

# Courte introduction à R

Florent Benaych-Georges

Université Paris Descartes

Janvier 2015

*Séries temporelles pratiques*

2 étapes :

- 1 Télécharger et installer **R** depuis  
`http://cran.r-project.org`
- 2 Utiliser **R studio**, un software open-source pour R :  
`http://www.rstudio.com`

# Principaux types

Essentiellement 3 types d'objets :

char, logical, numeric

```
1 x=5
2 class(x)
3
4 y=(x<6)
5 y
6 class(y)
7
8 (z="Hello_world") #parentheses: result display
9 class(z)
```

(file: Types.R)

Essentiellement 5 types de conteneurs :

- vector
- list
- factor
- data.frame
- ts

Toutes les entrées ont le **même type**

```
1 (v=c(2,5,7,8)) #c=concatenates entries or vectors
2                 #in a single vector
3
4 (u=c(v,9,15:18)); u[u<10]
5
6 (w=letters[1:7]); class(w) #vectors: not only with
7                             #numerics as entries
8 w[3] #third coordinate
9 w[length(v4)] #last entry
10 w[c(1,4,5)] #1st, 4th and 5th entries
11 w[-c(1,4,5)] #all but 1st, 4th and 5th entries
12
13 (x=seq(from=1, to=2, by=.2))
14
15 (y=rep(x,2))
```

(file: Vectors.R)

Conteneur similaire en 2 dimensions : **matrix**

Les coordonnées peuvent avoir des **types différents** et des **noms**

```
1 L=list (name="Alban", age=1, country="France",  
2         size_in_cm=83)  
3 str(L)  
4 L$name  
5 L[[2]]
```

(file: Lists .R)

Adaptés pour les données **qualitatives** (i.e. non quantitatives)

```
1 C=c("blue","blue","red","red","blue","green",  
2     "yellow","green","red") #vector of characters  
3 (my_factor1=factor(C))  
4 str(my_factor1)  
5 table(my_factor1)  
6  
7 L2=c("blue","red","green","yellow","orange")  
8 (my_factor2=factor(C,levels=L2))  
9 table(my_factor2)
```

(file: Factors.R)

**data.frame** : liste de vecteurs de **même taille** mais éventuellement de classes différentes

```
1 name=c("Clo", "Tony", "Alban", "Tom", "Lou")
2 city=c("Paris", "Paris", "Lyon", "NYC", "Shanghai")
3 size=c(170,181,83,185,179)
4
5 (d = data.frame(name, city, size))
6 str(d)
7 dim(d); nrow(d); ncol(d)
8 d$name
9 d[2,3] #2nd entry of 3rd column
10 d$year = c(1983,1980,2013,1988,1986) #add column
11 d[2,] #2nd row
12 d[,c(3,1)] #3rd and 1st columns
13 d[,-c(2,4)] #all but 2nd and 4th columns
14 d[d$city=="Paris",] #extracts a subset
15 colnames(d)[4]="birth_year" #changes column name
16 summary(d)
```



## ts (Time Series) (1)

Objet de type **ts** : un **vecteur** de valeurs, un vecteur **start** (année, moment), un vecteur **end** (année, moment) et une **fréquence**

```
1 # save a numeric vector containing 72 monthly observations
2 myvector = .1*(1:72)+sin((1:72)*pi/12) +
3   rnorm(72,mean=0,sd=.3)
4
5 # from Jan 2009 to Dec 2014 as a time series object
6 myts = ts(myvector, start=c(2009, 1), end=c(2014, 12),
7   frequency=12)
8 class(myts)
9 plot(myts, col="red")
10
11 # subset the time series (June 2014 to December 2014)
12 myts2 <- window(myts, start=c(2010, 6), end=c(2014, 12))
13 class(myts2)
14 plot(myts2, col="blue")
```

(file: Example\_ts\_1.R)

## ts (Time Series) (2)

Objet de type **ts** : un **vecteur** de valeurs, un vecteur **start** (année, moment), un vecteur **end** (année, moment) et une **fréquence**

```
1 # save a numeric vector containing 72 monthly observations
2 myvector = .1*(1:72)+sin((1:72)*pi/12) +
3   rnorm(72, mean=0, sd=.3)
4
5 # from Jan 2009 to Dec 2014 as a time series object
6 myts = ts(myvector, start=c(2009, 1), end=c(2014, 12),
7   frequency=12)
8
9 # Seasonal decomposition
10 fit = stl(myts, s.window="period")
11 plot(fit, col="purple")
12
13 # additional plots
14 monthplot(myts)
15 library(forecast)
16 seasonplot(myts)
```

## ts (Time Series, Corrélation) (2)

Objet de type **ts** : un **vecteur** de valeurs, un vecteur **start** (année, moment), un vecteur **end** (année, moment) et une **fréquence**

```
1 data(USAccDeaths)
2
3 USAccDeaths
4
5 bruit.blanc=ts(rnorm(100), frequency = 1, start = c(1), end=c(100))
6
7 plot( bruit.blanc )
8
9 acf( bruit.blanc , type="correlation" ) #ou type="covariance" o
```

(file: Example\_ts\_3.R)

## Autres types, conversions

Beaucoup d'autres types en R (date, matrix, etc...). Conversions possibles grâce à `as.mon_type(mon_object)`

```
1 (d1=as.Date("1999-09-17"))
2 class(d1)
3
4 (d2=as.Date("19980901", format="%Y%m%d"))
5 (my_dates=seq(d1, d2, by="-1month"))
6
7 (x=as.numeric(d1))
8
9 (my_ch=as.character(d2))
10
11 (my_factor1=factor(c("blue", "red", "green", "red")))
12 (v=as.character(my_factor1))
```

(file: Dates\_and\_conversions.R)

## Obtenir et modifier le répertoire de travail

```
1  getwd() #get working directory
2  setwd("~/Desktop/") #set working directory
3  #setwd works with linux absolute or relative pathnames:
4
5  # ~/ returns you to your login directory
6
7  # /home/ takes you to the home directory
8  #where user login directories are usually stored
9
10 # ../ moves you up one directory
11
12 # ./my_sub_dir moves you down to the subdirectory
13    #named "my_sub_dir"
14
15 getwd() #get working directory
16 dir() #list of directory files
17 ls() #list of R objects
```

(file: Linux\_pathnames.R)

Une étape clé dans l'analyse de données Souvent : **compliqué !!**

```
1 (x=scan("my_data_file.dat"))
2 class(x) #scan gives a vector back
3
4 (y=read.table("my_table.dat", header=1))
5 class(y) #read.table gives a data.frame back
6
7 ta=read.csv(file="exo1.csv", header=1, skip=6)
8 #here, read.table also works
9 #"skip" allows to skip a certain number of rows
10
11 h=read.table("http://www.cmapx.polytechnique.fr/~benaych/201
12             header=T, fill=T)#fill: blank fields
13                               #added when rows
14                               #have missing values
15 names(h)
16 str(h)
17 summary(h)
```

Beaucoup plus facile que l'importation

```
1 write(runif(5), file="my_writing_file.dat")
2
3 y=read.table("my_table.dat", header=1)
4 write.table(y, file="my_writing_file_for_table.dat")
5
6 write.csv(y, file="my_writing_file_for_table.csv")
7 (z=read.table("my_writing_file_for_table.dat",
8             header=1))
9
10 ta=read.csv(file="exo1.csv", header=1, skip=6)
11 write.csv(ta, file = "../my_csv_file.csv")
```

(file: Export\_data.R)

## Import/export d'objects R

```
1 myvector = .1*(1:72)+sin((1:72)*pi/12) +
2   rnorm(72,mean=0,sd=.3)
3
4 # from Jan 2009 to Dec 2014 as a time series object
5 myts = ts(myvector, start=c(2009, 1), end=c(2014, 12),
6           frequency=12)
7
8 save(myts, file="my_rda_file.rda")
9 rm(list=ls())
10 print(myts)
11 load("my_rda_file.rda")
12 print(myts)
13 plot(myts)
```

(file: Import\_export\_R\_objects.R)



# Fonctions

```
1 f = function(x) {x*x}
2 f(5)
3 f(1:7)
4
5 PF = function(n=1,proba=.5) {
6     as.integer(runif(n)<proba)}
7 PF(10)
8 PF(10,.1)
9
10 PF2 = function(n=1,proba=.5) {
11     ech=as.integer(runif(n)<proba)
12     freq=mean(ech)
13     return(list(sample=ech,frequency=freq))}
14 PF2(10,.6)
```

(file: Functions.R)

# Structures de contrôle

```
1 x=runif(1); y=c(); z=c(); #x: random in [0,1]
2 for(i in 1:10) y[i]=sin(pi*i/10)
3 y==sin(pi*(1:10)/10)
4
5 for(i in 1:10) {
6   if (y[i]<.5) {z[i]=y[i]} else {z[i]=y[i]+1}}
7 z==((y<.5)*y+(y>=.5)*(y+1))
8
9 while(x<.5) {x=runif(1)}
10 x
11
12 #way quicker : matrix approach:
13 v=runif(1E5) #10^5 random numbers in [0,1]
14 t=Sys.time()
15 for (i in 1:length(v)){v[i]=v[i]+5}
16 t1=Sys.time()-t
17 t=Sys.time()
18 v=v+5
19 t2=Sys.time()-t; as.numeric(t1)/as.numeric(t2)
```

## Graphiques : plot, points, lines

```
1 graphics.off() #erases previous graphics
2 x=seq(.4,4,.1); y=log(x^2+1/x^2)
3 plot(x,y,type="p") #p for points
4 plot(x,y,type="l") #l for line
5 plot(x,y,type="s") #l for step
6 plot(x,y,type="p",pch=2) #pch : type of points
7 ##ylim=c(ay,by), xlim=c(ax,bx) : bounds of the axes
8 ##lty : type of line, col : color, main : title
9 x=rnorm(20); y=rexp(20)
10 plot(x,y,col="red",ylim=c(-0.5,5),xlim=c(-3,2))
11 points(x+.5,y-1,pch=2,col="blue") #add points
12 legend("topleft",c("plot_1","plot_2"),col=c("red","blue"),
13        pch=c(1,2))
14 lines(x,y,lty=2,col="green") #add line
15 abline(h=3,col="yellow",lty=1)#add horizontal line
16 abline(v=-1.3,col="aquamarine2")#add vertical line
17 abline(a=1,b=1,col="coral3")#add oblique line
18 text(-2,2,"commentary"); title("Example┘curves")
```

## Les nombreuses "méthodes" de plot

```
1 #Images produced by plot(x) depend on the class of x:
2 methods(plot)
3
4 #Examples:
5 class(sin); plot(sin, xlim=c(-pi, pi), xlab="x", ylab="sin(x)")
6
7 x=rpois(1000, 1); y=table(x); class(y); plot(y)
8
9 # save a numeric vector containing 72 monthly observations
10 myvector = .1*(1:72)+sin((1:72)*pi/12) +
11   rnorm(72, mean=0, sd=.3)
12
13 # from Jan 2009 to Dec 2014 as a time series object
14 myts = ts(myvector, start=c(2009, 1), end=c(2014, 12),
15   frequency=12)
16 class(myts); plot(myts, col="red")
17
18 r=rnorm(1E3, 1); w=density(r); class(w)
19 plot(w)
```

# Probabilités, variables éleatoires, distributions

```
1 #Discrete laws: binom, hyper, pois, geom, sample, ...
2 #Continuous: norm, unif, t, chisq, ...
3
4 # General syntax, for "my_law" a probability distribution:
5 # rmy_law(n,...) : sample
6 # dmy_law(x,...) : density of my_law at x
7 # pmy_law(x,...) : P( X ≤ x) for X distributed as my_law
8 # qmy_law(a,...) : quantile with level a
9
10 x=seq ( -5 ,5 ,.01)
11 u=seq ( 0 , 1 , .01)
12 par(mfrow=c(2,2))#defines an array of plots
13 plot(rnorm(500,0,1),pch=20,col="red") #sample
14 plot(x,dnorm(x,0,1),type="l",col="blue") #density
15 plot(x,pnorm(x,0,1),type="l",col="purple") #P( X ≤ x)
16 plot(u,qnorm(u,0,1),type="l",col="green") #quantiles
17 graphics.off() #erases graphics
```

(file: Probability \_ Random \_ variables \_ Distributions.R)

## Description graphique des données : hist, boxplot, pie, barplot...

```
1 data(iris); str(iris); summary(iris); head(iris)
2 x=iris$Sepal.Length
3 mean(x); sd(x); var(x)
4 quantile(x,c(.25 ,.5 ,.75))
5 boxplot(x)
6 hist(x)
7 hist(x,breaks = 12)
8
9 records=sample(c("clear", "cloudy", "rainy"),365,
10               replace=T,prob=c(5,4,1))
11 pie(table(records))
12 barplot(table(records))
```

(file: Graphical\_description\_of\_data-hist-boxplot-pie-barplot.R)

## Description graphique des données : distributions conditionnelles

```
1 xyplot(Ozone~Wind | Month, data=airquality , layout=c(5,1))
2 histogram(~Sepal.Length|Species, data=iris)
3 bwplot(Sepal.Length~Species, data=iris)
4
5 require(UsingR) #opens the functions and data of
6 #the package UsingR
7 attach(kid.weights) #allows to write height
8 #instead of kid.weights$height...
9 histogram(~weight|gender, data=kid.weights, layout=c(1,2))
10 histogram(~height|gender, data=kid.weights, layout=c(1,2))
11 mean(weight[gender=="M"])
12 mean(weight[gender=="F"])
13 mean(height[gender=="M"])
14 mean(height[gender=="F"])
15
16 H=read.csv(file="male_female_over_20_heights.csv")
17 str(H)
18 histogram(~H$height|H$gender, layout=c(1,2))
```

## Graphical description of data : correlations

```
1 require(UsingR) #opens the functions and data of
2                 #the package UsingR
3 data(kid.weights)
4 X=kid.weights; str(X)
5 pairs(X[,1:3] , col=X[,4] , pch=20)
6
7 ## BMI : body mass index
8 attach(kid.weights) #allows to write height
9                 #instead of kid.weights$height...
10 ht_in_m = height*2.54/100 # 2.54 cm per inch
11 weight_in_kg = weight / 2.2046 # 2.2046 lbs. per kg
12 bmi = weight_in_kg/ht_in_m^2
13 histogram(bmi)
14 histogram(~bmi|gender)
15
16 library(car)
17 states=as.data.frame(state.x77[,c("Murder", "Population",
18 "Illiteracy", "Income", "Frost")])
19 scatterplotMatrix(states , spread=T, pch=20, lty .smooth=2)
```

(file: Graphical description of data—correlations.R)



## Help, demos, examples

```
1 help("plot")
2
3 ?plot
4
5 help.search("plot")
6
7 ??plot
8
9 demo() # demonstrations list
10
11 demo( graphics )
12
13 example(plot) #the example function usually runs
14               #the examples that one can find at
15               #the end of the help files
```

(file: help\_demo\_examples.R)