**Instituto Superior de Engenharia de Lisboa**
**Mestrado em Engenharia Informática e de Computadores**
**Mestrado em Engenharia Informática e Multimédia**
**Mestrado em Matemática Aplicada para a Indústria**

**Big data mining (MDLE)**

**Laboratory Class #1 — R language and environment setup**
$2^{nd}$ **semester, 2022/2023 (March, 15)**

**Code and Report are due by March, 22**

---

1. **Data Resources and Software Tools.**
   For this laboratory class, you will need the following software:

   - Access to RServer (`http://datalys.dyn.fil.e.ipl.pt:8787`) or,

   - R, https://cran.r-project.org/, and

   - R Studio, https://posit.co/downloads/.

   As a preparation to the lab, you should:

   - download the Sparklyr and RStudio IDE cheat sheets from https://posit.co/resources/cheatsheets/.

---

2. **R Primer**.

   (a) Go to http://www.r-tutor.com/r-introduction and read the basic data types section.

   (b) Open R Studio, execute and explain each of the following statements, one by one (and check the values on the 'Global Environment' tab). Beware with the cut+paste:

   i. Scalars and Operators.

   ```
   var1 <- 3
   show(var1)
   var2 <- var1 * var1
   var3 <- var1 ** 2
   var4 <- var1 ^ 2
   var1 < var1
   var3 != var4
   var2 == var4
   var2 <- var2 - var2
   var5 <- var3 / var2
   var5 + 1
   ```

   ii. Manipulating vectors.

   ```
   a <- c(1,2,5.3,6,-2,4,3.14159265359)
   a
   b <- c("1","2","3")
   "2" %in% b
   "5" %in% b
   ```

```
 6  c <- c(TRUE,TRUE,FALSE,TRUE)
 7  a <- c
 8  a[0]
 9  a[1]
10  a[-1]
11  a[8]
12  a[c(1,3)]
13  a[c(3,1)]
14  a[a > 2]
```

### iii. Manipulating matrices.

```
 1  m <- matrix(1:6, nrow=3,ncol=2)
 2  show(m)
 3  n <- matrix(2:7, nrow=2,ncol=3)
 4  n
 5  m[,2]
 6  n[1,]
 7  m[2:3,1:2]
 8  n %*% m
 9  m %*% n
10  n %*% n
11  n^2
12  sqrt(n)
```

### iv. Manipulating data frames.

```
 1  d <- c(1,2,3,4)
 2  e <- c("Bob", "Alice", NA,"Joe")
 3  f <- c(TRUE,TRUE,FALSE,TRUE)
 4  my.data <- data.frame(d,e,f,stringsAsFactors=FALSE)
 5  show(my.data)
 6  names(my.data) <- c("ID","Name","Passed")
 7  View(my.data)
 8  my.data$Name
 9  my.data <- cbind(my.data, Failed=!f)
10  View(my.data)
11  my.data <- rbind(my.data, c(5,"Carol",FALSE,TRUE))
12  View(my.data)
13  ?data.frame
```

### v. Manipulating factors.

```
 1  colour <- c(rep("red",20), rep("blue", 30))
 2  colour <- factor(colour)
 3  summary(colour)
 4  dimensions <- c("large", "medium", "small")
 5  show(dimensions)
 6  dimensions <- ordered(dimensions)
 7  show(dimensions)
```

### vi. Some useful functions.

```
 1  length(colour)
 2  length(a)
 3  class(colour)
 4  class(dimensions)
 5  class(a)
 6  class(my.data)
 7  nrow(my.data)
 8  ncol(my.data)
 9  str(my.data)
10  sessionInfo()
11  ls()
12  rm(a)
13  glimpse(my.data)
```

vii. Packages.

```
1   library()
2   search()
3   library("MASS")
4   search()
5   .libPaths()
6   install.packages("e1071")
7   install.packages("funModeling")
8
9   x<-c("MASS","dplyr", "e1071")
10  lapply(x, require, character.only = TRUE)
11
12  require(funModeling)
```

viii. Basic descriptive statistics and more.

```
1   iris
2   summary(iris)
3   fivenum(iris$Sepal.Length)
4
5   status(iris)
```

(c) Open the file "**ControlFlow.R**", handed with this guide, and:

   i. Set a breakpoint at line 6 and run the code delimited by region *EX.1*, step by step.
   ii. Explain the purpose of the three examples, named *EX.1*, *EX.2*, and *EX.3*.

(d) Implement and show a code that gets the first two columns of *my.data*, using the range operator indexation.

(e) Implement and show a code that gets the first two columns of *my.data*, using vector indexation.

(f) Install the package `Hmisc`, and use the function `describe` on *my.data* variable.

3. **Setup Spark**.

(a) Install the package `sparklyr` using `install.packages("sparklyr")`. It may take some minutes, so be patient.

(b) Install *Spark* using `library(sparklyr) spark_install(version = '3.3.2', hadoop_version = '3')`

4. **Spark primer**.
Describe the commands used (if not listed), the problems (if any), and the results of the following points, line by line:

(a) Check, programmatically, if `sparklyr` package is loaded. Next, connect to *Spark* using
`ss <- spark_connect('local', version = '3.3.2', hadoop_version = '3', config = list())`.

(b) View the content of variable `ss`.

(c) Use function `copy_to` from package `dplyr` to copy the *iris* data set to *Spark*.

```
1   library(dplyr)
2   df <- copy_to(ss, iris)
3   show(df)
```

(d) Show some sample data from the loaded data set.

```
1   head( select( df, Petal_Width, Species))
2   head( filter( df, Petal_Width > 0.3))
3   df %>% head
```

(e) Using SQL

```
1   library(DBI)
2   df_sql <- dbGetQuery(ss, "SELECT * FROM iris WHERE Petal_Width > 0.3 LIMIT 5")
3   show(df_sql)
```

(f) Get data from *Spark* nodes.

```
1  local_df <- collect(df)
2  show(local_df)
3  show(df)
```

(g) Disconnect spark using `spark_disconnect(ss)`. Confirm, programmatically, that *Spark* is closed.

(h) Indicate, as an example, a <u>supervised learning</u> algorithm implemented on *Spark*, and available through `sparklyr`.

(i) Indicate, as an example, a <u>unsupervised learning</u> algorithm implemented on *Spark*, and available through `sparklyr`.

(j) Tell how *Spark* can support the *Big data process pipeline*. Try to start at https://spark.rstudio.com