

MAT4681 - Statistique pour les sciences

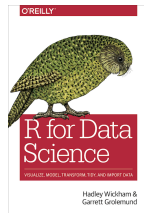
Arthur Charpentier

02 - Introduction langage R

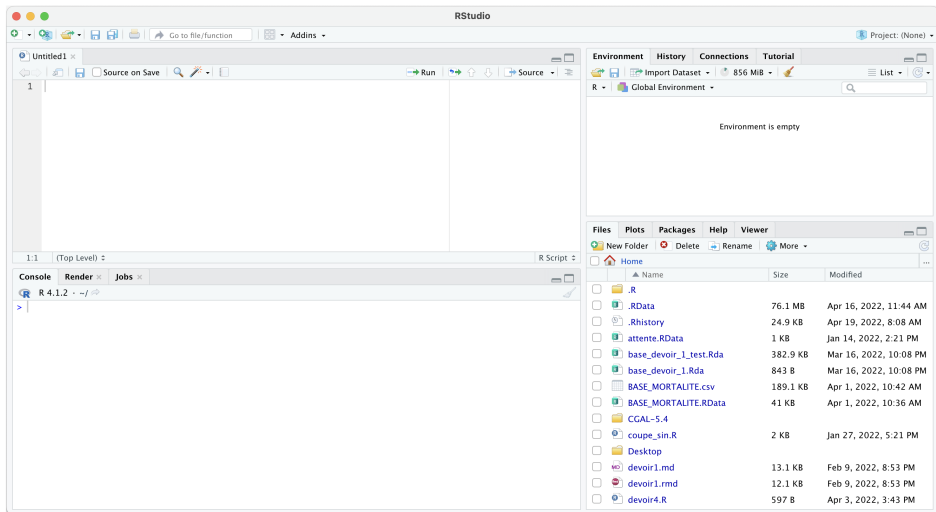
été 2022

Références

- ▶ Introduction à la programmation en R de Vincent Goulet
- ▶ Notes de cours de R d'Ewen Gallic
- ▶ R pour les débutants d'Emmanuel Paradis
- ▶ R pour la statistique et la science des données de Pierre-André Cornillon, Arnaud Guyader, François Husson, Nicolas Jégou, Julie Josse, Nicolas Klutchnikoff, Erwan Le Pennec, Eric Matzner-Lober, Laurent Rouvière, Benoît Thieurmél
- ▶ R for data science d'Hadley Wickham et Garrett Golemund
- ▶ YaRrr! The Pirate's Guide to R de Nathaniel Phillips



R Studio



Nombres et vecteurs I

```
1 > set.seed(1)
2 > U <- runif(20)
3 > U[1:4]
4 [1] 0.2655087 0.3721239 0.5728534 0.9082078
5 > options(digits = 3)
6 > U[1:4]
7 [1] 0.266 0.372 0.573 0.908
8 > options(digits = 22)
9 > U[1:4]
10 [1] 0.2655086631420999765396 0.3721238996367901563644
11 [3] 0.5728533633518964052200 0.9082077899947762489319
12 > x <- exp(1)
13 > y <- x
14 > x <- 2
15 > y
16 [1] 2.72
17 > class(x)
18 [1] "numeric"
```

Nombres et vecteurs II

```
1 > (3/10-1/10)
2 [1] 0.2
3 > (3/10-1/10)==(7/10-5/10)
4 [1] FALSE
5 > (3/10-1/10)-(7/10-5/10)
6 [1] 2.78e-17
7 > all.equal((3/10-1/10),(7/10-5/10))
8 [1] TRUE
9 > 0/0
10 [1] NaN
11 > 1/0
12 [1] Inf
```

Nombres et vecteurs III

```
1 > x <- rnorm(8)
2 > names(x) <- letters[1:6]
3 > x
4           a           b           c           d           e           f
5  1.5118  0.3898 -0.6212 -2.2147  1.1249 -0.0449
6 > x[2:4]
7           b           c           d
8  0.390 -0.621 -2.215
9 > x[c("b", "c", "d")]
10          b           c           d
11  0.390 -0.621 -2.215
```

Matrices I

```
1 > M <- 1:24
2 > dim(M) <- c(6,4)
3 > M
4           [,1] [,2] [,3] [,4]
5 [1,]         1     7    13    19
6 [2,]         2     8    14    20
7 [3,]         3     9    15    21
8 [4,]         4    10    16    22
9 [5,]         5    11    17    23
10 [6,]         6    12    18    24
11 > str(M)
12 int [1:6, 1:4] 1 2 3 4 5 6 7 8 9 10 ...
13 > colnames(M)=letters[1:4]
14 > rownames(M)=LETTERS[10:15]
15 > M
16      a  b  c  d
17 J 1  7 13 19
18 K 2  8 14 20
19 L 3  9 15 21
20 M 4 10 16 22
```

Matrices II

```
21  N 5 11 17 23
22  0 6 12 18 24
23  > M["K", ]
24      a  b  c  d
25      2  8 14 20
26  > M[c("K", "N"), ]
27      a  b  c  d
28  K 2  8 14 20
29  N 5 11 17 23
30  > M[c(2,5), ]
31      a  b  c  d
32  K 2  8 14 20
33  N 5 11 17 23
```


Facteur - variable catégorielle I

```
1 > x <- factor(c("b","a","b"))
2 > levels(x)
3 [1] "a" "b"
4 > x[3] <- "c"
5 Warning in '[<-.factor'('(*tmp*', 3, value = "c"):
   invalid factor level, NA generated
6 > x
7 [1] b      a      <NA>
8 Levels: a b
9 > factor(x,levels=c("b","a"))
10 [1] b      a      <NA>
11 Levels: b a
12 > x[1]
13 [1] b
14 Levels: a b
```

Facteur - variable catégorielle I

```
1
2
3 U <- runif(20)
4 cut(U,breaks=2)
5
6 ##      [1] (0.145,0.544] (0.145,0.544] (0.145,0.544]
7         (0.544,0.944] (0.544,0.944]
8 ##      [6] (0.145,0.544] (0.145,0.544] (0.145,0.544]
9         (0.544,0.944] (0.145,0.544]
10 ##     [11] (0.145,0.544] (0.145,0.544] (0.145,0.544]
11         (0.544,0.944] (0.544,0.944]
12 ##     [16] (0.145,0.544] (0.145,0.544] (0.544,0.944]
13         (0.145,0.544] (0.544,0.944]
14 ## Levels: (0.145,0.544] (0.544,0.944]
15
16 cut(U,breaks=2,labels=c("small","large"))
17
18 ##      [1] small small small large large small small
19         small large small small
```

Facteur - variable catégorielle II

```
15 ## [12] small small large large small small large
    small large
16 ## Levels: small large
17
18 cut(U,breaks=c(0,.3,.8,1),labels=c("small","medium","
    large"))
19
20 ## [1] medium small medium large medium medium
    medium small large medium
21 ## [11] medium small medium medium large medium
    medium medium medium large
22 ## Levels: small medium large
23
24 table(cut(U,breaks=c(0,.3,.8,1),labels=c("small","
    medium","large")))
25
26 ##
27 ## small medium large
28 ## 3 13 4
```

Facteur - variable catégorielle I

```
1 "Be carefull of 'quotes '"
2 'Be carefull of "quotes"'
3
4 cities <- c("New York, NY", "Los Angeles, CA", "Boston
5           , MA")
6 substr(cities, nchar(cities)-1, nchar(cities))
7 ## [1] "NY" "CA" "MA"
8
9 unlist(strsplit(cities, ", "))[seq(2,6,by=2)]
10
11 ## [1] "NY" "CA" "MA"
12
13 some.dates <- as.Date(c("16/10/12", "19/11/12"), format
14                       = "%d/%m/%y")
15 diff(some.dates)
16 ## Time difference of 34 days
```

Facteur - variable catégorielle I

```
1 x <- list(1:5,c(1,2,3,4,5),a="test",
2 b=c(TRUE,FALSE),rpois(5,8))
3 x
4
5 ## [[1]]
6 ## [1] 1 2 3 4 5
7 ##
8 ## [[2]]
9 ## [1] 1 2 3 4 5
10 ##
11 ## $a
12 ## [1] "test"
13 ##
14 ## $b
15 ## [1] TRUE FALSE
16 ##
17 ## [[5]]
18 ## [1] 6 13 9 4 7
19
20 f <- function(x) { return(x*(1-x)) }
```

Facteur - variable catégorielle II

```
21 optimize(f, interval=c(0, 1), maximum=TRUE)
22
23 ## $maximum
24 ## [1] 0.5
25 ##
26 ## $objective
27 ## [1] 0.25
28
29 set.seed(1)
30 u <- runif(1)
31 if(u>.5) {"greater than 50%"} else {"smaller than
32       50%"}
33 ## [1] "smaller than 50%"
34
35 ifelse(u>.5,("greater than 50%"),("smaller than 50%"))
36
37 ## [1] "smaller than 50%"
38
39 u
```

Facteur - variable catégorielle III

```
40
41 ## [1] 0.266
42
43 df <- data.frame(x=1:3,y=letters[1:3])
44 str(df)
45
46 ## 'data.frame':    3 obs. of  2 variables:
47 ## $ x: int  1 2 3
48 ## $ y: Factor w/ 3 levels "a","b","c": 1 2 3
49
50 typeof(df)
51
52 ## [1] "list"
53
54 class(df)
55
56 ## [1] "data.frame"
57
58 df$z<-5:3
59 df
```

Facteur - variable catégorielle IV

```
60
61 ##      x y z
62 ## 1 1 a 5
63 ## 2 2 b 4
64 ## 3 3 c 3
65
66 set.seed(1)
67 df[sample(nrow(df)),]
68
69 ##      x y z
70 ## 1 1 a 5
71 ## 3 3 c 3
72 ## 2 2 b 4
```


Facteur - variable catégorielle I

```
1 download.file("http://freakonometrics.free.fr/  
  superheroes.RData","superheroes.RData")  
2 load("superheroes.RData")  
3 superheroes  
4  
5 ##      name alignment gender      publisher  
6 ## 1  Magneto      bad   male      Marvel  
7 ## 2   Storm     good female      Marvel  
8 ## 3 Mystique     bad female      Marvel  
9 ## 4  Batman     good   male        DC  
10 ## 5   Joker     bad   male        DC  
11 ## 6 Catwoman    bad female        DC  
12 ## 7  Hellboy     good   male Dark Horse Comics  
13  
14 publishers  
15  
16 ##      publisher yr_founded  
17 ## 1          DC      1934  
18 ## 2      Marvel      1939  
19
```

Facteur - variable catégorielle II

```
20 library(dplyr, verbose=FALSE)
21
22 inner_join(publishers, superheroes)
23
24 ## Joining by: "publisher"
25
26 ## Warning in inner_join_impl(x, y, by$x, by$y):
    joining factors with
27 ## different levels, coercing to character vector
28
29 ##      publisher yr_founded      name alignment gender
30 ## 1           DC      1934   Batman      good   male
31 ## 2           DC      1934    Joker      bad    male
32 ## 3           DC      1934 Catwoman      bad female
33 ## 4      Marvel      1939   Magneto      bad   male
34 ## 5      Marvel      1939    Storm      good female
35 ## 6      Marvel      1939 Mystique      bad female
36
37 merge(superheroes, publishers, all = TRUE)
38
```

Facteur - variable catégorielle III

```
39 ##           publisher      name alignment gender
    yr_founded
40 ## 1 Dark Horse Comics  Hellboy      good   male
    NA
41 ## 2                DC    Batman      good   male
    1934
42 ## 3                DC     Joker      bad    male
    1934
43 ## 4                DC Catwoman      bad female
    1934
44 ## 5                Marvel  Magneto      bad   male
    1939
45 ## 6                Marvel   Storm      good female
    1939
46 ## 7                Marvel Mystique      bad female
    1939
47 ## 8                Image    <NA>      <NA>   <NA>
    1992
48
49 left_join(superheroes , publishers)
```

Facteur - variable catégorielle IV

```
50
51 ## Joining by: "publisher"
52
53 ## Warning in left_join_impl(x, y, by$x, by$y):
54   joining factors with different
55   levels, coercing to character vector
56 ##           name alignment gender           publisher
57   yr_founded
58 ## 1  Magneto           bad   male           Marvel
59   1939
60 ## 2   Storm           good female           Marvel
61   1939
62 ## 3 Mystique           bad female           Marvel
63   1939
64 ## 4   Batman           good   male              DC
65   1934
66 ## 5    Joker           bad   male              DC
67   1934
```

Facteur - variable catégorielle V

```
62 ## 6 Catwoman      bad female      DC
    1934
63 ## 7 Hellboy      good   male Dark Horse Comics
    NA
64
65 left_join(publishers, superheroes)
66
67 ## Joining by: "publisher"
68
69 ## Warning in left_join_impl(x, y, by$x, by$y):
    joining factors with different
70 ## levels, coercing to character vector
71
72 ##      publisher yr_founded      name alignment gender
73 ## 1          DC      1934    Batman      good   male
74 ## 2          DC      1934     Joker      bad    male
75 ## 3          DC      1934 Catwoman      bad female
76 ## 4      Marvel      1939    Magneto      bad   male
77 ## 5      Marvel      1939     Storm      good female
78 ## 6      Marvel      1939  Mystique      bad female
```

Facteur - variable catégorielle VI

79	##	7	Image	1992	<NA>	<NA>	<NA>
----	----	---	-------	------	------	------	------

Facteur - variable catégorielle I

```
1 download.file("http://freakonometrics.free.fr/
  gapminderDataFiveYear.txt", "gapminderDataFiveYear.
  txt")
2 gdf <- read.delim("gapminderDataFiveYear.txt")
3 head(gdf, 4)
4
5 ##           country year      pop continent lifeExp
6 ## 1 Afghanistan 1952  8425333      Asia      28.8
7   779
8 ## 2 Afghanistan 1957  9240934      Asia      30.3
9   821
10 ## 3 Afghanistan 1962 10267083      Asia      32.0
11   853
12 ## 4 Afghanistan 1967 11537966      Asia      34.0
13   836
14
15 str(gdf)
16
17 ## 'data.frame':    1704 obs. of  6 variables:
```

Facteur - variable catégorielle II

```
14 ## $ country : Factor w/ 142 levels "Afghanistan
    ",...: 1 1 1 1 1 1 1 1 1 1 ...
15 ## $ year : int 1952 1957 1962 1967 1972 1977
    1982 1987 1992 1997 ...
16 ## $ pop : num 8425333 9240934 10267083
    11537966 13079460 ...
17 ## $ continent: Factor w/ 5 levels "Africa","Americas
    ",...: 3 3 3 3 3 3 3 3 3 3 ...
18 ## $ lifeExp : num 28.8 30.3 32 34 36.1 ...
19 ## $ gdpPercap: num 779 821 853 836 740 ...
20
21 subset(gdf, lifeExp < 30)
22
23 ## country year pop continent lifeExp
    gdpPercap
24 ## 1 Afghanistan 1952 8425333 Asia 28.8
    779
25 ## 1293 Rwanda 1992 7290203 Africa 23.6
    737
26
```


Facteur - variable catégorielle III

```
27 gdf[gdf$lifeExp < 30,]
28
29 ##           country year      pop continent lifeExp
      gdpPercap
30 ## 1      Afghanistan 1952 8425333      Asia      28.8
      779
31 ## 1293      Rwanda 1992 7290203      Africa      23.6
      737
32
33 gdf[gdf$country == "Italy", c("year", "lifeExp")]
34
35 ##      year lifeExp
36 ## 769 1952      65.9
37 ## 770 1957      67.8
38 ## 771 1962      69.2
39 ## 772 1967      71.1
40 ## 773 1972      72.2
41 ## 774 1977      73.5
42 ## 775 1982      75.0
43 ## 776 1987      76.4
```

Facteur - variable catégorielle IV

```
44 ## 777 1992      77.4
45 ## 778 1997      78.8
46 ## 779 2002      80.2
47 ## 780 2007      80.5
48
49 small_df <- df[df$country %in% c("France","Italy","
    Spain"), c("country","year", "lifeExp")]
50
51 aggregate(small_df$lifeExp,FUN=max,by=list(
    small_df$country))
52
53 ##      Group.1      x
54 ## 1   France 80.7
55 ## 2    Italy 80.5
56 ## 3    Spain 80.9
```

Facteur - variable catégorielle I

```
1 library(gamair)
2 data(chicago)
3 head(chicago)
4
5 ##      death pm10median pm25median o3median so2median
6 ## 1      130      -7.434          NA      -19.6        1.928
7 ##      -2556 31.5
8 ## 2      150          NA          NA      -19.0       -0.986
9 ##      -2556 33.0
10 ## 3      101     -0.827          NA     -20.2       -1.891
11 ##      -2554 33.0
12 ## 4      135      5.566          NA     -19.7        6.139
13 ##      -2554 29.0
14 ## 5      126          NA          NA     -19.2        2.278
15 ##      -2552 32.0
16 ## 6      130      6.566          NA     -17.6        9.859
17 ##      -2552 40.0
18
19 base=data.frame(death=chicago$death,
```

Facteur - variable catégorielle II

```
14         temp_F=chicago$tmpd ,
15         o3=chicago$o3median ,
16         date=seq(as.Date("1987-01-01") ,
17                 as.Date("2000-12-31"),by=1))
18 base$temp_C <- (base$temp_F-32)/1.8
19 base$year <- substring(base$date,1,4)
20
21 date2season <- function(date){
22     m <- as.numeric(format(as.Date(date, format = "%d/%m
23                             /%Y"), "%m"))
24     d <- as.numeric(format(as.Date(date, format = "%d/%m
25                             /%Y"), "%d"))
26     s <- NA
27     if(m %in% c(1,2) | ((m==12)&(d>=21)) | ((m==3)&(d
28         <21))) s <- "winter"
29     if(m %in% c(4,5) | ((m==3)&(d>=21)) | ((m==6)&(d
30         <21))) s <- "spring"
31     if(m %in% c(7,8) | ((m==6)&(d>=21)) | ((m==9)&(d
32         <21))) s <- "summer"
```

Facteur - variable catégorielle III

```
28   if(m %in% c(10,11) | ((m==9)&(d>=21)) | ((m==12)&(d
    <21))) s <- "autumn"
29   return(s)}
30 base$season <- sapply(base$date, date2season)
31 head(base)
32
33 ##      death temp_F      o3      date temp_C year season
34 ## 1      130    31.5  -19.6 1987-01-01  -0.278 1987 winter
35 ## 2      150    33.0  -19.0 1987-01-02   0.556 1987 winter
36 ## 3      101    33.0  -20.2 1987-01-03   0.556 1987 winter
37 ## 4      135    29.0  -19.7 1987-01-04  -1.667 1987 winter
38 ## 5      126    32.0  -19.2 1987-01-05   0.000 1987 winter
39 ## 6      130    40.0  -17.6 1987-01-06   4.444 1987 winter
40
41 plot(base[,c("date", "temp_C")])
```