rbind(), cbind(), dim(), ...
- However, contrary to matrices, data frames may include data of a different type in each column.
- So, a data frame is a special type of lists, where a list represents only one row (in-order) and a data frame can be one row or more.

DATA FRAMES (2)



ID, Name, Age

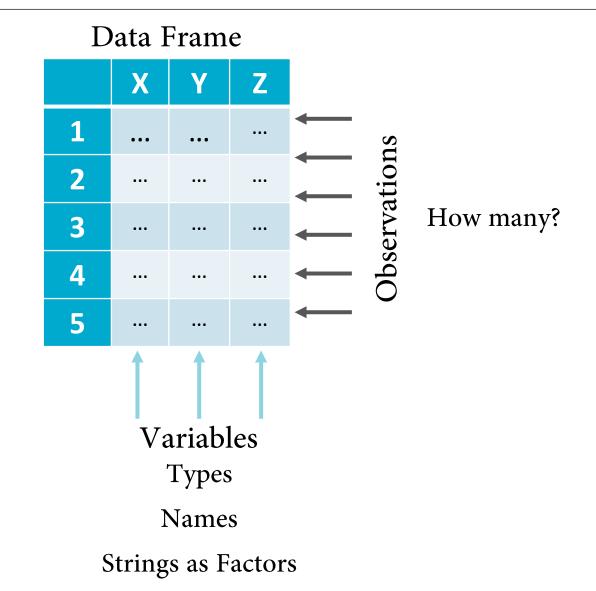
23424, Ana, 45

11234, Charles, 23

77654, Susanne, 76

data.csv

How to read this using the base functions in R?



DATA FRAMES (3)



• Create data frames from vectors:

```
# create dataframe from vectors
employee <- c('John Doe', 'Peter Gynn', 'Jolie Hope')
salary <- c(21000, 23400, 26800)
startdate <- as.Date(c('2010-11-1', '2008-3-25', '2007-3-14'))
employ.data <- data.frame(employee, salary, startdate)</pre>
```

Check structure of data frame

```
str(employ.data)
```

• Characters by the default are set as Factors, however this can be changed using:

```
employ.data <- data.frame(employee, salary, startdate,
stringsAsFactors = FALSE)</pre>
```

• Full example in the next slide

DATA FRAMES (4)



```
# create dataframe from vectors
       employee <- c('John Doe','Peter Gynn','Jolie Hope')</pre>
       salary <- c(21000, 23400, 26800)
       startdate <- as.Date(c('2010-11-1','2008-3-25','2007-3-14'))
       employ.data <- data.frame(employee, salary, startdate)</pre>
       # check structure of data frame
       str(employ.data)
Example
       # keep characters as characters
       employ.data <- data.frame(employee, salary, startdate, stringsAsFactors = FALSE)</pre>
       str(employ.data)
   13
Code
   14
       # using the accessor $ to access the variables of the data frame
   15
       employ.data$employee
   16
       employ.data$salary
   17
       employ.data$startdate[1]
   18
       # what is the difference between the following expressions
   20
       aa <- employ.data["salary"]
       |bb <- employ.data[["salary"]]
   21
   22
       aa
   23
       class(aa)
   24
       bb
   25
       class(bb)
```

©Dr. Ibrahim Radwan – University of Canberra

DATA FRAMES (5)



• Sub-setting data frame using [] or accessor operator, \$:

• Sub-setting based on conditions:

```
my.dataset[my.dataset$pH > 7, ]
my.dataset[my.dataset$pH > 7, "site"]
my.dataset[my.dataset$season == "Summer", c("site", "pH")]
```

- Is it possible to refer to the columns directly?
 - Like we say, my.dataset[pH > 7,]
 - No this is not possible unless we attach. Let us discuss the full example in the next slide.

DATA FRAMES (6)



```
1 # create and manipulate data frame
        my.dataset <- data.frame(site=c("A", "B", "A", "A", "B"),
                                  season=c("Winter", "Summer", "Summer", "Spring", "Fall"),
                                  pH=c(7.4,6.3,8.6,7.2,8.9), stringsAsFactors = FALSE)
        my.dataset
      6 my.dataset[3, 2]
        my.dataset$pH
     8 # subsetting with conditions
     9 my.dataset[my.dataset$pH > 7, ]
    10 my.dataset[my.dataset$pH > 7, "site"]
Code Example
     11 my.dataset[my.dataset$season == "Summer", c("site", "pH")]
    12 # Is it possible to refer to the columns directly?
    13 my.dataset[pH > 7, ]
    14 # This is only possible if you `attach` the dataset
    15 attach(my.dataset)
    16 my.dataset[pH > 7, ]
    17 season
    18 my.dataset[site=='B',]
    19 # To go back, use 'detach'
     20 detach(my.dataset)
    21 season
     22 # it is much safer to use `subset`
     23 subset(my.dataset, pH > 8)
     24 subset(my.dataset, season=="Summer", select=c(season,pH))
        # change values in the dataframe, e.g., sum 1 to the pH values of all summer rows
        my.dataset[my.dataset$season == "Summer", 'pH'] <-
           my.dataset[my.dataset$season == "Summer", 'pH'] + 1
     27
     28
        subset(my.dataset, season=="Summer", select=c(season,pH))
```

©Dr. Ibrahim кадwan – University of Canberra

DATA FRAMES (7)



• Some other operations on R:

```
#1- add new column to the data frame, must be the same number of rows
my.dataset$N03 \leftarrow c(234.5, 256.6, 654.1, 356.7, 776.4)
#2- check number of rows
nrow (my.dataset)
#3- check number of columns
ncol (my.dataset)
#4- check dimension
dim (my.dataset)
#5- edit the existing dataset
my.dataset <- edit(my.dataset)</pre>
#6- create new dataset and open it in the edit mode
new.data <- edit(data.frame())</pre>
#7- check names of the columns
names (my.dataset)
#8- change names of the columns
names(my.dataset) <- c("area", "season", "pH", "N03")</pre>
```

KEY TAKEAWAYS



- Matrices and Arrays are tabular forms that hold data from the same types
- Lists are like vectors, but for heterogeneous data types
- Data frames are the most commonly used data structure in R
- Data frames combine features of lists and matrices together

RECOMMENDED READING



- You are recommended to read section 3.5 from the following online book:
 - https://rafalab.github.io/dsbook/r-basics.html#data-frames
- Also, you may check some of the following links for some tutorial about data frames:
 - https://www.datamentor.io/r-programming/data-frame/
 - http://www.r-tutor.com/r-introduction/data-frame
 - https://www.tutorialspoint.com/r/r data frames.htm

ANNOUNCEMENTS



• Unit readiness test is available, due on Sunday 7th of March

• The solutions of the weekly lab exercises will be available every subsequent Monday.

• The census date is due this Friday, 5th of March