# Programming in Data science
## Factors, Missing values and Data frames

Asad Waqar Malik

asad.malik@um.edu.my

Department of Information Systems
Faculty of Computer Science and Information Technology
University of Malaya, Malaysia

Oct 03, 2019

UNIVERSITY
OF MALAYA

## Outline

## Factors

- Factors are used to represent categorical data.
- One can think of a factor as an integer vector where each integer has a label.
- Factors are important in statistical modeling.
- Using factors with labels is better than using integers because factors are self-describing.

> **Examples**
>
> 1. Having a variable that has values "Male" and "Female" is better than a variable that has values 1 and 2.
> 2. A data field such as marital status may contain only values from single, married, separated, divorced, or widowed.

## Factors

▶ Factor objects can be created with the factor(..) function.

```
1  > x <- factor(c("single", "married", "married", "↩
      single"));
2  > x
3  [1] single  married married  single
4  Levels: married single
```

▶ Levels may be predefined even if not used.

```
1  > x <- factor(c("single", "married", "married", "↩
      single"),
2  levels = c("single", "married", "divorced"));
3  > x
4  [1] single  married married  single
5  Levels: single married divorced
```

## Factors

### Sample Code

```
1  > x <- factor(c("single","married","married","↩
       single"))
2  > str(x)
3  Factor w/ 2 levels "married","single": 2 1 1 2
```

▶ Levels are stored in a character vector and the individual
  elements are actually stored as indices.

## Access components of a factor

▶ Accessing components of a factor is very similar to that of vectors.

### Sample Code

```
1  > x[3]                    # access 3rd element
2  >  x[c(2, 4)]             # access 2nd and 4th element
3  > x[-1]                   # access all but 1st element
4  > x[c(TRUE, FALSE, FALSE, TRUE)] # using logical ↩
        vector
```

## Modify Factors

▶ Components of a factor can be modified using simple assignments.

### Sample Code

```
1  > x
2  [1] single   married married  single
3  Levels: single married divorced
4  > x[2] <- "divorced"     # modify second element;
5  > x
6  [1] single    divorced married   single
7  Levels: single married divorced
```

## Modify Factors

### Invalid assignment

```
1  > x[3] <- "widowed"     # cannot assign values ←
         outside levels
2  Warning message:
3  In factor(*tmp*, 3, value = "widowed") :
4  invalid factor level, NA generated
5  > x
6  [1] single   divorced <NA>      single
7  Levels: single married divorced
```

► However, cannot choose values outside of predefined levels.

```
1  > k <- factor(c("Male"))
2  > k[1] <- "Female"
```

## Modify Factors

▶ A workaround to this is to add the value to the level first.

### Sample Code

```
1  >k <- factor(c("Male"))
2  > k
3  [1] Male
4  Levels: Male
5
6  > levels(k) <- c(levels(k), "Female")
7  > k
8  [1] Male
9  Levels: Male Female
```

▶ Often factors will be automatically created for you when you read a dataset in using a function like read.table(..).

## Missing Values

- Missing values are denoted by NA (not available) or NaN (not a number) for undefined mathematical operations.

  - NA used to represent a value that is not known, as a placeholder.
  - NA says no result was available or the result is missing. It can be used in a matrix to fill in a value of a vector.
  - NaN implies a result that cannot be calculated for whatever reason.

### Example – Sqrt of a negative number

```
1   sqrt(-1)
2   [1] NaN
3   Warning message:
4   In sqrt(-1) : NaNs produced
```

## Missing Values

- Following functions can be used to test objects.
  - is.na(..) is used to test objects if they are NA
  - is.nan(..) is used to test for NaN
  - NA values have a class also, so there are integer NA, character NA, etc

is

```
1  ## Create a vector with NAs in it
2  x <- c(1, 2, NA, 10, 3)
3  ## Return a logical vector indicating which ←
       elements are NA
4  is.na(x)
5  [1] FALSE FALSE TRUE FALSE FALSE
6  ## Return a logical vector indicating which ←
       elements are NaN
7  is.nan(x)
8  [1] FALSE FALSE FALSE FALSE FALSE
```

## Missing Values

▶ **Why is.na(..) can detect NA and NaN whereas,
is.nan(..) only detect NaN?**

### Sample Code

```
1  ## Now create a vector with both NA and NaN values
2  x <- c(1, 2, NaN, NA, 4)
3  is.na(x)
4  [1] FALSE FALSE TRUE TRUE FALSE
5  is.nan(x)
6  [1] FALSE FALSE TRUE FALSE FALSE
```

## Missing Values

▶ Excluding Missing Values from Analyses

### Sample Code

```
1  > x <- c(1,2,NA,3)
2  > mean(x)
3  [1] NA
4  > mean(x, na.rm=TRUE)
5  [1] 2
```

## Exclude Missing Values Using complete.cases(..)

► Excluding Missing Values from a vector using
  complete.cases(..)

### Sample Code

```
 1  > vec <- c(a=c(1,2,3,4,NA), b=c(3,4,5,6,NA))
 2  > vec
 3  a1 a2 a3 a4 a5 b1 b2 b3 b4 b5
 4   1  2  3  4 NA  3  4  5  6 NA
 5
 6  > k <- complete.cases(vec)
 7  > k
 8  [1]   TRUE   TRUE   TRUE   TRUE FALSE   TRUE   TRUE   TRUE↩
          TRUE FALSE
 9
10  > vec [k]
11  a1 a2 a3 a4 b1 b2 b3 b4
12   1  2  3  4  3  4  5  6
```

## Getting started with Data Frame

- ▶ Data frames are used to store tabular data in R.

- ▶ It is an Important type of object in R, used in a variety of statistical modeling applications.

- ▶ Unlike matrices, data frames can store different classes of objects in each column.

- ▶ Data frames have a special attribute called row.names which indicate information about each row of the data frame.

## Data Frame

- Data frames are usually *created by reading in a dataset* using the `read.table(..)` or `read.csv(..)`.
- Data frames can also be *created explicitly with the* `data.frame(..)` *function* or they can be *coerced from other types* of objects like lists.
- Data frames can be *converted to a matrix* by calling `data.matrix(..)`.

```
1   CustomerId,Name,Country,Salary
2   1,AA,MY,$30
3   1,AB,AU,$50
4   1,AC,USA,$60
5   1,AD,FR,$40
6   1,AE,IN,$20
7   1,AF,PK,$15
8
```

Figure 1: CSV format

## Creating a Data Frame

### Example

```
1  >x <- data.frame(foo = 1:4, bar = c(T, T, F, F))
2  >x
3  foo bar
4  1 TRUE
5  2 TRUE
6  3 FALSE
7  4 FALSE
```

```
1  >nrow(x)
2  [1] 4
3  >ncol(x)
4  [1] 2
```

Factors  Missing Values  **Data Frames**  Built-in data frame  Data Import  Names  Attributes  Summary  Exercise
0000000  00000  0000●0  000000  00

○
○
○○

## Creating Data Frames – Explicitly like a list

▶ Creating data frame

```
1 > measrs <- data.frame(gender = c("M", "M","F"), ht←
       = c(172, 186.5, 165), wt = c(91,99, 74))
2 > measrs
3     gender   ht        wt
4  1    M        172.0     91
5  2    M        186.5     99
6  3    F        165.0     74
```

▶ Entries in a data.frame are indexed like a matrix:

```
1  measrs[1, 2]
```

## Creating Data Frames from Vectors

▶ A data frame is used for storing data tables. It is a list of vectors of equal length.

### Sample Code

```
1  > n = c(2, 3, 5)
2  > s = c("aa", "bb", "cc")
3  > b = c(TRUE, FALSE, TRUE)
4  > df = data.frame(n, s, b)    # df is a data frame
```

## Data Frame Attributes

▶ Accessing data frame attributes using `names(..)` function.

```
1  > names(measrs)
2  [1] "gender" "ht" "wt"
```

▶ Assign names to row.

```
1  > rownames(measrs) <- c("S1", "S2", "S3")
```

▶ Access column inside the data frame.

```
1  > measrs$ht
2  [1] 172.0 186.5 165.0
```

## Data frame compenents as Vectors

▶ The components of a data frame can be extracted as a vector as in a list:

```
1  > height <- measrs$ht
2  > height
3  [1] 172.0 186.5 165.0
```

▶ Character vectors in a data frame are always stored as a factor. It's assumed that's what you should do.

```
1  > class(measrs$gend)
2  [1] "factor"
```

**Expanding Data Frames**

▶ Components can be added to a data frame in the natural way.

```
1 > measrs$age <- c(28, 55, 43)
2 > measrs
3      gender  ht       wt age
4 S1   M       172.0    91 28
5 S2   M       186.5    99 55
6 S3   F       165.0    74 43
```

## Expanding Data Frames

▶ If you expand the experiment to add data, use row binding to expand.

```
1  > m2 <- data.frame(gender = c("M", "F"), ht = c↩
       (170, 166), wt = c(68, 72), age = c(38, 22))
2  > rownames(m2) <- c("S4", "S5")
3  > measrs2 <- rbind(measrs, m2)
```

▶ If other data are kept on the same samples in another data frame it can be combined with the original using cbind.

## Built-in Data Frame

▶ Lets explore the built-in data frames in R. For example, here is a built-in data frame in R, called mtcars.

### mtcars data frame

```
1  > mtcars
2                        mpg cyl  disp   hp drat
3  Mazda RX4            21.0   6 160.0  110 3.90
4  Mazda RX4 Wag        21.0   6 160.0  110 3.90
5  Datsun 710           22.8   4 108.0   93 3.85
6  Hornet 4 Drive       21.4   6 258.0  110 3.08
7  Hornet Sportabout    18.7   8 360.0  175 3.15
8  Valiant              18.1   6 225.0  105 2.76
9  Duster 360           14.3   8 360.0  245 3.21
```

▶ Top line of the table called the header contains the column names.

▶ Horizontal line afterward denotes a data row.

## Retrieve data in a cell

- ▶ To retrieve data in a cell, use row and column coordinates in the single square bracket "[]" operator, separated by a comma.
- ▶ The coordinates begins with row position, followed by comma, and ends with the column position.

### Sample Code

```
1  > mtcars [1, 2]
2  [1] 6
```

- ▶ We can use the row and column names instead of the numeric coordinates.

```
1  > mtcars ["Mazda RX4", "cyl"]
2  [1] 6
```

Factors    Missing Values    Data Frames    **Built-in data frame**    Data Import    Names    Attributes    Summary    Exercise
ooooooo    ooooo             ooooo           ooooooo                   oo             

                                             o
                                             o
                                             oo

## nrow and ncol functions

▶ Lastly, the number of data rows in the data frame is given by
the nrow function.

**Sample code**

```
1  > nrow(mtcars)      # number of data rows
2  [1] 32
```

▶ And the number of columns of a data frame is given by the
ncol function.

```
1  > ncol(mtcars)      # number of columns
2  [1] 11
```

## Data Frame Column Vector

- We reference a data frame column with the double square bracket "[[]]" operator.
- To retrieve the ninth column vector of the built-in data set mtcars, we write mtcars[[9]].

### Sample Code

```
1  > mtcars[[9]]      # access 9th column
2    [1]  1 1 1 0 0 0 0 0 0 0 0 ...
```

- We can retrieve the same column vector by its name.

```
1  > mtcars[["am"]]
2    [1]  1 1 1 0 0 0 0 0 0 0 0 ...
```

Factors
0000000

Missing Values
00000

Data Frames
00000
0
0
00

Built-in data frame
0000●0

Data Import
00

Names

Attributes

Summary

Exercise

## Data Frame Column Slice

▶ We retrieve a data frame column slice with the single square bracket "[]" operator.

### Sample Code

```
1  > mtcars[1]
2                     mpg
3  Mazda RX4          21.0
4  Mazda RX4 Wag      21.0
5  Datsun 710         22.8
```

▶ We can retrieve the same column slice by its name.

```
1  > mtcars["mpg"]               # access single column
2  > mtcars[c("mpg", "hp")]      # access multiple ←
       column
```

## Data Frame Row Slice

▶ Retrieve rows from a data frame with the single square bracket operator, in additional to an index vector of row positions, append an extra comma character.

```
1  > mtcars[24,]
2              mpg cyl disp  hp drat   wt ...
3  Camaro Z28 13.3   8  350 245 3.73 3.84 ...
```

▶ To retrieve more than one rows, use a numeric index vector.

```
1  > mtcars[c(3, 24),]    # reterive multiple row
2  > mtcars["Camaro Z28",] #retrieve a row by its name←
       .
3  > mtcars[c("Datsun 710", "Camaro Z28"),]  # ←
       reterive multiple row
```

## Data Import from CSV File

- ▶ CSV : The data is in comma separated format.
- ▶ The first row represents column names.

```
1  Col1 , Col2 , Col3
2  100 , a1 , b1
3  200 , a2 , b2
4  300 , a3 , b3
```

- ▶ Code to read from "mydata.csv", using `read.csv(..)`.

```
1  > mydata = read.csv("mydata.csv")  # read csv file
2  > mydata
3     Col1 Col2 Col3
4  1  100    a1   b1
5  2  200    a2   b2
6  3  300    a3   b3
```

## Working directory

▶ The data files must be located in the R working directory,
  which can be found and set with the following functions.

```
1  > getwd()                   # get current working ←
       directory
2  > setwd("<new path>")     # set working directory
3  > setwd("C:/MyDoc")       # Note that the forward ←
       slash should be used as the path separator even←
        on Windows platform.
```

## Names

▶ R objects can have names, which is very useful for writing readable code and self-describing objects.

### Sample Code

```
1  > x <- 1:3
2  > names(x)
3  NULL
4  > names(x) <- c("New York", "Seattle", "Los Angeles↩
       ")
5  > x
6  New York Seattle Los Angeles
7  1              2        3
8  >names(x)
9  [1] "New York" "Seattle" "Los Angeles"
```

## Names

▶ Lists can also have names, which is often very useful.

### Example

```
1  > x <- list("Los Angeles" = 1, Boston = 2, London =←↩
      3)
2  > x
3  $ Los Angeles
4  [1] 1
5  $Boston
6  [1] 2
7  $London
8  [1] 3
9
10 >names(x)
11 [1] "Los Angeles" "Boston" "London"
```

## Names

- Matrices can have both column and row names.

### Example

```
1  > m <- matrix(1:4, nrow = 2, ncol = 2)
2  > dimnames(m) <- list(c("a", "b"), c("c", "d"))
3  > m
4     c d
5   a 1 3
6   b 2 4
```

## Names in Data frame

▶ for data frames, there is a separate function for setting the row names, the row.names() function.

▶ data frames do not have column names, they just have names (like lists).

▶ to set the column names of a data frame just use the names() function.

| Object | Set column names | Set row names |
|--------|------------------|---------------|
| data frame | names() | row.names() |
| matrix | colnames() | rownames() |

Table 1: Quick Summary.

## Attributes

- R objects can have attributes, which are like metadata for the object.
- metadata can be very useful in that they help to describe the object.
- For example, column names on a data frame help to tell us what data are contained in each of the columns.

### Example

- names, dimnames
- dimensions (e.g. matrices, arrays)
- class (e.g. integer, numeric)
- length
- other user-defined attributes/metadata

## Reading from Excel

```
1  > install.packages("readxl")
2  > library("readxl")
3  > my_data <- read_excel("C:/myfile.xlsx")


1  #Select file through browsing
2  >my_data <- read_excel(file.choose())


1  # Specify sheet by its name
2  my_data <- read_excel("my_file.xlsx", sheet = "data↩
       ")
```

Factors  Missing Values  Data Frames  Built-in data frame  Data Import  Names  Attributes  **Summary**  Exercise
0000000  00000       00000      000000        00

Summary

There are a variety of different builtin-data types in R. In this section we have reviewed the following.

▶ atomic classes: numeric, logical, character, integer, complex

▶ vectors, lists

▶ factors

▶ missing values

▶ data frames and matrices

**Practise questions**

$\Rightarrow$ Create the following data frame, afterwards invert Sex for all
individuals.

```
1           Age Height Weight Sex
2  Alex      25  177    57     F
3  Lily      20  165    69     F
4  Mark      21  145    70     M
5  Oliver     24  165    65     M
6  Martha     22  175    71     F
```

$\Rightarrow$ How many rows and columns does the new data frame have?
$\Rightarrow$ What class of data is in each column?

**Exc**

$\Rightarrow$ Create this data frame

```
1       Working
2  Alex      YES
3  Lily      NO
4  Mark      YES
5  Oliver    NO
6  Martha    YES
```

$\Rightarrow$ Add this data frame column-wise to the previous one.
$\Rightarrow$ How many rows and columns does the new data frame have?
$\Rightarrow$ What class of data is in each column?

**Exc**

$\Rightarrow$ Check what class of data is the (built-in data set) state.center and convert it to data frame. Hine, see as.data.frame?