

Introduction au logiciel R

et à la pratique des statistiques

en vue de l'analyse de données issues de la biologie

Julien Chiquet

École doctorale « *du génome aux organismes* »
Université d'Évry

30 janvier – 3 février 2012

http://stat.genopole.cnrs.fr/members/jchiquet/teachings/initiation_r

Laboratoire « Statistique & Génome »

<http://stat.genopole.cnrs.fr>



Christophe Ambroise



Professeur
statistiques

Julien Chiquet



Maître de Conférences
statistiques

`prenom.nom@genopole.cnrs.fr`

Laboratoire « Statistique & Génome »

<http://stat.genopole.cnrs.fr>

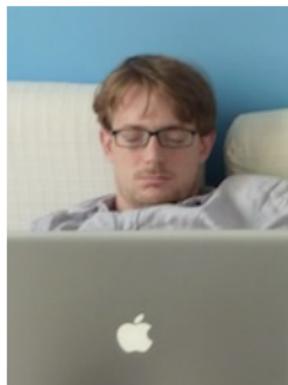


Christophe Ambroise



Professeur
statistiques

Julien Chiquet



Maître de Conférences
statistiques

`prenom.nom@genopole.cnrs.fr`

Agenda (théorique) de la semaine

Lundi Introduction & Syntaxe

Mardi Entrées / Sorties, Éléments de statistiques

Mercredi Tests d'hypothèses / ANOVA sous R

Jeudi Modèle linéaire sous R

Vendredi Outils R pour la biologie (vos données ?)

Cours

1. Introduction
2. Chapitre 1 : Structures de données
3. Chapitre 2 : Développer avec R
4. Chapitre 3 : Entrées / Sorties
5. Chapitre 4 : Statistiques et outils connexes

... et travaux dirigés !

Agenda (théorique) de la semaine

Lundi Introduction & Syntaxe

Mardi Entrées / Sorties, Éléments de statistiques

Mercredi Tests d'hypothèses / ANOVA sous R

Jeudi Modèle linéaire sous R

Vendredi Outils R pour la biologie (vos données ?)

Cours

1. Introduction
2. Chapitre 1 : Structures de données
3. Chapitre 2 : Développer avec R
4. Chapitre 3 : Entrées / Sorties
5. Chapitre 4 : Statistiques et outils connexes

... et travaux dirigés !

Agenda (théorique) de la semaine

Lundi Introduction & Syntaxe

Mardi Entrées / Sorties, Éléments de statistiques

Mercredi Tests d'hypothèses / ANOVA sous R

Jeudi Modèle linéaire sous R

Vendredi Outils R pour la biologie (vos données ?)

Cours

1. Introduction
2. Chapitre 1 : Structures de données
3. Chapitre 2 : Développer avec R
4. Chapitre 3 : Entrées / Sorties
5. Chapitre 4 : Statistiques et outils connexes

... et *travaux dirigés* !

Première partie I

Introduction

Avant de démarrer

Installation et premiers contacts

Une session exemple

Avant de démarrer

Installation et premiers contacts

Une session exemple

Qu'est-ce que R ?

En deux mots,

*R est un logiciel de développement scientifique spécialisé dans le calcul et l'**analyse statistique**.*

R est aussi

- ▶ un langage,
- ▶ un environnement,
- ▶ un projet open source (projet GNU),
- ▶ un logiciel multi-plateforme (Linux, Mac, Windows),

Qu'est-ce que R ?

En deux mots,

*R est un logiciel de développement scientifique spécialisé dans le calcul et l'**analyse statistique**.*

R est aussi

- ▶ un langage,
- ▶ un environnement,
- ▶ un projet open source (projet GNU),
- ▶ un logiciel multi-plateforme (Linux, Mac, Windows),
- ▶ la 18^e lettre de l'alphabet ☺.

Qu'est-ce que R ?

En deux mots,

*R est un logiciel de développement scientifique spécialisé dans le calcul et l'**analyse statistique**.*

R est aussi

- ▶ un langage,
- ▶ un environnement,
- ▶ un projet open source (projet GNU),
- ▶ un logiciel multi-plateforme (Linux, Mac, Windows),
- ▶ la 18^e lettre de l'alphabet ☹.

1. Gestionnaire de données
 - ▶ Lecture, manipulation, stockage.
2. Algèbre linéaire
 - ▶ Opérations classiques sur vecteurs, tableaux et matrices
3. Statistiques et analyse de données
 - ▶ Dispose d'un *grand* nombre de méthodes d'analyse de données (des plus anciennes et aux plus récentes)
4. Moteur de sorties graphiques
 - ▶ Sorties écran ou fichier
5. Système de modules
 - ▶ Alimenté par la communauté (+ de 2000 extensions !)
6. Interface « facile » avec C/C++, Fortran,...

1. Gestionnaire de **données**
 - ▶ Lecture, manipulation, stockage.
2. Algèbre linéaire
 - ▶ Opérations classiques sur vecteurs, tableaux et matrices
3. **Statistiques et analyse de données**
 - ▶ Dispose d'un *grand* nombre de méthodes d'analyse de données (des plus anciennes et aux plus récentes)
4. Moteur de **sorties graphiques**
 - ▶ Sorties écran ou fichier
5. Système de modules
 - ▶ Alimenté par la communauté (+ de 2000 extensions !)
6. Interface « facile » avec C/C++, Fortran,...

Approche chronologique

- 1970s développement de S au Bell labs.
- 1980s développement de S-PLUS au AT&T. Lab
- 1993 développement de R sur le modèle de S par Robert Gentleman et Ross Ihaka au département de statistique de l'université d'Auckland.
- 1995 dépôts des codes sources sous licence GNU/GPL
- 1997 élargissement du groupe
- 2002 la fondation R dépose ses statuts sous la présidence de Gentleman et Ihaka

Développement entièrement bénévole

- ▶ « R development core team » (12aine de personnes)
- ▶ Participation de *nombreux* chercheurs (>2000 packages)

Approche chronologique

- 1970s développement de S au Bell labs.
- 1980s développement de S-PLUS au AT&T. Lab
- 1993 développement de R sur le modèle de S par Robert Gentleman et Ross Ihaka au département de statistique de l'université d'Auckland.
- 1995 dépôts des codes sources sous licence GNU/GPL
- 1997 élargissement du groupe
- 2002 la fondation R dépose ses statuts sous la présidence de Gentleman et Ihaka

Développement entièrement bénévole

- ▶ « R development core team » (12aine de personnes)
- ▶ Participation de *nombreux* chercheurs (>2000 packages)

1. La page web de la **fondation R**
 - ▶ les statuts, des liens, des références.
 - ▶ <http://www.r-project.org/>
2. La page web du **CRAN** (Comprehensive R Arxiv Network)
 - ▶ binaires d'installation, packages, documentations, ...
 - ▶ <http://cran.r-project.org/>
3. La **conférence** des utilisateurs de R :
 - ▶ annuelle, prochaine édition à Nashville
 - ▶ <http://biostat.mc.vanderbilt.edu/wiki/Main/UseR-2012>
4. *The R journal* propose des articles sur
 - ▶ de nouvelles extensions, des applications, des actualités.
 - ▶ <http://journal.r-project.org/>
5. <http://onertipaday.blogspot.com/>

1. La page web de la **fondation R**
 - ▶ les statuts, des liens, des références.
 - ▶ <http://www.r-project.org/>
2. La page web du **CRAN** (Comprehensive R Arxiv Network)
 - ▶ binaires d'installation, packages, documentations, ...
 - ▶ <http://cran.r-project.org/>
3. La **conférence** des utilisateurs de R :
 - ▶ annuelle, prochaine édition à Nashville
 - ▶ <http://biostat.mc.vanderbilt.edu/wiki/Main/UseR-2012>
4. *The R journal* propose des articles sur
 - ▶ de nouvelles extensions, des applications, des actualités
 - ▶ <http://journal.r-project.org/>
5. <http://onertipaday.blogspot.com/>

1. La page web de la **fondation R**
 - ▶ les statuts, des liens, des références.
 - ▶ <http://www.r-project.org/>
2. La page web du **CRAN** (Comprehensive R Arxiv Network)
 - ▶ binaires d'installation, packages, documentations, ...
 - ▶ <http://cran.r-project.org/>
3. La **conférence** des utilisateurs de R : *useR!*
 - ▶ annuelle, prochaine édition à Nashville
 - ▶ <http://biostat.mc.vanderbilt.edu/wiki/Main/UseR-2012>
4. *The R journal* propose des articles sur
 - ▶ de nouvelles extensions, des applications, des actualités.
 - ▶ <http://journal.r-project.org/>
5. <http://onertipaday.blogspot.com/>

1. La page web de la **fondation R**
 - ▶ les statuts, des liens, des références.
 - ▶ <http://www.r-project.org/>
2. La page web du **CRAN** (Comprehensive R Arxiv Network)
 - ▶ binaires d'installation, packages, documentations, ...
 - ▶ <http://cran.r-project.org/>
3. La **conférence** des utilisateurs de R : *useR!*
 - ▶ annuelle, prochaine édition à Nashville
 - ▶ <http://biostat.mc.vanderbilt.edu/wiki/Main/UseR-2012>
4. *The R journal* propose des articles sur
 - ▶ de nouvelles extensions, des applications, des actualités.
 - ▶ <http://journal.r-project.org/>
5. <http://onertipaday.blogspot.com/>

1. La page web de la **fondation R**
 - ▶ les statuts, des liens, des références.
 - ▶ <http://www.r-project.org/>
2. La page web du **CRAN** (Comprehensive R Arxiv Network)
 - ▶ binaires d'installation, packages, documentations, ...
 - ▶ <http://cran.r-project.org/>
3. La **conférence** des utilisateurs de R : *useR!*
 - ▶ annuelle, prochaine édition à Nashville
 - ▶ <http://biostat.mc.vanderbilt.edu/wiki/Main/UseR-2012>
4. *The R journal* propose des articles sur
 - ▶ de nouvelles extensions, des applications, des actualités.
 - ▶ <http://journal.r-project.org/>
5. <http://onertipaday.blogspot.com/>

Plus ☺

1. Libre et gratuit,
2. Richesse des modules (en statistique),
3. Rapidité d'exécution,
4. Développement rapide (langage de scripts),
5. Syntaxe intuitive et compact,
6. Nombreuses possibilités graphiques.

Moins ☹

1. Aide intégrée succincte,
2. Debugger un peu sec,
3. Code parfois illisible (compacité),
4. Personnalisation des graphiques un peu lourde.

Plus 😊

1. Libre et gratuit,
2. Richesse des modules (en statistique),
3. Rapidité d'exécution,
4. Développement rapide (langage de scripts),
5. Syntaxe intuitive et compact,
6. Nombreuses possibilités graphiques.

Moins ☹️

1. Aide intégrée succincte,
2. Debugger un peu sec,
3. Code parfois illisible (compacité),
4. Personnalisation des graphiques un peu lourde.

Les logiciels de développement scientifique sont spécialisés en

1. algèbre linéaire

- ▶ Matlab (Mathworks), la référence,
- ▶ Scilab (INRIA), l'alternative libre,
- ▶ Octave (GNU), l'alternative open source ☺,

2. statistiques

- ▶ SAS (SAS Inc.), la référence,
- ▶ S-PLUS (TIBCO), le concurrent,
- ▶ R (GNU), l'alternative open source ☺,

3. calcul symbolique

- ▶ Mathematica (Wolfram), la référence,
- ▶ Maple (Maplesoft), la référence aussi,
- ▶ Maxima (GNU), l'alternative open source ☺.

Les logiciels de développement scientifique sont spécialisés en

1. algèbre linéaire

- ▶ Matlab (Mathworks), la référence,
- ▶ Scilab (INRIA), l'alternative libre,
- ▶ Octave (GNU), l'alternative open source ☺,

2. statistiques

- ▶ SAS (SAS Inc.), la référence,
- ▶ S-PLUS (TIBCO), le concurrent,
- ▶ R (GNU), l'alternative open source ☺,

3. calcul symbolique

- ▶ Mathematica (Wolfram), la référence,
- ▶ Maple (Maplesoft), la référence aussi,
- ▶ Maxima (GNU), l'alternative open source ☺.

Les logiciels de développement scientifique sont spécialisés en

1. algèbre linéaire

- ▶ Matlab (Mathworks), la référence,
- ▶ Scilab (INRIA), l'alternative libre,
- ▶ Octave (GNU), l'alternative open source ☺,

2. statistiques

- ▶ SAS (SAS Inc.), la référence,
- ▶ S-PLUS (TIBCO), le concurrent,
- ▶ R (GNU), l'alternative open source ☺,

3. calcul symbolique

- ▶ Mathematica (Wolfram), la référence,
- ▶ Maple (Maplesoft), la référence aussi,
- ▶ Maxima (GNU), l'alternative open source ☺.

► Obtenir de l'aide

<code>help -i</code>	<code>help.start ()</code>
<code>help</code>	<code>help(help)</code>
<code>help sort</code>	<code>help(sort) _or_ ?sort</code>

► Séquence de vecteurs

<code>1:10</code>	<code>1:10 _or_ seq(10)</code>
<code>1:3:10</code>	<code>seq(1,10,by=3)</code>
<code>10:-1:1</code>	<code>10:1</code>
<code>linspace(1,10,7)</code>	<code>seq(1,10,length=7)</code>

► Manipulation de vecteurs

<code>a=[2 7 8 5]</code>	<code>a <- c(2,7,8,5)</code>
<code>a=a[3:4]</code>	<code>a <- a[c(3,4)]</code>
<code>adash=[2 3 4 5]'</code>	<code>adash <- t(c(2,3,4,5))</code>

Avant de démarrer

Installation et premiers contacts

Une session exemple

Installation

Rendez-vous sur la page du CRAN <http://cran.r-project.org/>

Macintosh

Télécharger `R-2.14.1.pkg`, cliquer.

Windows

Télécharger `R-2.14.1-win32.exe`, cliquer (prier).

Linux

Systèmes supportants `apt` (Debian, Ubuntu, ...)

```
$ sudo apt-get update  
$ sudo apt-get install r-base
```

Premiers pas

```
$ R
R version 2.14.0 (2011-10-31)
Copyright (C) 2011 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
[...]
Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.
> 1+1
[1] 2
```

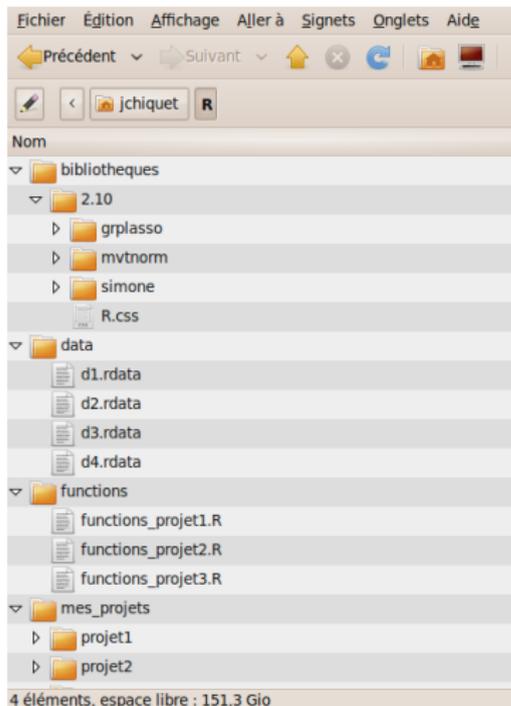
Sortez moi de là !

```
> q()
Save workspace image? [y/n/c]:y
```

↪ Sauve l'environnement et le réouvre la prochaine fois

Organiser un projet R

À calquer lors des travaux dirigés



- ▶ Dans un répertoire R, placer
 - ▶ un répertoire data
 - ▶ un répertoire mes_projets
 - ▶ un répertoire fonctions
- ▶ Créer un répertoire par projet
 - ▶ sauvegarde des données
`save.image(file = "f.RData")`
 - ▶ sauvegarde des instructions
`savehistory(file = "f.Rhistory")`
- ▶ bibliothèques contient les extensions installées.

FIGURE: Arborescence type

Environnement de travail sous Linux

Un bureau de développement avec R

```
File Edit Options Buffers Tools Imenu-S ESS Help
rm(list=ls())
library(mvtnorm)
source("fonctions.R")
source("fonctions_group_l1.R")

set.seed(1002)

## données simulés (settings de Yuan et Lin - papier
er de 2006)
n <- 100
Sigma <- matrix(c(1,0.5,0.5,1),2,2)
X <- rmvnorm(n, Sigma)
y <- X[,1]^3 + X[,1]^2 - 2 * X[,1] + (1/3)*X[,2]^3
- 0*X[,2]^2 + (2/3)*X[,2] + rnorm(n,0,3)
U:~ check_CoopLasso.R Top (11.0) SVN-385 (ESS[S][none])--15:50 0.47--
```

```
Fichier Édition Affichage Terminal Onglets Aide
Terminal Terminal Terminal
15:03 jchiquet@term14 ~/svn/notiid/branches/regressionCoop/R R >
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

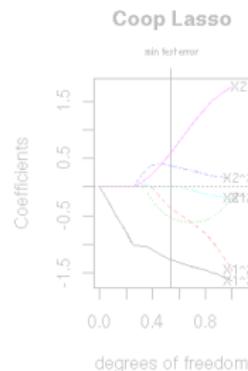
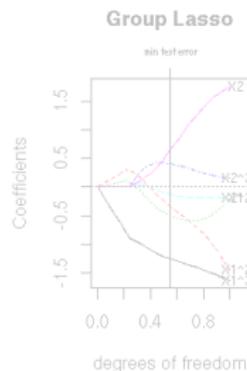
R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

> source("check_CoopLasso.R")
```

1. un éditeur de texte
2. un terminal avec R
3. des sorties graphiques



Environnement de travail sous Linux

Un bureau de développement avec R

```
File Edit Options Buffers Tools Imenu-S ESS Help
rm(list=ls())
library(mvtnorm)
source("functions.R")
source("functions_group_ll.R")

set.seed(1002)

## données simulés (settings de Yuan et Lin - papier
er de 2006)
n <- 100
Sigma <- matrix(c(1,0.5,0.5,1),2,2)
X <- rmvnorm(n, Sigma)
y <- X[,1]^3 + X[,1]^2 - 2 * X[,1] + (1/3)*X[,2]^3
  - 0*X[,2]^2 + (2/3)*X[,2] + rnorm(n,0,3)
U:~ check_CoopLasso.R Top (11.0) SVN-385 (ESS[S][none])--15:50 0.47~
```

```
Fichier Édition Affichage Terminal Onglets Aide
Terminal Terminal Terminal
15:03 jchiquet@term14 ~/svn/notiid/branches/regressionCoop/R R >
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

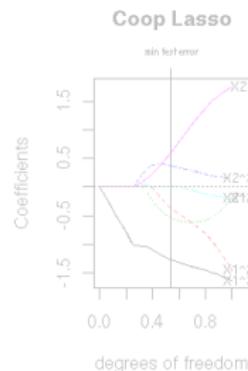
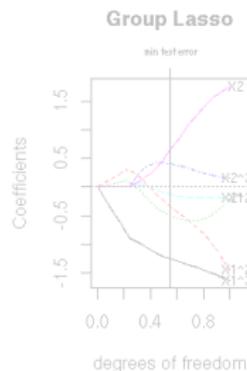
R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

> source("check_CoopLasso.R")
```

1. un éditeur de texte
2. un terminal avec R
3. des sorties graphiques



Environnement de travail sous Linux

Un bureau de développement avec R

```
File Edit Options Buffers Tools Imenu-S ESS Help
[Icons]
rm(list=ls())
library(mvtnorm)
source("functions.R")
source("functions_group_ll.R")

set.seed(1002)

## données simulés (settings de Yuan et Lin - papier
er de 2006)
n <- 100
Sigma <- matrix(c(1,0.5,0.5,1),2,2)
X <- rmvnorm(n, Sigma)
y <- X[,1]^3 + X[,1]^2 - 2 * X[,1] + (1/3)*X[,2]^3
  - 0*X[,2]^2 + (2/3)*X[,2] + rnorm(n,0,3)
U:~ check_CoopLasso.R Top (11.0) SVN-385 (ESS[S][none])--15:50 0.47~
```

```
Fichier Édition Affichage Terminal Onglets Aide
Terminal Terminal Terminal
15:03 jchiquet@term14 ~/svn/notiid/branches/regressionCoop/R R >
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

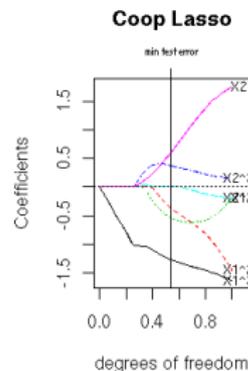
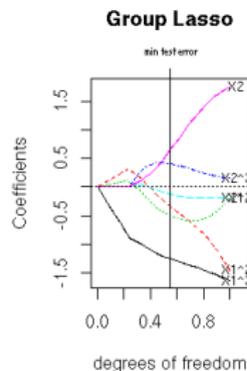
R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

> source("check_CoopLasso.R")
```

1. un éditeur de texte
2. un terminal avec R
3. des sorties graphiques



Depuis R

- ▶ `help(str)` : lance l'aide associée à la commande `str`,
- ▶ `help.search("factorial")` : cherche les commandes contenant le mot-clé `factorial`,
- ▶ `help.start()` : lance l'aide HTML.

Sur le Web

- ▶ **Le site du CRAN** : beaucoup (trop ?) de guides d'utilisations sont répertoriés (y compris ceux en français),
- ▶ le site du module en propose une sélection.

À tout moment

- ▶ la liste des commandes usuelles,
- ▶ le prof (pas infailible mais rapide d'accès).

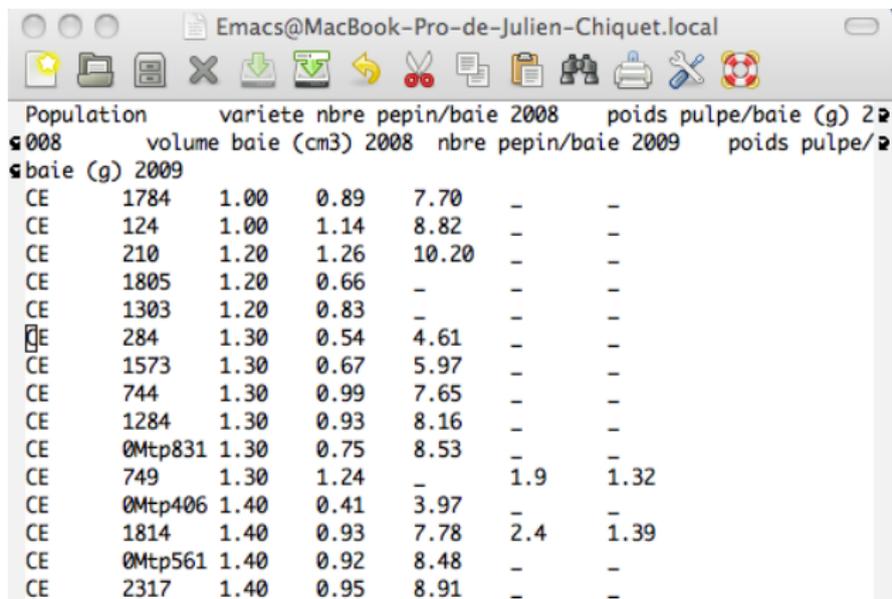
Avant de démarrer

Installation et premiers contacts

Une session exemple

Analyse élémentaire d'un jeu de données

Quelle tête ont les données ? On ouvre avec Emacs :



The image shows a screenshot of the Emacs text editor window. The title bar reads "Emacs@MacBook-Pro-de-Julien-Chiquet.local". The menu bar contains standard icons for Emacs. The main content area displays a table with the following data:

Population	variete	nbre pepin/baie 2008	volume baie (cm3) 2008	nbre pepin/baie 2009	poids pulpe/baie (g) 2008	poids pulpe/baie (g) 2009
CE	1784	1.00	0.89	7.70	-	-
CE	124	1.00	1.14	8.82	-	-
CE	210	1.20	1.26	10.20	-	-
CE	1805	1.20	0.66	-	-	-
CE	1303	1.20	0.83	-	-	-
CE	284	1.30	0.54	4.61	-	-
CE	1573	1.30	0.67	5.97	-	-
CE	744	1.30	0.99	7.65	-	-
CE	1284	1.30	0.93	8.16	-	-
CE	0Mtp831	1.30	0.75	8.53	-	-
CE	749	1.30	1.24	-	1.9	1.32
CE	0Mtp406	1.40	0.41	3.97	-	-
CE	1814	1.40	0.93	7.78	2.4	1.39
CE	0Mtp561	1.40	0.92	8.48	-	-
CE	2317	1.40	0.95	8.91	-	-

FIGURE: données baies de vignes 2008/2009

Je remplace tous les _ par du vide (R le comprendra mieux) !

Les commandes `getwd()` et `setwd()` gèrent le répertoire de travail :

```
> setwd("~/Documents/documents_in_progress/myteachings/Intro_R/cours_R/")
> getwd()

[1] "/Users/jchiquet/Documents/documents_in_progress/myteachings/Intro_R/cours_R"
```

Les données possèdent un entête et sont délimitées par des tabulations :

```
> donnees <- read.delim("mesures_baie_raisin_2008-2009.txt")
```

Qu'est-ce qui se trouve dorénavant dans mon itinéraire de recherche ?

```
> ls()

[1] "donnees"

> objects()

[1] "donnees"
```

Quelle tête ont mes données ?

```
> head(donnees)
```

```
  Population variete  nbre.pepin.baie.2008 poids.pulpe.baie..g..2008
1          CE      1784                1.0                0.89
2          CE       124                1.0                1.14
3          CE       210                1.2                1.26
4          CE      1805                1.2                0.66
5          CE      1303                1.2                0.83
6          CE       284                1.3                0.54
 volume.baie..cm3..2008  nbre.pepin.baie.2009 poids.pulpe.baie..g..2009
1                    7.70                    NA                    NA
2                    8.82                    NA                    NA
3                   10.20                    NA                    NA
4                     NA                    NA                    NA
5                     NA                    NA                    NA
6                    4.61                    NA                    NA
```

Je les mets dans mon itinéraire de recherche

```
> attach(donnees)
```

Et les attributs ?

```
> str(donnees)
```

```
'data.frame':      245 obs. of  7 variables:
 $ Population      : Factor w/  3 levels "CE","CO","TE": 1 1 1 1 1 1 1 1 1 1 ...
 $ variete         : Factor w/ 245 levels "OMtp1004","OMtp1005",...: 113 66 ...
 $ nbre.pepin.baie.2008 : num  1 1 1.2 1.2 1.2 1.3 1.3 1.3 1.3 1.3 ...
 $ poids.pulpe.baie..g..2008: num  0.89 1.14 1.26 0.66 0.83 0.54 0.67 0.99 0.93 0.7 ...
 $ volume.baie..cm3..2008  : num  7.7 8.82 10.2 NA NA 4.61 5.97 7.65 8.16 8.53 ...
 $ nbre.pepin.baie.2009   : num  NA ...
 $ poids.pulpe.baie..g..2009: num  NA ...
```

Le « nécessaire » résumé statistique :

```
> summary(donnees)
```

```
Population      variete      nbre.pépin.baie.2008 poids.pulpe.baie..g..2008
CE:84          OMtp1004: 1      Min.   : 0.000          Min.   : 0.390
CO:89          OMtp1005: 1      1st Qu.: 1.400          1st Qu.: 0.820
TE:72          OMtp1033: 1      Median : 1.800          Median : 1.060
               OMtp1068: 1      Mean   : 1.858          Mean   : 1.212
               OMtp1072: 1      3rd Qu.: 2.400          3rd Qu.: 1.360
               OMtp1073: 1      Max.   : 3.200          Max.   : 3.750
               (Other) :239    NA's   :27.000          NA's   :28.000

volume.baie..cm3..2008 nbre.pépin.baie.2009 poids.pulpe.baie..g..2009
Min.   : 3.44          Min.   : 1.000          Min.   : 0.560
1st Qu.: 7.14          1st Qu.: 1.500          1st Qu.: 0.875
Median : 8.91          Median : 2.000          Median : 1.230
Mean   :10.47          Mean   : 2.082          Mean   : 1.513
3rd Qu.:11.90          3rd Qu.: 2.600          3rd Qu.: 1.607
Max.   :33.80          Max.   : 3.600          Max.   : 4.340
NA's   :44.00          NA's   :196.000         NA's   :203.000
```

Et si je veux le nombre de baies moyen en 2008 pour chaque population ?

```
> tapply(nbre.pepin.baie.2008,Population,mean,na.rm=TRUE)
```

```
      CE      CO      TE  
1.850649 2.034667 1.666667
```

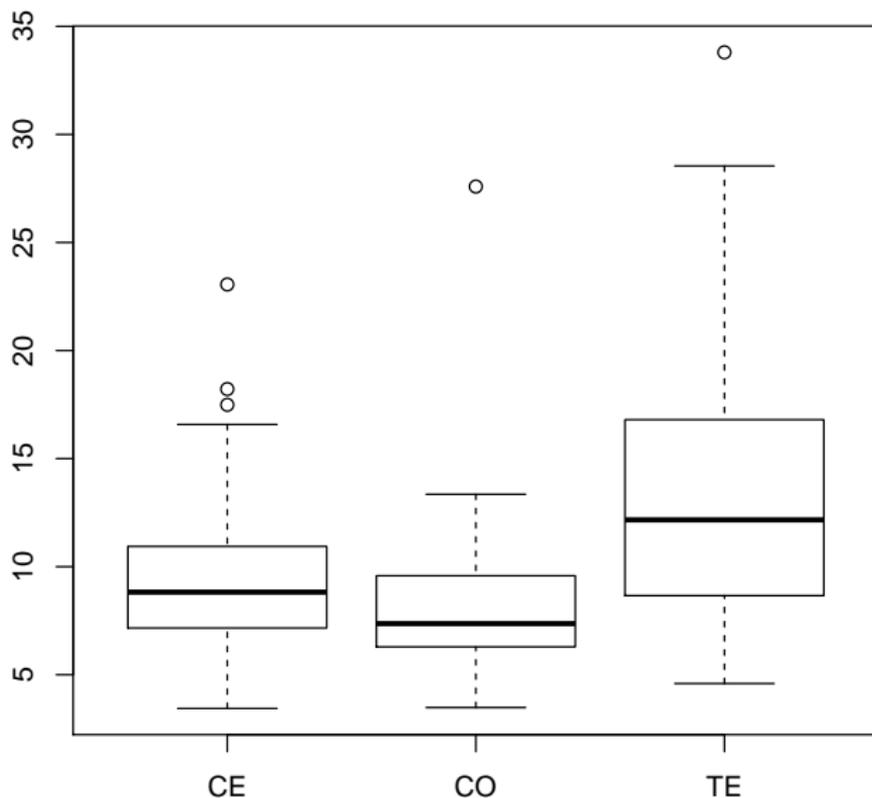
Et le volume moyen des baies ?

```
> tapply(volume.baie..cm3..2008,Population,mean,na.rm=TRUE)
```

```
      CE      CO      TE  
9.667971 8.003000 14.132097
```

Ça a l'air disparate. Qu'est-ce que ça donne graphiquement ?

```
> boxplot(volume.baie..cm3..2008 ~ Population)
```



Finalement, y-a-t-il un effet « population » pour le volume des baies ?

```
> anova(lm(volume.baie..cm3..2008 ~ Population))
```

Analysis of Variance Table

Response: volume.baie..cm3..2008

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Population	2	1301.9	650.94	29.45	6.357e-12 ***
Residuals	198	4376.5	22.10		

Signif. codes: 0

etc. *beaucoup* d'autres choses sont possibles. . .

Deuxième partie II

Structures de données

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Propriétés

- ▶ objet le plus **élémentaire** sous R,
- ▶ collection d'entités **de même nature**,
- ▶ **mode** (ou type) défini par la nature des entités qui le composent.

Les modes possibles

1. numérique (`numeric`),
2. caractère (`character`),
3. logique (`boolean`).

Propriétés

- ▶ objet le plus **élémentaire** sous \mathbb{R} ,
- ▶ collection d'entités **de même nature**,
- ▶ **mode** (ou type) défini par la nature des entités qui le composent.

Les modes possibles

1. numérique (`numeric`),
2. caractère (`character`),
3. logique (`boolean`).

Propriétés

- ▶ objet le plus **élémentaire** sous \mathbb{R} ,
- ▶ collection d'entités **de même nature**,
- ▶ **mode** (ou type) défini par la nature des entités qui le composent.

Les modes possibles

1. numérique (`numeric`),
2. caractère (`character`),
3. logique (`boolean`).

1. Numérique

```
> x0 <- 0
> x1 <- c(-1, 23, 98.7)
> mode(x0)

[1] "numeric"
```

2. Caractère

```
> y0 <- "bonjour"
> y1 <- c("Pomme", "Flore", "Alexandre")
> mode(y1)

[1] "character"
```

3. Logique

```
> z0 <- TRUE
> z1 <- c(FALSE, TRUE, FALSE, TRUE, TRUE)
> z2 <- c(T, F, F)
> mode(z2)

[1] "logical"
```

Définition (affectation)

C'est l'opération qui consiste à *attribuer une valeur* à une variable.

En R, plusieurs choix sont possibles :

- ▶ l'opérateur usuel est '`<-`' (signe inférieur suivi du signe moins)

```
> jo <- "l'indien"  
> jo  
[1] "l'indien"
```

- ▶ l'opérateur '`=`' peut être utilisé la plupart du temps

```
> nb.max.d.annees.pour.faire.une.these = 3  
> nb.max.d.annees.pour.faire.une.these  
[1] 3
```

- ▶ la commande `assign` permet cette opération (d'où l'anglicisme *assignation*)

```
> assign("x", c(8, 9, -pi, sqrt(2)))  
> x  
[1] 8.000000 9.000000 -3.141593 1.414214
```

VARIABLES RÉSERVÉES PAR R

- ▶ NA est le code R pour les valeurs manquantes (absentes des données),
- ▶ NaN est le code de R pour signifier un résultat numérique aberrant ,
- ▶ Inf et -Inf sont les valeurs réservées pour plus et moins ∞ ,
- ▶ NULL est l'objet nul.

```
> c(4, 2, NA, 5)
```

```
[1] 4 2 NA 5
```

```
> 0/0
```

```
[1] NaN
```

```
> 1/0
```

```
[1] Inf
```

```
> names(1)
```

```
NULL
```

Vecteurs

- Les modes ou typages

- Opérations élémentaires**

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Soient x, y tels que

```
> x <- c(1, 2, -3, -4)
```

```
> y <- c(-5, -6, 9, 0)
```

'+' addition des éléments de deux vecteurs

```
> x + y
```

```
[1] -4 -4 6 -4
```

'-' soustraction des éléments de deux vecteurs

```
> x - y
```

```
[1] 6 8 -12 -4
```

'*' multiplication des éléments de deux vecteurs

```
> x * y
```

```
[1] -5 -12 -27 0
```

'/' division des éléments de deux vecteurs

```
> x/y
```

```
[1] -0.2000000 -0.3333333 -0.3333333 -Inf
```

Le « recyclage » des éléments du vecteur

Lors d'une opération entre vecteurs, les vecteurs trop courts sont ajustés pour atteindre la taille du plus grand vecteur en recyclant les données.

Exemple

```
> x <- c(10, 100, 1000)
> y <- c(1, 2)
> 2 * x + y - 1
[1] 20 201 2000
```

↪ souvent pratique mais **attention aux effets de bords!**

Fonctions numériques élémentaires

`floor`, `ceiling`, `round`.

```
> floor(2/3)
```

```
[1] 0
```

```
> ceiling(2/3)
```

```
[1] 1
```

```
> round(2/3, 3)
```

```
[1] 0.667
```

Fonctions arithmétiques élémentaires

`^`, `%%`, `%/%`, `abs`, `log`, `exp`, `log10`, `sqrt`, `cos`, `tan`, `sin`... s'appliquent toutes **terme-à-terme**.

```
> log10(c(10, 100, 1000))
```

```
[1] 1 2 3
```

```
> cos(c(pi/2, pi))^2 + sin(c(pi/2, pi))^2
```

```
[1] 1 1
```

Fonctions caractérisant un vecteur

prod, sum, max, min, range, which.min, which.max, length

```
> x <- c(-8, 1.5, 3)
```

```
> prod(x)
```

```
[1] -36
```

```
> sum(x)
```

```
[1] -3.5
```

```
> length(x)
```

```
[1] 3
```

```
> max(x)
```

```
[1] 3
```

```
> min(x)
```

```
[1] -8
```

```
> range(x)
```

```
[1] -8 3
```

```
> which.max(x)
```

```
[1] 3
```

```
> which.min(x)
```

```
[1] 1
```

Pour le minimum / maximum terme-à-terme : `pmin`, `pmax`.

Fonctions appliquées le long du vecteur

`cumsum`, `cumprod`, `cummin`, `cummax`

```
> x <- c(-2, 1, -3, 2)
```

```
> cumprod(x)
```

```
[1] -2 -2 6 12
```

```
> cumsum(x)
```

```
[1] -2 -1 -4 -2
```

```
> cummax(x)
```

```
[1] -2 1 1 2
```

```
> cummin(x)
```

```
[1] -2 -2 -3 -3
```

Fonctionnent pour tous les modes

`unique`, `intersect`, `union`, `setdiff`, `setequal`, `is.element`

```
> unique(c("banane", "citron", "banane"))
```

```
[1] "banane" "citron"
```

```
> intersect(c("banane", "citron"), c("orange", "banane"))
```

```
[1] "banane"
```

```
> union(c("banane", "citron"), c("orange", "banane"))
```

```
[1] "banane" "citron" "orange"
```

```
> setequal(c("banane", "citron"), c("orange", "banane"))
```

```
[1] FALSE
```

```
> is.element(1, sample(c(1, 2, 3), 2))
```

```
[1] TRUE
```

Vecteurs

Les modes ou typages

Opérations élémentaires

Génération de vecteurs

Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

Définition, création

Manipulation de matrices

Opérateurs d'algèbre linéaire

Listes et Tableaux de données

`from:to`

Génère une séquence par pas de un depuis le nombre `from` jusqu'à `to` (si possible).

```
> -5:5
```

```
[1] -5 -4 -3 -2 -1  0  1  2  3  4  5
```

```
> 5:-5
```

```
[1]  5  4  3  2  1  0 -1 -2 -3 -4 -5
```

```
> pi:6
```

```
[1] 3.141593 4.141593 5.141593
```

```
> 1:6/2
```

```
[1] 0.5 1.0 1.5 2.0 2.5 3.0
```

```
> 1:(6/2)
```

```
[1] 1 2 3
```

Plusieurs schémas possibles

- ▶ `seq(from,to)`
- ▶ `seq(from,to,by=)`
- ▶ `seq(from,to,length.out=)`

```
> seq(1, 10)
```

```
[1] 1 2 3 4 5 6 7 8 9 10
```

```
> seq(2, 10, by = 2)
```

```
[1] 2 4 6 8 10
```

```
> seq(2, 10, length.out = 6)
```

```
[1] 2.0 3.6 5.2 6.8 8.4 10.0
```

Fonctionne pour tous les modes

- ▶ `rep(x, times)`, où `times` peut être un vecteur,
- ▶ `rep(x, each)`.

```
> rep(1, 3)
```

```
[1] 1 1 1
```

```
> rep("Mercy", 3)
```

```
[1] "Mercy" "Mercy" "Mercy"
```

```
> rep(c("A", "B", "C"), c(3, 2, 4))
```

```
[1] "A" "A" "A" "B" "B" "C" "C" "C" "C"
```

```
> rep(c(TRUE, FALSE), each = 2)
```

```
[1] TRUE TRUE FALSE FALSE
```

Obtenus par conditions avec

- ▶ les opérateurs logiques '`<`', '`<=`', '`>`', '`>=`', '`==`' '`!=`'
- ▶ le ET, le OU, NON, OU exclusif : '`&`' (intersection), '`|`' (union), '`!`' (négation), `xor`.

```
> note1 <- c(8, 9, 14, 3, 17.5, 11)
> note2 <- c("C", "B", "A", "B", "E", "B")
> admis <- (note1 >= 10) & (note2 == "A" | note2 == "B")
> mention <- (note1 >= 15) & (note2 == "A")
> admis
[1] FALSE FALSE TRUE FALSE FALSE TRUE
> sum(admis)
[1] 2
> sum(mention)
[1] 0
```

Avec 'c()'

L'opérateur 'c()' peut s'appliquer à n'importe quoi pourvu que l'on concatène des vecteurs de même type.

```
> c(c(1, 2), c(3, 4))
```

```
[1] 1 2 3 4
```

```
> round(c(seq(-pi, pi, len = 4), rep(c(1:3), each = 2), 0), 2)
```

```
[1] -3.14 -1.05 1.05 3.14 1.00 1.00 2.00 2.00 3.00 3.00 0.00
```

Remarque

Dans le second exemple, les entiers composants c(1:3) ont été forcés au typage flottant.

Avec `paste`

Concaténation de chaînes de caractères. Convertit en caractères les éléments passés en argument avant toute opération.

```
> paste("R", "c'est", "bien")
```

```
[1] "R c'est bien"
```

```
> paste(2:4, "ieme")
```

```
[1] "2 ieme" "3 ieme" "4 ieme"
```

```
> paste("A", 1:5, sep = "")
```

```
[1] "A1" "A2" "A3" "A4" "A5"
```

```
> paste("A", 1:5, sep = "", collapse = "")
```

```
[1] "A1A2A3A4A5"
```

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs**

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Principe

- ▶ Permet la **sélection d'un sous-ensemble** du vecteur x .
- ▶ Le sous-ensemble est spécifié **entre crochets** $x[\text{subset}]$.

L'objet `subset` peut prendre 4 types différents :

1. **un vecteur logique**, qui doit être de la même taille que le vecteur x ;
2. **un vecteur numérique aux composantes positives**, qui spécifie les valeurs à inclure ;
3. **un vecteur numérique aux composantes négatives**, qui spécifie les valeurs à exclure ;
4. **un vecteur de chaînes de caractères**, qui spécifie les noms des éléments de x à conserver.

Principe

- ▶ Permet la sélection d'un sous-ensemble du vecteur x .
- ▶ Le sous-ensemble est spécifié **entre crochets** $x[\text{subset}]$.

L'objet `subset` peut prendre 4 types différents :

1. **un vecteur logique**, qui doit être de la même taille que le vecteur x ;
2. **un vecteur numérique aux composantes positives**, qui spécifie les valeurs à inclure ;
3. **un vecteur numérique aux composantes négatives**, qui spécifie les valeurs à exclure ;
4. **un vecteur de chaînes de caractères**, qui spécifie les noms des éléments de x à conserver.

Principe

- ▶ Permet la sélection d'un sous-ensemble du vecteur x .
- ▶ Le sous-ensemble est spécifié **entre crochets** $x[\text{subset}]$.

L'objet `subset` peut prendre 4 types différents :

1. **un vecteur logique**, qui doit être de la même taille que le vecteur x ;
2. **un vecteur numérique aux composantes positives**, qui spécifie les valeurs à inclure ;
3. **un vecteur numérique aux composantes négatives**, qui spécifie les valeurs à exclure ;
4. **un vecteur de chaînes de caractères**, qui spécifie les noms des éléments de x à conserver.

Principe

- ▶ Permet la sélection d'un sous-ensemble du vecteur x .
- ▶ Le sous-ensemble est spécifié **entre crochets** $x[\text{subset}]$.

L'objet `subset` peut prendre 4 types différents :

1. **un vecteur logique**, qui doit être de la même taille que le vecteur x ;
2. **un vecteur numérique aux composantes positives**, qui spécifie les valeurs à inclure ;
3. **un vecteur numérique aux composantes négatives**, qui spécifie les valeurs à exclure ;
4. **un vecteur de chaînes de caractères**, qui spécifie les noms des éléments de x à conserver.

Principe

- ▶ Permet la sélection d'un sous-ensemble du vecteur x .
- ▶ Le sous-ensemble est spécifié **entre crochets** $x[\text{subset}]$.

L'objet `subset` peut prendre 4 types différents :

1. **un vecteur logique**, qui doit être de la même taille que le vecteur x ;
2. **un vecteur numérique aux composantes positives**, qui spécifie les valeurs à inclure ;
3. **un vecteur numérique aux composantes négatives**, qui spécifie les valeurs à exclure ;
4. **un vecteur de chaînes de caractères**, qui spécifie les noms des éléments de x à conserver.

Vecteurs logiques

```
> x <- c(3, 6, -2, 9, NA, sin(-pi/6))
```

```
> x[x > 0]
```

```
[1] 3 6 9 NA
```

```
> x[!is.na(x)]
```

```
[1] 3.0 6.0 -2.0 9.0 -0.5
```

```
> x[!is.na(x) & x > 0]
```

```
[1] 3 6 9
```

```
> mean(x, na.rm = TRUE)
```

```
[1] 3.1
```

```
> x[x <= mean(x, na.rm = TRUE)]
```

```
[1] 3.0 -2.0 NA -0.5
```

Vecteurs aux composantes positives ou négatives

```
> x
[1] 3.0 6.0 -2.0 9.0 NA -0.5

> x[2]
[1] 6

> x[1:5]
[1] 3 6 -2 9 NA

> x[-c(1, 5)]
[1] 6.0 -2.0 9.0 -0.5

> x[-(1:5)]
[1] -0.5
```

Vecteurs de chaînes de caractères

```
> names(x) <- c("var1", "var2", "var3", "var4", "var5", "var6")
```

```
> x
```

```
var1 var2 var3 var4 var5 var6  
 3.0  6.0 -2.0  9.0  NA -0.5
```

```
> x[c("var1", "var3")]
```

```
var1 var3  
  3   -2
```

1. Classer

- ▶ `sort` renvoie le vecteur classé par ordre croissant ou décroissant,
- ▶ `order` renvoie les indices d'ordre des éléments par ordre croissant ou décroissant,

2. Extraire

- ▶ `which` renvoie les indices de `x` vérifiant une condition ;

3. Échantillonner

- ▶ `sample` échantillonne aléatoirement dans un vecteur `x`, avec ou sans remise.

1. Classer

- ▶ `sort` renvoie le vecteur classé par ordre croissant ou décroissant,
- ▶ `order` renvoie les indices d'ordre des éléments par ordre croissant ou décroissant,

2. Extraire

- ▶ `which` renvoie les indices de `x` vérifiant une condition ;

3. Échantillonner

- ▶ `sample` échantillonne aléatoirement dans un vecteur `x`, avec ou sans remise.

1. Classer

- ▶ `sort` renvoie le vecteur classé par ordre croissant ou décroissant,
- ▶ `order` renvoie les indices d'ordre des éléments par ordre croissant ou décroissant,

2. Extraire

- ▶ `which` renvoie les indices de x vérifiant une condition ;

3. Échantillonner

- ▶ `sample` échantillonne aléatoirement dans un vecteur x , avec ou sans remise.

Exemples

```
> x <- -5:5
> y <- sample(x)
> sort(y)

[1] -5 -4 -3 -2 -1  0  1  2  3  4  5

> order(y)

[1]  2  6  4  9  8  3 10  1  5 11  7

> y[order(y)]

[1] -5 -4 -3 -2 -1  0  1  2  3  4  5

> y[order(y, decreasing = TRUE)]

[1]  5  4  3  2  1  0 -1 -2 -3 -4 -5

> which(sample(x, 4) > 0)

[1] 1
```

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Définition

*Un facteur est un vecteur de **variables catégorielles**. Les niveaux du facteur peuvent être ordonnés ou pas.*

Utilisation

les facteurs s'utilisent pour **catégoriser les données** d'un vecteur (ce qui s'avère très utile pour la gestion des variables qualitatives).

↪ un facteur est souvent associé à d'autres vecteurs pour en définir une **partition**.

Définition

*Un facteur est un vecteur de **variables catégorielles**. Les niveaux du facteur peuvent être ordonnés ou pas.*

Utilisation

les facteurs s'utilisent pour **catégoriser les données** d'un vecteur (ce qui s'avère très utile pour la gestion des variables qualitatives).

↪ un facteur est souvent associé à d'autres vecteurs pour en définir une **partition**.

Définition

*Un facteur est un vecteur de **variables catégorielles**. Les niveaux du facteur peuvent être ordonnés ou pas.*

Utilisation

les facteurs s'utilisent pour **catégoriser les données** d'un vecteur (ce qui s'avère très utile pour la gestion des variables qualitatives).

↪ un facteur est souvent associé à d'autres vecteurs pour en définir une **partition**.

Création : la fonction factor

```
> factor(sample(1:3, 10, replace = TRUE))
```

```
[1] 3 3 1 2 1 2 3 1 1 1
```

```
Levels: 1 2 3
```

```
> factor(sample(1:3, 10, replace = TRUE), levels = 1:5)
```

```
[1] 3 1 2 2 3 2 1 2 2 1
```

```
Levels: 1 2 3 4 5
```

Gestion : nlevels, levels, table

```
> x <- factor(sample(c("thésard", "CR", "MdC"), 15, replace = TRUE))
```

```
> cat(nlevels(x), "niveaux:", levels(x))
```

```
3 niveaux: CR MdC thésard
```

```
> table(x)
```

```
x
```

CR	MdC	thésard
4	8	3

Un exemple de facteur associé à un vecteur

Un exemple de facteur associé à un vecteur

Données

Chacun me donne son âge et son grade ¹

```
> age <- c(25, 35, 32, 27, 32, 40, 26, 25, 26, 28, 30, NA, 36,  
+         30, 30)  
> grd <- c("thd", "CR", "MdC", "thd", "thd", "MdC", "MdC", "thd",  
+         "thd", "MdC", "CR", "MdC", "CR", "thd", "thd")
```

Question : nombre d'individus par catégorie ?

```
> table(grd)
```

```
grd  
CR MdC thd  
 3   5   7
```

1. sauf un qui refuse :'(

Utilisation

Applique une fonction sur un vecteur partitionné en groupes.

Question : âge moyen / écart-type par catégorie ?

```
> tapply(age, grd, mean, na.rm = TRUE)
```

```
      CR      MdC      thd  
33.66667 31.50000 27.85714
```

```
> tapply(age, grd, sd, na.rm = TRUE)
```

```
      CR      MdC      thd  
3.214550 6.191392 2.794553
```

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Définition (objet array)

*Un tableau est un vecteur muni d'un attribut dimension (*dim*), lui même défini par un vecteur. Il est défini par la commande `array(data, dim, dimnames=)`*

```
> array(1:8, c(2, 2, 2))
```

```
, , 1
```

	[,1]	[,2]
[1,]	1	3
[2,]	2	4

```
, , 2
```

	[,1]	[,2]
[1,]	5	7
[2,]	6	8

Définition (objet `matrix`)

Une matrice est un tableau à deux dimensions. Elle est définie par la commande

```
matrix(data, nrow=, ncol=, byrow)
```

En conséquence

- ▶ Un objet `array` à deux dimensions est automatiquement converti en `matrix`
- ▶ Un vecteur auquel on ajoute un attribut `dimension` est automatiquement converti en `matrix`

```
> class(array(1:4, c(2, 2)))
```

```
[1] "matrix"
```

```
> x <- c(1, 2, 3, 4)
```

```
> dim(x) <- c(2, 2)
```

```
> class(x)
```

```
[1] "matrix"
```

1. R range les éléments d'une matrice par défaut par **colonne**.

```
> matrix(1:6, nrow = 2)
```

```
      [,1] [,2] [,3]  
[1,]    1    3    5  
[2,]    2    4    6
```

```
> matrix(1:6, nrow = 2, byrow = TRUE)
```

```
      [,1] [,2] [,3]  
[1,]    1    2    3  
[2,]    4    5    6
```

2. Lors de la création d'une matrice, R **recycle** les éléments jusqu'à ce que les contraintes de dimension soient vérifiées.

```
> matrix(1:3, nrow = 2, ncol = 2)
```

```
      [,1] [,2]  
[1,]    1    3  
[2,]    2    1
```

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Étant donné qu'une matrice est un vecteur pourvu d'une dimension, on a la proposition suivante :

Proposition

La plupart des opérateurs vectorielles s'appliquent (arithmétiques/mathématiques, ensemblistes, d'indexation).

```
> a <- matrix(sample(-4:4, 9), 3, 3)
> cat(max(a), sum(a), prod(a))
```

```
4 0 0
```

```
> which(a > 0)
```

```
[1] 1 2 8 9
```

```
> cumsum(a[a > 0])
```

```
[1] 1 4 6 10
```

```
> order(a)
```

```
[1] 6 5 7 4 3 1 8 2 9
```

```
> round(exp(a), 4)
```

```
      [,1]  [,2]  [,3]
[1,]  2.7183 0.3679  0.1353
[2,] 20.0855 0.0498  7.3891
[3,]  1.0000 0.0183 54.5982
```

Opérateurs matriciels usuels

- ▶ `+`, `/`, `*`, `^` sont les opérateurs usuels terme-à-terme,
- ▶ `%*%` est le produit matriciel,
- ▶ `crossprod()` est le produit scalaire,
- ▶ `t()` transpose une matrice,
- ▶ `diag()` extrait / spécifie la diagonale.

```
> a <- matrix(sample(-4:4, 9), 3, 3)
> b <- matrix(sample(a), 3, 3)
> diag(a)
[1] -3 -4  4
> diag(a) <- diag(b) <- 1
> diag(a)
[1] 1 1 1
> a + t(b) %*% b
      [,1] [,2] [,3]
[1,]   22  -11   1
[2,]  -10   22  -2
[3,]    2   -5  12
```

Concaténation de matrices

Trois fonctions selon l'effet voulu :

1. `c()` concatène les éléments de plusieurs matrices en un vecteur,
2. `cbind()` empile **horizontalement** plusieurs matrices,
3. `rbind()` empile **verticalement** plusieurs matrices.

```
> a <- matrix(1, 2, 3)
> b <- matrix(2, 2, 3)
> c(a, b)

[1] 1 1 1 1 1 1 1 2 2 2 2 2 2
```

```
> cbind(a, b)

      [,1] [,2] [,3] [,4] [,5] [,6]
[1,]    1    1    1    2    2    2
[2,]    1    1    1    2    2    2
```

```
> rbind(a, b)

      [,1] [,2] [,3]
[1,]    1    1    1
[2,]    1    1    1
[3,]    2    2    2
[4,]    2    2    2
```

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Résolution de systèmes linéaires, inversion matricielle

La commande `solve` résout

$$Ax = b,$$

```
> A <- matrix(c(4, 2, 8, -3), 2, 2)
> b <- c(2, 3)
> solve(A, b)
```

```
[1] 1.0714286 -0.2857143
```

ou inverse une matrice :

```
> round(solve(A) %*% A, 8)
```

```
      [,1] [,2]
[1,]    1    0
[2,]    0    1
```

R dispose des outils classiques d'algèbre linéaire

- ▶ `det` : calcule le **déterminant** d'une matrice ;
- ▶ `chol` : factorisation de **Cholesky** ($A = C^T C$, avec A symétrique, C triangulaire supérieure) ;
- ▶ `qr` : factorisation **QR** ($A = QR$ avec Q orthogonale, R triangulaire supérieure) ;
- ▶ `eigen` : calcule valeurs propres et **vecteurs propres** d'une matrice ;
- ▶ `svd` : calcule la décomposition en **valeurs singulières**.
- ▶ ...

Vecteurs

- Les modes ou typages

- Opérations élémentaires

- Génération de vecteurs

- Manipulation de vecteurs

Facteurs

Matrices (et tableaux)

- Définition, création

- Manipulation de matrices

- Opérateurs d'algèbre linéaire

Listes et Tableaux de données

Définition (objet list)

Une liste est une *collection d'objets hétérogènes*. Elle est définie par la commande `list(e11=, e12=, ...)`. Les éléments d'une liste peuvent posséder un nom.

```
> list(c(1,2,3),c("robert","johnson"),matrix(rnorm(4),2,2))
```

```
[[1]]  
[1] 1 2 3
```

```
[[2]]  
[1] "robert" "johnson"
```

```
[[3]]  
      [,1]      [,2]  
[1,] 0.2913862 0.3303319  
[2,] -1.4849716 2.4071609
```

```
> list(numero = c(1,2,3), noms = c("robert","johnson"), mat = matrix(rnorm(4),2,2))
```

```
$numero  
[1] 1 2 3
```

```
$noms  
[1] "robert" "johnson"
```

```
$mat  
      [,1]      [,2]  
[1,] 1.4355363 -0.8544773  
[2,] 0.6386084 1.9803265
```

Deux situations

1. Les éléments de la liste **ne sont pas nommés** : on accède au i^{e} élément par indexation `nom_liste[[i]]` uniquement.
2. Les éléments de la liste **sont nommés** : on peut y accéder comme ci-dessus ou en utilisant le nom de l'élément `nom_liste$nom_elt`.

```
> maliste <- list(numero = c(1, 2, 3), noms = c("robert", "johnson"),
+   mat = matrix(rnorm(4), 2, 2))
> maliste$nom
[1] "robert" "johnson"
> maliste$nom[2]
[1] "johnson"
> maliste[[2]]
[1] "robert" "johnson"
> maliste[[2]][2]
[1] "johnson"
```

Sélectionner des éléments

Fonctionne (presque) comme pour les vecteurs

```
> l1 <- list(1:2, c("a", "c", "g", "t"))
> l1[[-2]]

[1] 1 2
```

Commande lapply

Applique une fonction à chaque élément d'une liste

```
> lapply(maliste, length)

$numero
[1] 3

$noms
[1] 2

$mat
[1] 4
```

Commande `c()`

Permet de concaténer deux listes.

```
> c(list(1:2, c("a", "c", "g", "t")), list(rnorm(3), "yop"))
```

```
[[1]]
```

```
[1] 1 2
```

```
[[2]]
```

```
[1] "a" "c" "g" "t"
```

```
[[3]]
```

```
[1] -0.03599559 0.15875263 1.44373837
```

```
[[4]]
```

```
[1] "yop"
```

Définition (objet `data.frame`)

C'est une liste à laquelle on impose certaines contraintes², afin de rassembler vecteurs et facteurs sous la forme d'un tableau de données.

- ▶ *Pratiquement, un tableau de données est une matrice dont les colonnes sont de mode différent,*
- ▶ C'est l'objet idéal pour la **manipulation de données** (**forcez-vous** à l'utiliser).

2. que je vous épargne

Définition (objet `data.frame`)

C'est une liste à laquelle on impose certaines contraintes², afin de rassembler vecteurs et facteurs sous la forme d'un tableau de données.

- ▶ *Pratiquement, un tableau de données est une matrice dont les colonnes sont de mode différent,*
- ▶ C'est l'objet idéal pour la **manipulation de données** (**forcez-vous** à l'utiliser).

2. que je vous épargne

Syntaxe

On peut spécifier le nom des colonnes par le vecteur `row.names` ou directement comme pour une liste :

```
data.frame(e1=,e2=,...,row.names=)
```

```
> age <- c(25, 35, 32, 27, 32, 40, 26, 25, 26, 28, 30, NA, 36,
+         30, 30)
> grd <- c("thd", "CR", "MdC", "thd", "thd", "MdC", "MdC", "thd",
+         "thd", "MdC", "CR", "MdC", "CR", "thd", "thd")
> sex <- factor(sample(c(rep("M", 3), rep("F", 12))))
> donnees <- data.frame(age = age, grade = grd, sexe = sex)
> head(donnees)
```

	age	grade	sexe
1	25	thd	F
2	35	CR	M
3	32	MdC	F
4	27	thd	F
5	32	thd	F
6	40	MdC	F

Manipulation des éléments du tableau de données

- ▶ Comme une liste !
- ▶ les commandes `attach()` / `detach` placent / ôtent les éléments du tableaux de données dans l'itinéraire de recherche.

```
> donnees$age
```

```
[1] 25 35 32 27 32 40 26 25 26 28 30 NA 36 30 30
```

```
> attach(donnees, warn.conflicts = FALSE)
```

```
> grade
```

```
[1] thd CR MdC thd thd MdC MdC thd thd MdC CR MdC CR thd thd
```

```
Levels: CR MdC thd
```

```
> detach(donnees)
```

- ▶ beaucoup de fonctions prédéfinies
- ▶ penser aux fonctions `tapply` (ou `by`)

```
> summary(donnees)
      age      grade  sexe
Min.   :25.00   CR :3   F:12
1st Qu.:26.25   MdC:5   M: 3
Median :30.00   thd:7
Mean   :30.14
3rd Qu.:32.00
Max.   :40.00
NA's   : 1.00

> attach(donnees, warn.conflicts = FALSE)
> by(age, sexe, mean, na.rm = TRUE)

sexe: F
[1] 29.36364
-----

sexe: M
[1] 33

> by(age, grade, mean, na.rm = TRUE)

grade: CR
[1] 33.66667
-----

grade: MdC
[1] 31.5
-----

grade: thd
[1] 27.85714

> detach(donnees)
```

Idéal pour analyse statistique type anova

Pour plus tard. . .

```
> anova(lm(age ~ sexe * grade, data = donnees))
```

Analysis of Variance Table

Response: age

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
sexe	1	31.169	31.169	1.9211	0.1991
grade	2	50.041	25.021	1.5421	0.2655
sexe:grade	1	36.480	36.480	2.2484	0.1680
Residuals	9	146.024	16.225		

Troisième partie III

Développer avec R

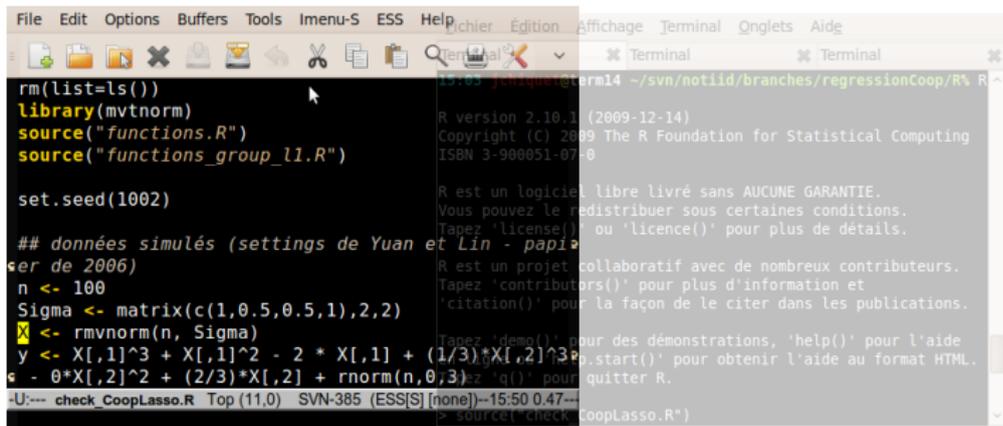
Structures de contrôle

Les fonctions

Les packages

Le module Sweave

Programmer en pratique : rappels



```
File Edit Options Buffers Tools Imenu-S ESS Help
15.05 jchiquet@term14 ~/svn/notiid/branches/regressionCoop/R% R ~
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
ou 'start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

-U:-- check_CoopLasso.R Top (11.0) SVN-385 (ESS[S][none])--15:50 0.47--
source("check_CoopLasso.R")

rm(list=ls())
library(mvtnorm)
source("fonctions.R")
source("fonctions_group_ll.R")

set.seed(1002)

## données simulés (settings de Yuan et Lin - papier de 2006)
n <- 100
Sigma <- matrix(c(1,0.5,0.5,1),2,2)
X <- rmvnorm(n, Sigma)
y <- X[,1]^3 + X[,1]^2 - 2 * X[,1] + (1/3)*X[,2]^3
e <- 0*X[,2]^2 + (2/3)*X[,2] + rnorm(n,0,3)
```

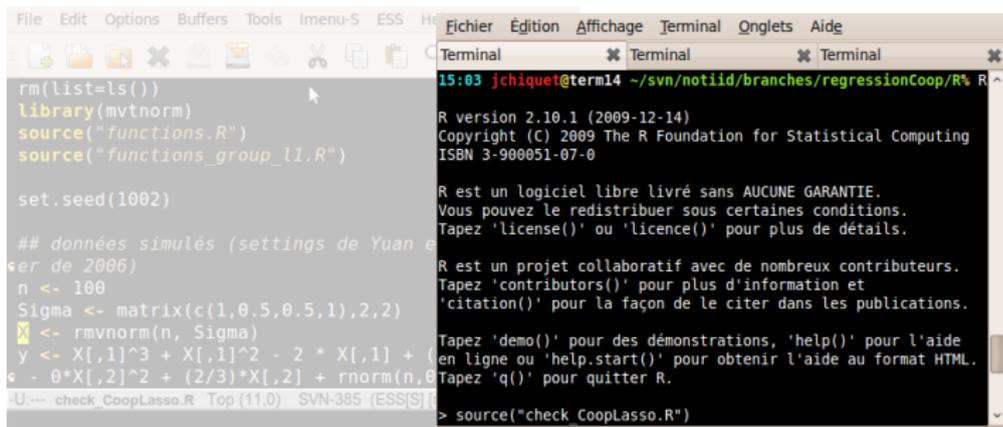
1. un éditeur de texte (vos fonctions / scripts)
2. un terminal avec R (tester, « sourcer »)

« Sourcer »

`source("un_script.R")` : exécute la série de commandes de `mon_script.R`

`source("mes_fonctions.R")` : charge le contenu (les fonctions) de `mes_fonctions.R`

Programmer en pratique : rappels



```
File Edit Options Buffers Tools Imenu-S ESS H...
rm(list=ls())
library(mvtnorm)
source("fonctions.R")
source("fonctions_group_l1.R")

set.seed(1002)

## données simulés (settings de Yuan et al. de 2006)
n <- 100
Sigma <- matrix(c(1,0.5,0.5,1),2,2)
X <- rmvnorm(n, Sigma)
y <- X[,1]^3 + X[,1]^2 - 2 * X[,1] + (
  - 0*X[,2]^2 + (2/3)*X[,2] + rnorm(n,0
-U--- check_CoopLasso.R Top (11.0) SVN-385 (ESS[S])
Fichier Édition Affichage Terminal Onglets Aide
Terminal Terminal Terminal
15:03 jchiquet@term14 ~/svn/notiid/branches/regressionCoop/R% R ^
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.
> source("check_CoopLasso.R")
```

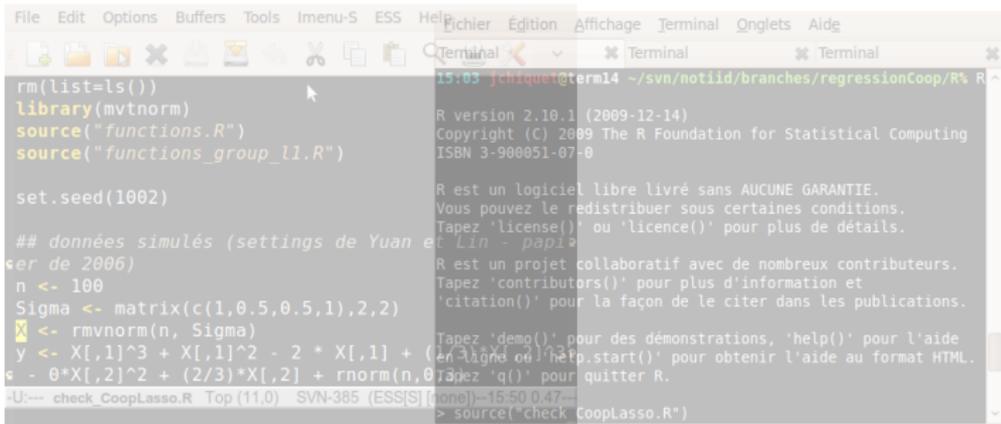
1. un éditeur de texte (vos fonctions / scripts)
2. un terminal avec R (tester, « sourcer »)

« Sourcer »

`source("un_script.R")` : exécute la série de commandes de `mon_script.R`

`source("mes_fonctions.R")` : charge le contenu (les fonctions) de `mes_fonctions.R`

Programmer en pratique : rappels



```
File Edit Options Buffers Tools Imenu-S ESS Help
Terminal
15:03 jchiquet@term14 ~/svn/notiid/branches/regressionCoop/R R ~
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.
en ligne où 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

rm(list=ls())
library(mvtnorm)
source("fonctions.R")
source("fonctions_group_l1.R")

set.seed(1002)

## données simulés (settings de Yuan et al. 2006)
n <- 100
Sigma <- matrix(c(1,0.5,0.5,1),2,2)
X <- rmvnorm(n, Sigma)
y <- X[,1]^3 + X[,1]^2 - 2 * X[,1] + (1-3)*X[,2]^2
e - 0*X[,2]^2 + (2/3)*X[,2] + rnorm(n,0,0.1)
-U-- check_CoopLasso.R Top (11.0) SVN-385 (ESS[S] [pane])--15:50 0.47--
> source("check_CoopLasso.R")
```

1. un éditeur de texte (vos fonctions / scripts)
2. un terminal avec R (tester, « sourcer »)

« Sourcer »

`source("un_script.R")` : exécute la série de commandes de `mon_script.R`

`source("mes_fonctions.R")` : charge le contenu (les fonctions) de `mes_fonctions.R`

Structures de contrôle

Les fonctions

Les packages

Le module Sweave

Syntaxe

```
{expr_1; expr_2; ...; expr_n }
```

ou

```
{  
  expr_1  
  ...  
  expr_n  
}
```

Remarques sur les groupes

- ▶ La dernière valeur du groupe est retournée ;
- ▶ un groupe d'expressions peut être passé à une fonction, réutilisé dans une expression plus grande, etc.

Syntaxe

```
if (condition) {  
  expr_1  
} else {  
  expr_2  
}
```

ou

```
ifelse(condition, a, b)
```

Remarques

- ▶ `condition` est une valeur logique : penser à `&`, `|`, `!`, ...
- ▶ le `else` est optionnel,
- ▶ `elseif` permet d'imbriquer les conditionnements.

Syntaxe

```
for (var in set) {  
  expr(var)  
}
```

ou

```
for (var in set)  
  expr(var)
```

à fuir pour éviter les effets de bords sournois!

Remarques sur la boucle `for`

- ▶ `var` est la variable incrémentée,
- ▶ `set` est un vecteur définissant les valeurs successives,
- ▶ *lente* comparée aux opérateurs matriciels.

Syntaxe

```
while (condition) {  
  expr  
}
```

ou

```
repeat {  
  expr  
}
```

Remarque

- ▶ Comme pour `for`, éviter les imbrications sources de lenteur.

Exemples d'utilisation

```
repeat {  
  expr  
  if (condition) {break}  
}
```

OU

```
while (condition1){  
  expr_1  
  if (condition2) {next}  
  expr_2  
}
```

Remarque

- ▶ `break` est la seule manière d'interrompre une boucle `repeat`.

Structures de contrôle

Les fonctions

Les packages

Le module Sweave

Syntaxe

```
nom_de_la_fonction <- function(arg1,arg2, ...) {  
  expression  
  
  return(var)  
}
```

Remarques

- ▶ `return` peut être omis (à éviter) : dans ce cas la dernière valeur calculée est renvoyée.
- ▶ peut être tapée directement dans l'interpréteur ou dans un fichier externe `fonctions.R`, chargé par `source`.

Moyenne empirique d'un vecteur

Avec suppression des valeurs manquantes !

```
> moyenne <- fonction(x) {  
+   ## suppression des valeurs manquantes  
+   x.not.na <- x[!is.na(x)]  
+   ## moyenne empirique  
+   resultat <- sum(x.not.na) / length(x.not.na)  
+  
+   return(resultat)  
+ }
```

Tests

```
> moyenne(rnorm(100))
```

```
[1] 0.06792586
```

```
> moyenne(c(1, -5, 3, NA, 8.7))
```

```
[1] 1.925
```

Propriétés

- ▶ les arguments peuvent être passés dans le **désordre** s'ils sont **nommés** : `var=object`,
- ▶ on peut définir une valeur par défaut pour n'importe quel argument lors de la définition de la fonction : `var=10`.
- ▶ en cas de **sorties multiples**, les sorties doivent être renvoyées sous forme de liste.

Remarques

- ▶ Les valeurs par défaut rendent la lecture des fonctions beaucoup plus aisée pour l'utilisateur : **imposer peu d'arguments obligatoires**.
- ▶ Les noms des éléments de la liste définie dans la fonction sont conservés à l'extérieur de la fonction.

Propriétés

- ▶ les arguments peuvent être passés dans le **désordre** s'ils sont **nommés** : `var=object`,
- ▶ on peut définir une valeur par défaut pour n'importe quel argument lors de la définition de la fonction : `var=10`.
- ▶ en cas de **sorties multiples**, les sorties doivent être renvoyées sous forme de liste.

Remarques

- ▶ Les valeurs par défaut rendent la lecture des fonctions beaucoup plus aisée pour l'utilisateur : **imposer peu d'arguments obligatoires**.
- ▶ Les noms des éléments de la liste définie dans la fonction sont conservés à l'extérieur de la fonction.

Propriétés

- ▶ les arguments peuvent être passés dans le **désordre** s'ils sont **nommés** : `var=object`,
- ▶ on peut définir une valeur par défaut pour n'importe quel argument lors de la définition de la fonction : `var=10`.
- ▶ en cas de **sorties multiples**, les sorties doivent être renvoyées sous forme de liste.

Remarques

- ▶ Les valeurs par défaut rendent la lecture des fonctions beaucoup plus aisée pour l'utilisateur : **imposer peu d'arguments obligatoires**.
- ▶ Les noms des éléments de la liste définie dans la fonction sont conservés à l'extérieur de la fonction.

Propriétés

- ▶ les arguments peuvent être passés dans le **désordre** s'ils sont **nommés** : `var=object`,
- ▶ on peut définir une valeur par défaut pour n'importe quel argument lors de la définition de la fonction : `var=10`.
- ▶ en cas de **sorties multiples**, les sorties doivent être renvoyées sous forme de liste.

Remarques

- ▶ Les valeurs par défaut rendent la lecture des fonctions beaucoup plus aisée pour l'utilisateur : **imposer peu d'arguments obligatoires**.
- ▶ Les noms des éléments de la liste définie dans la fonction sont conservés à l'extérieur de la fonction.

Résumé numérique d'un vecteur

```
> resume <- fonction(x, na.rm = TRUE, affiche = FALSE) {  
+   mu <- mean(x, na.rm = na.rm)  
+   s2 <- var(x, na.rm = na.rm)  
+   if (affiche) {  
+     cat("\nMoyenne:", mu, "et variance:", s2)  
+   }  
+   return(list(moyenne = mu, variance = s2))  
+ }
```

```
> out <- resume(rnorm(100))  
> out$variance
```

```
[1] 0.7163688
```

```
> out <- resume(affiche = TRUE, x = rexp(100, 0.5))
```

```
Moyenne: 1.896648 et variance: 4.912842
```

Structures de contrôle

Les fonctions

Les packages

Le module Sweave

Principe de Claerbout (Géophysicien, Stanford)

An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.

Démarche

1. Proposer une méthode et exposer dans un article ses propriétés,
2. Écrire et déposer un package sur CRAN,
3. Publier un article dans « journal of statistical software » ou une note dans « Bioinformatics ».

Principe de Claerbout (Géophysicien, Stanford)

An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.

Démarche

1. Proposer une méthode et exposer dans un article ses propriétés,
2. Écrire et déposer un package sur CRAN,
3. Publier un article dans « journal of statistical software » ou une note dans « Bioinformatics ».

Définir un objectif

Par exemple, SIMoNe : construire un graphe des interactions significatives entre gènes à partir de données du transcriptome.

Organisation type

1. Fichier DESCRIPTION
2. Répertoire R : fonctions R (fonction `inferGraph(data)`)
3. Répertoire man : documentation des fonctions
4. Répertoire data : données
5. (Répertoire src : pour les fichiers à compiler, header etc.)

Structures de contrôle

Les fonctions

Les packages

Le module Sweave

Produire un document scientifique, c'est

- ▶ **expérimenter**

(pipette + éprouvette + blouse) ou ordi \rightsquigarrow données

- ▶ **analyser les résultats**

méthode + logiciel + données \rightsquigarrow graphes

- ▶ **rédigier des observations**

idée + (bloc-note ou traitement de texte) \Rightarrow texte

- ▶ **mettre en forme**

texte + graphes + traitement de texte \Rightarrow article 

Produire un document scientifique, c'est

- ▶ **expérimenter**

(pipette + éprouvette + blouse) ou ordi \rightsquigarrow données

- ▶ analyser les résultats

méthode + logiciel + données \rightsquigarrow graphes

- ▶ rédiger des observations

idée + (bloc-note ou traitement de texte) \Rightarrow texte

- ▶ mettre en forme

texte + graphes + traitement de texte \Rightarrow article 

Produire un document scientifique, c'est

- ▶ **expérimenter**

(pipette + éprouvette + blouse) ou ordi \rightsquigarrow données

- ▶ **analyser les résultats**

méthode + logiciel + données \rightsquigarrow graphes

- ▶ rédiger des observations

idée + (bloc-note ou traitement de texte) \Rightarrow texte

- ▶ mettre en forme

texte + graphes + traitement de texte \Rightarrow article



Produire un document scientifique, c'est

- ▶ **expérimenter**

(pipette + éprouvette + blouse) ou ordi \rightsquigarrow données

- ▶ **analyser les résultats**

méthode + logiciel + données \rightsquigarrow graphes

- ▶ **rédigier des observations**

idée + (bloc-note ou traitement de texte) \Rightarrow texte

- ▶ **mettre en forme**

texte + graphes + traitement de texte \Rightarrow article



Produire un document scientifique, c'est

- ▶ **expérimenter**

(pipette + éprouvette + blouse) ou ordi \rightsquigarrow données

- ▶ **analyser les résultats**

méthode + logiciel + données \rightsquigarrow graphes

- ▶ **rédigier des observations**

idée + (bloc-note ou traitement de texte) \Rightarrow texte

- ▶ **mettre en forme**

texte + graphes + traitement de texte \Rightarrow article 

1.
 - ▶ analyser : MS Excel
 - ▶ rédiger : MS Word
 - ▶ mettre en forme : MS Word
2.
 - ▶ analyser : MatLab
 - ▶ rédiger : OpenOffice
 - ▶ mettre en forme : OpenOffice
3.
 - ▶ analyser : R
 - ▶ rédiger : Emacs
 - ▶ mettre en forme : L^AT_EX

mon opinion ^a : 

a. de geek très subjective

Le package `sweave` de L^AT_EX permet de faire appel à du code R dans le document

1.
 - ▶ analyser : MS Excel
 - ▶ rédiger : MS Word
 - ▶ mettre en forme : MS Word
2.
 - ▶ analyser : MatLab
 - ▶ rédiger : OpenOffice
 - ▶ mettre en forme : OpenOffice
3.
 - ▶ analyser : R
 - ▶ rédiger : Emacs
 - ▶ mettre en forme : \LaTeX

mon opinion ^a : ☹️

a. de geek très subjective

Le package `sweave` de \LaTeX permet de faire appel à du code R dans le document

- ▶ analyser : MS Excel
 - ▶ rédiger : MS Word
 - ▶ mettre en forme : MS Word
- ▶ analyser : MatLab
 - ▶ rédiger : OpenOffice
 - ▶ mettre en forme : OpenOffice
- ▶ analyser : R
 - ▶ rédiger : Emacs
 - ▶ mettre en forme : \LaTeX

mon opinion^a : 😊

a. de geek très subjective

Le package `sweave` de \LaTeX permet de faire appel à du code R dans le document

- ▶ analyser : MS Excel
 - ▶ rédiger : MS Word
 - ▶ mettre en forme : MS Word
- ▶ analyser : MatLab
 - ▶ rédiger : OpenOffice
 - ▶ mettre en forme : OpenOffice
- ▶ analyser : R
 - ▶ rédiger : Emacs
 - ▶ mettre en forme : \LaTeX

mon opinion^a : 😊

a. de geek très subjective

Le package `sweave` de \LaTeX permet de faire appel à du code R dans le document

Code latex

```
\documentclass[a4paper]{article}
```

```
\begin{document}
```

In this example we embed parts of the examples :

```
\texttt{kruskal.test} help page into a \LaTeX{} document: ...
```

```
\end{document}
```

Sortie pdf

In this example we embed parts of the examples from the kruskal.test help page into a LaTeX document : ...

```
\documentclass[a4paper]{article}
```

```
\begin{document}
```

In this example we embed parts of the examples from the `\texttt{kruskal.test}` help page into a `\LaTeX` document:

```
<<>>=
```

```
data(airquality)
```

```
kruskal.test(Ozone ~ Month, data = airquality)
```

```
@
```

which shows that the location parameter of the Ozone distribution varies significantly from month to month. Finally we include a boxplot of the data:

```
<<fig=TRUE,echo=FALSE>>=
```

```
boxplot(Ozone ~ Month, data = airquality)
```

```
@
```

```
\end{document}
```

```
\documentclass[a4paper]{article}
```

```
\usepackage{Sweave}
```

```
\begin{document}
```

In this example we embed parts of the examples from the `\texttt{kruskal.test}` help page into a \LaTeX document:

```
\begin{Sinput}
```

```
  R> data(airquality)
```

```
  R> kruskal.test(Ozone ~ Month, data = airquality)
```

```
\end{Sinput}
```

```
\begin{Soutput}
```

```
  Kruskal-Wallis rank sum test data: Ozone by Month Kruskal-Wallis  
  chi-squared = 29.2666, df = 4, p-value = 6.901e-06
```

```
\end{Soutput}
```

which shows that the location parameter of the Ozone distribution varies significantly from month to month.

Finally we include a boxplot of the data:

```
\includegraphics{example-1-002}
```

```
\end{document}
```

In this example we embed parts of the examples from the `kruskal.test` help page into a \LaTeX document :

```
> data(airquality)
> kruskal.test(Ozone ~ Month, data = airquality)
```

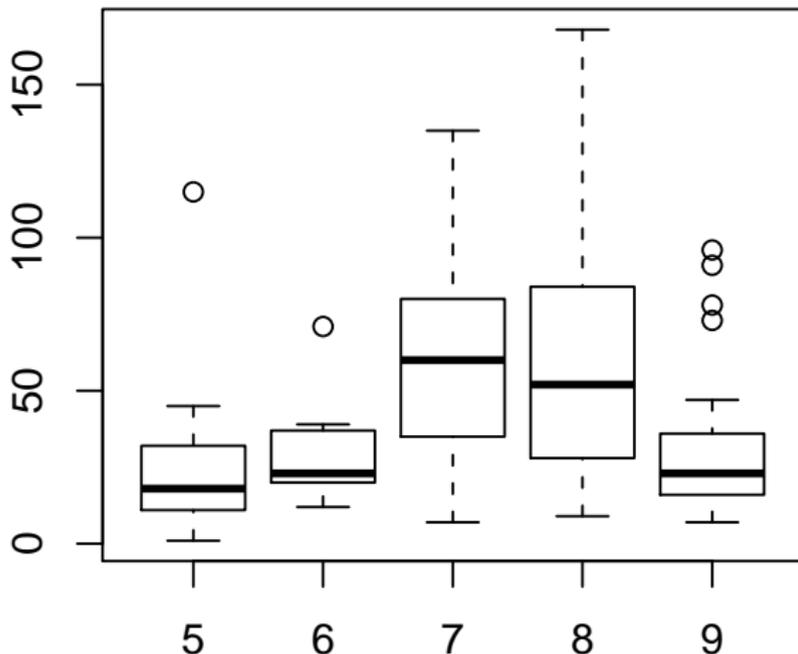
```
      Kruskal-Wallis rank sum test
```

```
data:  Ozone by Month
```

```
Kruskal-Wallis chi-squared = 29.2666, df = 4, p-value = 6.901
```

which shows that the location parameter of the Ozone distribution varies significantly from month to month. Finally we include a boxplot of the data :

Test de Kruskal



Quatrième partie IV

Entrées / Sorties

Charger des données

Les graphiques sous R

Charger des données

Les graphiques sous R

commande `scan`

Une utilisation élémentaire de `scan` permet une saisie plus agréable que la saisie directe des éléments d'un vecteur.

```
> x<-scan()
1: 1
2: 2
3: 3
4: 4
5: 5
6:
Read 5 items
>
> x
[1] 1 2 3 4 5
```

↪ valable pour les jeux de données d'au plus quelques dizaines d'éléments. . .

Éditer des données

commande `edit`

Permet d'éditer des données existantes à l'aide d'un mini-tableur. Utile pour faire de petites modifications.

```
> new.data <- edit(old.data)
```



delai.groupe
26 A
27 A
35 A
36 A
38 A
38 A
41 A

FIGURE: Éditeur Mac OS 10.6 / R 2.10

commandes `save` et `load`

La commande `save` permet de sauvegarder un sous ensemble des données de l'espace de travail dans un fichier binaire ; `load` permet de les recharger.

```
> x <- rnorm(125)
> y <- 1 + x + x^2
> save(file = "mes_simus", x, y)
> rm(list = ls())
> objects()
```

```
character(0)
```

```
> load(file = "mes_simus")
> objects()
```

```
[1] "x" "y"
```

commande data

R dispose d'une **collection de données prédéfinies** directement utilisables. La commande `data()` permet de les lister puis de les charger.

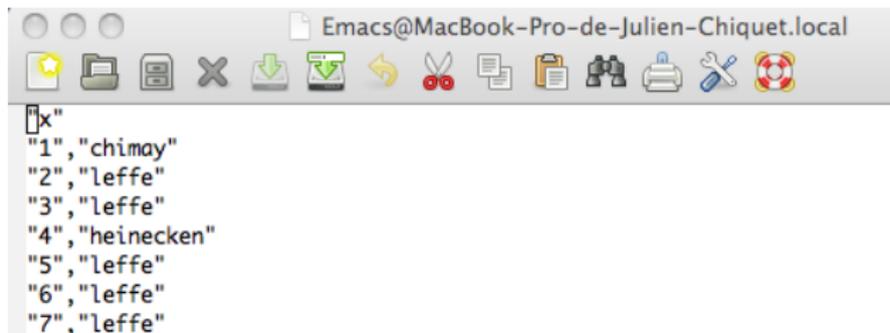
```
> data(iris3)
```

```
> head(iris3)
```

```
[1] 5.1 4.9 4.7 4.6 5.0 5.4
```

- ▶ La description d'un jeu de données est accessible dans l'aide.
- ▶ L'installation d'un nouveau package rend souvent disponibles de nouveaux jeux de données accessibles par `data`.

D'abord, un bon éditeur . . . il vous permet de constater le formatage d'un fichier texte et comment en « attaquer » l'importation.



```
Emacs@MacBook-Pro-de-Julien-Chiquet.local
"1", "chimay"
"2", "leffe"
"3", "leffe"
"4", "heinecken"
"5", "leffe"
"6", "leffe"
"7", "leffe"
```

FIGURE: Fichier au formatage "csv"

commande `read.table`

Elle permet de lire un fichier formaté sous forme de table.

`read.table` stocke les données sous forme d'objet `data.frame`.

```
> mes_donnees <- read.table("mesures_baie_raisin_2008-2009.txt",  
+   header = TRUE, sep = "\t")  
> head(mes_donnees)
```

	Population	variete	nbre.pepin.baie.2008	poids.pulpe.baie..g..2008	
1	CE	1784	1.0	0.89	
2	CE	124	1.0	1.14	
3	CE	210	1.2	1.26	
4	CE	1805	1.2	0.66	
5	CE	1303	1.2	0.83	
6	CE	284	1.3	0.54	
	volume.baie..cm3..2008	nbre.pepin.baie.2009	poids.pulpe.baie..g..2009		
1	7.70	NA	NA		NA
2	8.82	NA	NA		NA
3	10.20	NA	NA		NA
4	NA	NA	NA		NA
5	NA	NA	NA		NA
6	4.61	NA	NA		NA

commandes `read.csv` et `read.delim`

Ce sont des raccourcis pour la fonction `read.table`, spécialisés dans l'importation des données « `.csv` » (*comma-separated value*) ou tabulées (le séparateur est la tabulation).

commandes `write.table`, `write.csv` et `write.delim`

La fonction `write.table` permet d'imprimer les données issues d'une `data.frame` dans un fichier texte externe. `write.csv` et `write.delim` sont des raccourcis pour les données `csv` ou tabulée.

Beaucoup de choses sur l'importation des données dans



R Data Import /Export.

<http://cran.r-project.org/doc/manuals/R-data.pdf>

- ▶ Exemples avancés avec `read.table`,
- ▶ communication avec les bases de données (SQL),
- ▶ importation de données Excel,
- ▶ ...

Charger des données

Les graphiques sous R

Forme générique

La plupart des fonctions graphique s'utilisent par un appel du type

1. `nom.fonction(object, options),`
2. `nom.fonction(x, y , options).`

Parmi les options les plus courantes, on trouve :

- ▶ `type="p"` ; spécifie le type de tracé : "p" pour points, "l" pour lignes, "b" pour points liés par des lignes, "o" pour lignes superposées aux points. . .
- ▶ `xlim=` ; `ylim=`, spécifie les limites de axes x et y
- ▶ `xlab=` ; `ylab=`, annotation des axes x et y
- ▶ `main=` ; titre du graphe en cours
- ▶ `sub=` ; sous-titre du graphe en cours
- ▶ `add=FALSE` ; si TRUE superpose le graphe au précédent
- ▶ `axes=TRUE` ; si FALSE ne trace pas d'axes

Forme générique

La plupart des fonctions graphique s'utilisent par un appel du type

1. `nom.function(object, options),`
2. `nom.function(x, y , options).`

Parmi les options les plus courantes, on trouve :

- ▶ `type="p"` ; spécifie le type de tracé : "p" pour points, "l" pour lignes, "b" pour points liés par des lignes, "o" pour lignes superposées aux points. . .
- ▶ `xlim=` ; `ylim=`, spécifie les limites de axes x et y
- ▶ `xlab=` ; `ylab=`, annotation des axes x et y
- ▶ `main=` ; titre du graphe en cours
- ▶ `sub=` ; sous-titre du graphe en cours
- ▶ `add=FALSE` ; si TRUE superpose le graphe au précédent
- ▶ `axes=TRUE` ; si FALSE ne trace pas d'axes

commande `plot`

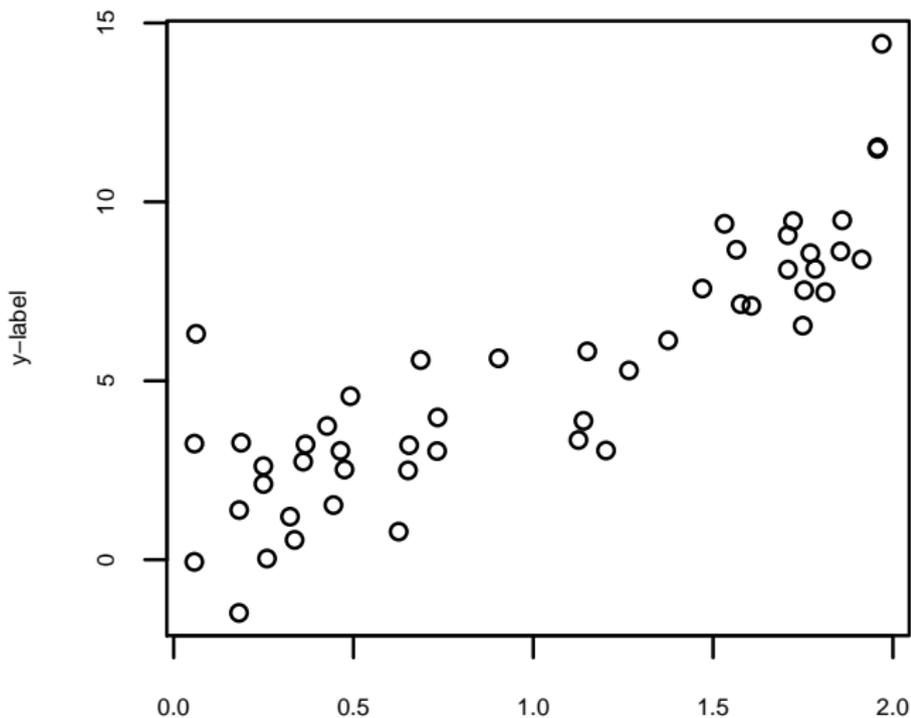
Fonction élémentaire de représentation graphique.

- ▶ `plot(vect)` représente le graphe des valeurs de `vect` sur l'axe des y .
- ▶ `plot(vect1, vect2)` représente le graphe des valeurs de `vect2` en fonction de `vect1`.

Par exemple, avec deux vecteurs :

```
> x <- runif(50, 0, 2)
> y <- 3 * x + 2 * x^2 + 1 + rnorm(50, sd = 1.5)
> plot(x, y, xlab = "x-label", ylab = "y-label", main = "mon premier graphe")
```

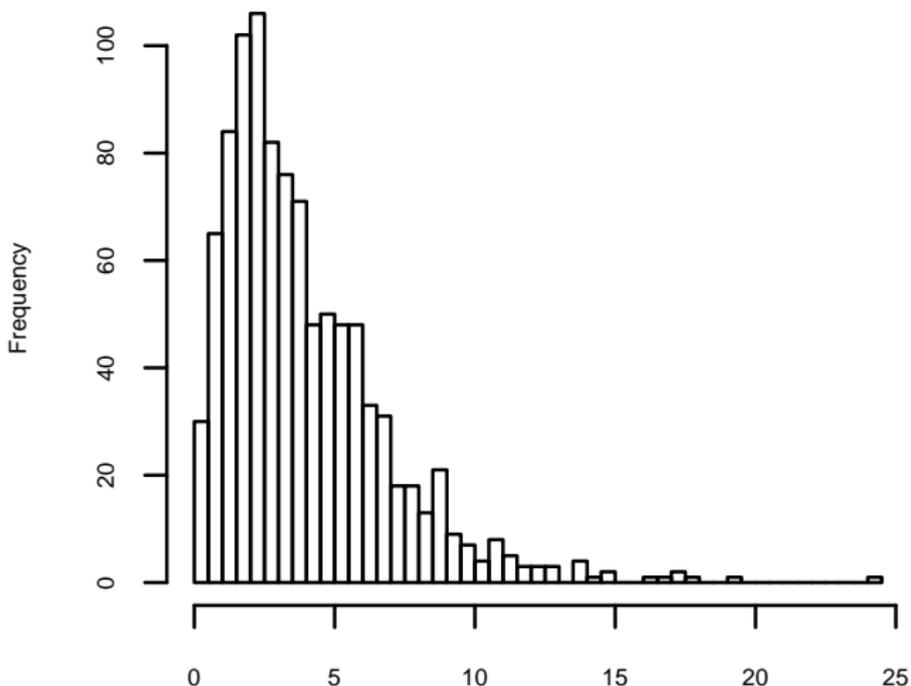
mon premier graphe



Beaucoup d'objet R accepte la commande `plot` ! En particulier, les histogrammes :

```
> mon_histo <- hist(rchisq(1000, df = 4), nclass = 75)
> plot(mon_histo, main = "distribution empirique du Khi-2")
```

Histogram of `rchisq(1000, df = 4)`

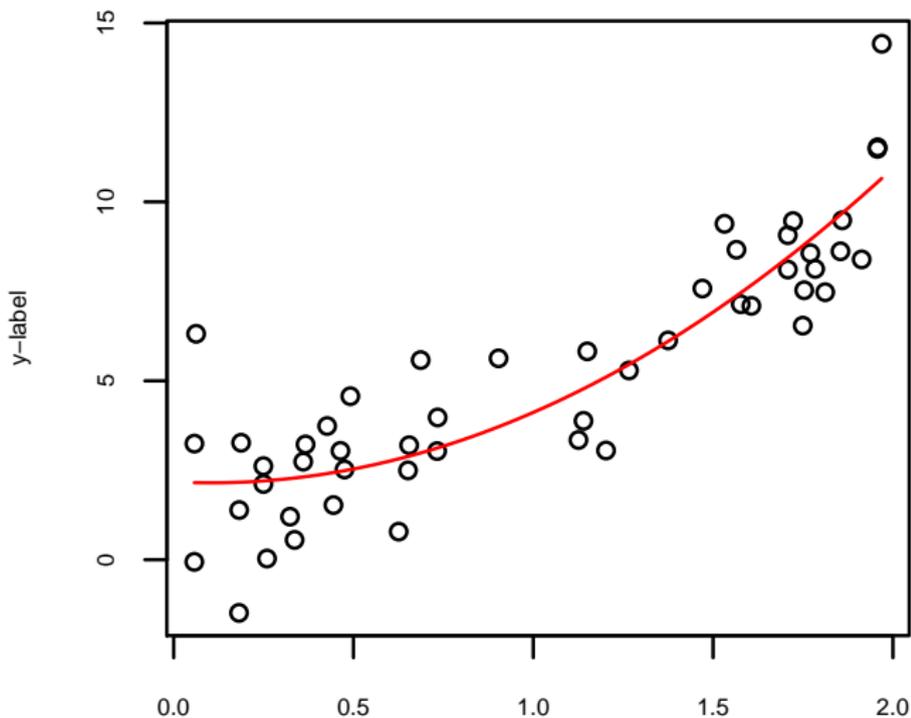


commande `curve`

Elle permet de tracer une fonction définie par une expression de x .

```
> plot(x, y, main = "données + modèle ajusté", xlab = "x-label",  
+       ylab = "y-label")  
> a <- coefficients(lm(y ~ 1 + x + I(x^2)))  
> curve(a[1] + a[2] * x + a[3] * x^2, add = TRUE, col = "red")
```

données + modèle ajusté

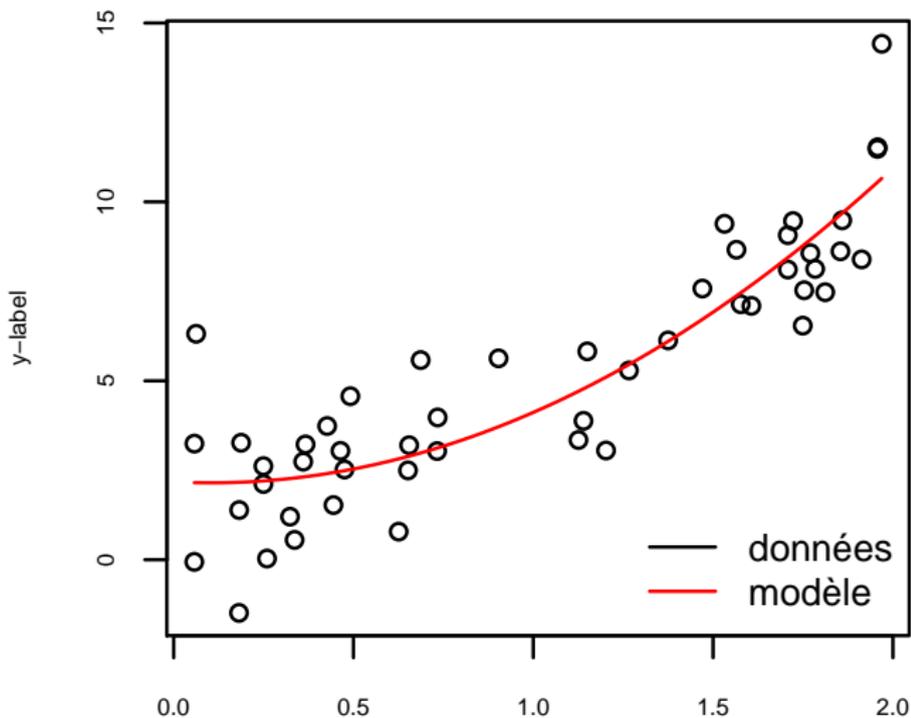


commande `legend`

Pour ajouter une légende. Attention aux options, assez nombreuses !

```
> plot(x, y, main = "données + modèle ajusté", xlab = "x-label",  
+      ylab = "y-label")  
> a <- coefficients(lm(y ~ 1 + x + I(x^2)))  
> curve(a[1] + a[2] * x + a[3] * x^2, add = TRUE, col = "red")  
> legend("bottomright", c("données", "modèle"), lty = c(1, 1),  
+      col = c("black", "red"), bty = "n")
```

données + modèle ajusté

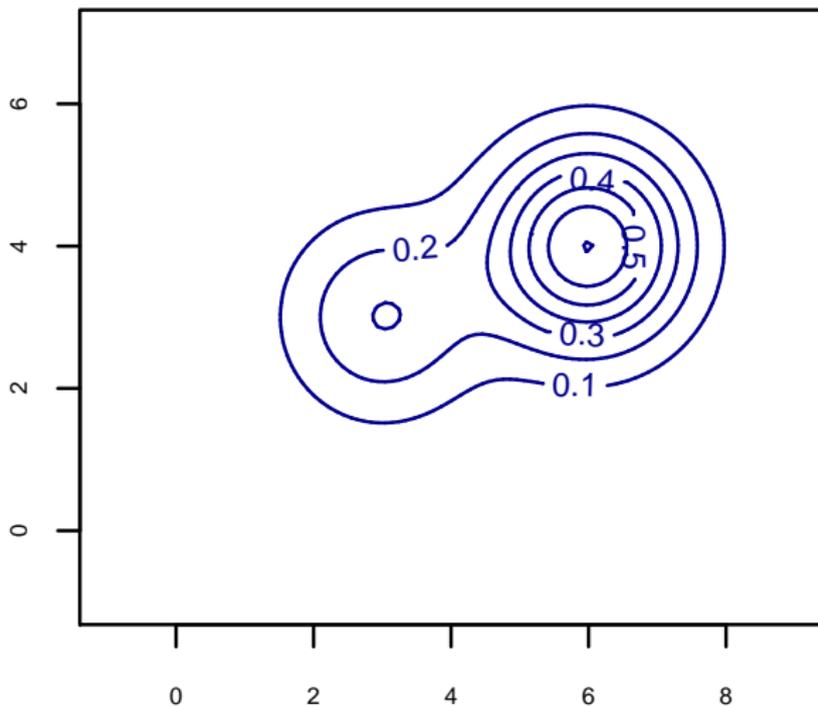


commande `contour`

`contour(x,y,z)` permet de tracer des courbes de niveaux : `x` et `y` sont des vecteurs et `z` une matrice telle que les dimensions de `z` soient `length(x)`, `length(y)`.

```
> x <- seq(-1, 9, length = 100)
> y <- seq(-1, 7, length = 100)
> z <- outer(x, y, function(x, y) 0.3 * exp(-0.5 * ((x - 3)^2 +
+ (y - 3)^2)) + 0.7 * exp(-0.5 * ((x - 6)^2 + (y - 4)^2)))
> contour(x, y, z, col = "blue4")
```

Représentation 3D (courbe de niveaux) II



commande `abline`

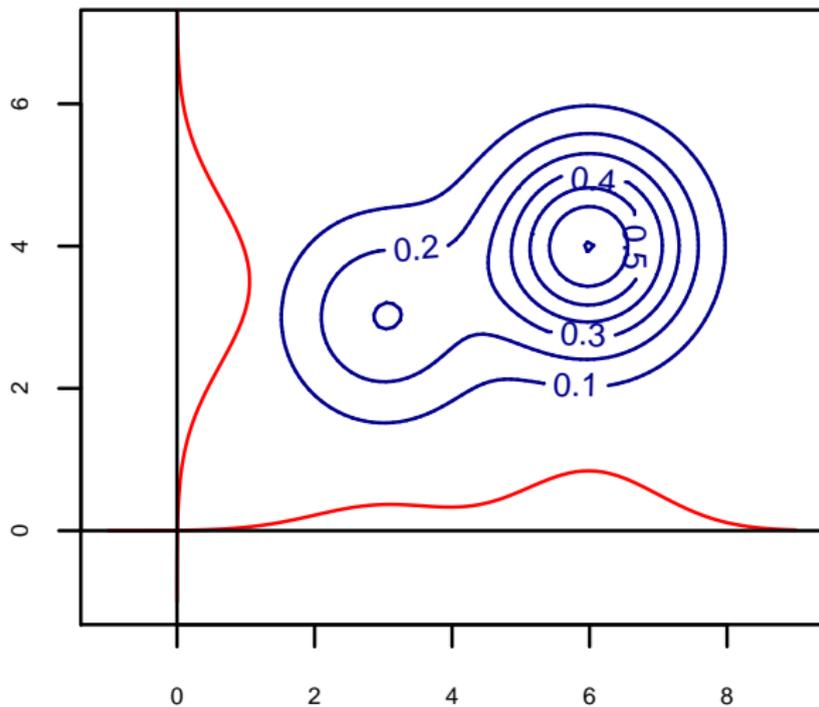
`abline` permet d'ajouter à un graphe courant

- ▶ des droites de décalage `a` et de coefficient directeur `b` avec `abline(a,b)`,
- ▶ des droites verticales avec `abline(v=)`,
- ▶ des droites horizontales avec `abline(h=)`.

commandes `lines` et `points`

Pour ajouter une courbe ou des points : s'utilisent de manière similaire à `plot`.

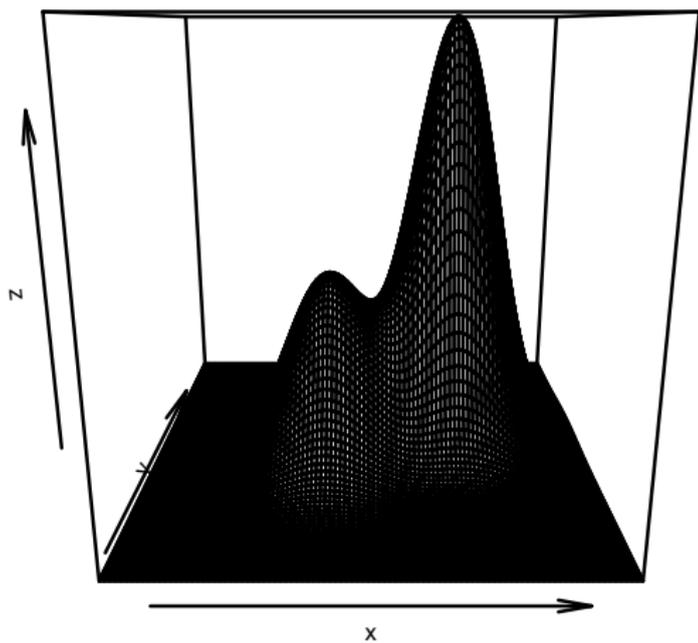
```
> contour(x, y, z, col = "blue4")
> curve((0.3 * dnorm(x, mean = 3) + 0.7 * dnorm(x, mean = 6)) *
+       3, -1, 9, col = "red", ylim = c(-1, 7), add = T)
> x <- seq(-1, 9, length = 100)
> lines((0.5 * dnorm(x, mean = 3) + 0.5 * dnorm(x, mean = 4)) *
+       3, x, col = "red")
> abline(h = 0)
> abline(v = 0)
```



commande `persp`

Fonctionne comme la fonction `contour` en proposant une représentation en perspective.

```
> persp(x, y, z)
```



Par défaut, R envoie les graphiques sur la sortie *écran*. De nombreuses

Exportation de graphes

Se réalise en encadrant les fonctions graphiques par les commandes `format_export(file="nom_fichier")` et `dev.off()`, où `format_fichier` peut prendre les valeurs `pdf`, `postscript`, `png`,

```
pdf(file="ma\_sortie.pdf")
plot(runif(20),runif(20))
dev.off()
```

Ouverture d'une nouvelle fenêtre graphique

Se fait, selon les plateformes, avec les commandes

- ▶ `x11()` pour Linux,
- ▶ `quartz()` ou `x11()` pour Mac OS,
- ▶ `windows()`.

Découpage d'une fenêtre

Plusieurs possibilités :

- ▶ `layout(mat,width=,height=)`, qui s'utilise en découpant l'écran via la matrice `mat`.
- ▶ `par(mfrow=vect)` OU `par(mfcol=vect)` qui découpent en n lignes et m colonne spécifiées par le vecteur `vect`. Le remplissage se fait par ligne ou par colonne selon la fonction choisie.

- ▶ D'autres fonctions de haut niveau dans la partie dédiée aux statistiques
- ▶ Utiliser la liste des commandes usuelles pour les options et fonctions secondaires,
- ▶ La commande `par` gère les options graphiques,
- ▶ Consulter le package `lattice`, *extrêmement* puissant.

 [Lattice : Multivariate Data Visualization with R](http://lmdvr.r-forge.r-project.org/)
Deepayan Sarkar
<http://lmdvr.r-forge.r-project.org/>

↪ Cette page web propose toutes les figures et tous les codes R correspondant à leur génération !

Cinquième partie V

Statistiques et outils connexes sous R

Statistiques descriptives

Les tests d'hypothèses

Introduction au modèle linéaire

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Statistiques descriptives

Tables statistiques

Analyse élémentaire d'une population

Les tests d'hypothèses

Principe du test d'hypothèse

Rappels de probabilités

Tests de Student

Test d'égalité des variances

Introduction au modèle linéaire

Les distributions disponibles

Distribution	R	Paramètres
beta	beta	
binomiale	binom	size, prob
binomiale négative	nbinom	
Cauchy	cauchy	
Chi-deux	chisq	df
Exponentielle	exp	rate
Fisher	f	df1, df2
Gamma	gamma	
géométrique	geom	
hypergéométrique	hyper	
log-normal	lnorm	
logistique	logis	
normale	norm	mean, sd
Poisson	pois	
Student	t	df
uniforme	unif	min, max
Weibull	weibull	
Wilcoxon	wilcox	

TABLE: Principales distributions

Les distributions disponibles

Distribution	R	Paramètres
beta	beta	
binomiale	binom	size, prob
binomiale négative	nbinom	
Cauchy	cauchy	
Chi-deux	chisq	df
Exponentielle	exp	rate
Fisher	f	df1, df2
Gamma	gamma	
géométrique	geom	
hypergéométrique	hyper	
log-normal	lnorm	
logistique	logis	
normale	norm	mean, sd
Poisson	pois	
Student	t	df
uniforme	unif	min, max
Weibull	weibull	
Wilcoxon	wilcox	

TABLE: Principales distributions

Les distributions disponibles

Distribution	R	Paramètres
beta	beta	
binomiale	binom	size, prob
binomiale négative	nbinom	
Cauchy	cauchy	
Chi-deux	chisq	df
Exponentielle	exp	rate
Fisher	f	df1, df2
Gamma	gamma	
géométrique	geom	
hypergéométrique	hyper	
log-normal	lnorm	
logistique	logis	
normale	norm	mean, sd
Poisson	pois	
Student	t	df
uniforme	unif	min, max
Weibull	weibull	
Wilcoxon	wilcox	

TABLE: Principales distributions

Forme générique : `r+distrib(n,...)`

`r` pour « random » : `n` donne la taille de l'échantillon et ... sont les paramètres requis selon la forme de `distrib`.

```
> rexp(10,rate=1/5)
```

```
[1] 1.8040438 8.5448108 8.4189741 10.3532942 1.0451114 2.5462225  
[7] 2.3830787 0.2378782 3.0679736 0.1143410
```

```
> rchisq(10,df=5)
```

```
[1] 8.9831561 4.2673919 1.3806250 5.3756070 0.8732194 2.2755234 6.7166650  
[8] 6.0161748 5.3812403 4.7332064
```

```
> runif(10,min=-2,max=2)
```

```
[1] -0.85684337 0.31422513 -0.33458314 1.36633366 0.06321932 1.24735060  
[7] 0.87420950 -1.18111914 1.10258613 1.88585268
```

```
> mean(rbinom(1000,10,prob=1/2))
```

```
[1] 5.013
```

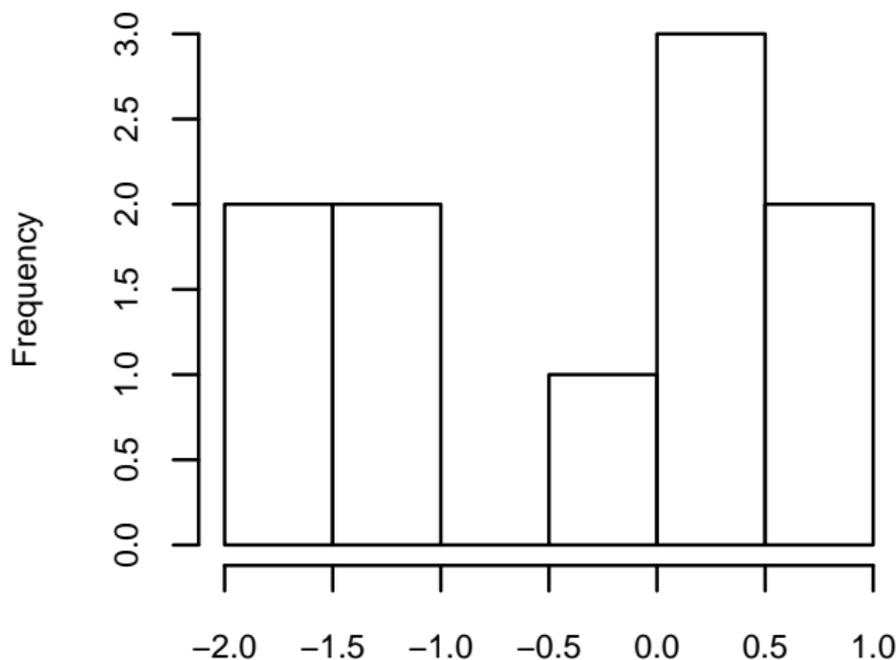
```
> var(rnorm(1000,mean=5,sd=2))
```

```
[1] 3.844386
```

Exemple avec la loi normale : histogramme

Avec $n = 10$

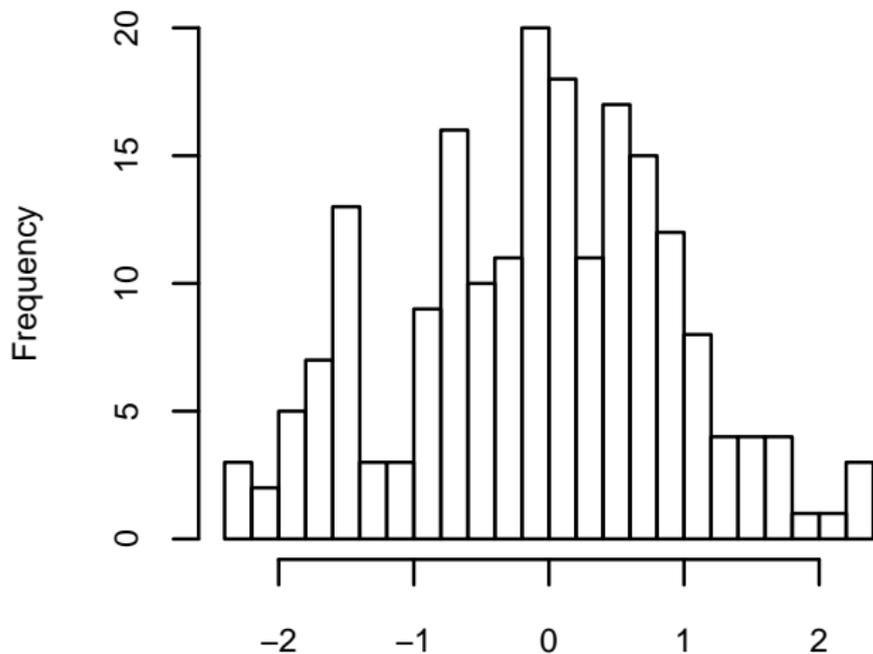
taille de l'échantillon = 10



Exemple avec la loi normale : histogramme

Avec $n = 200$

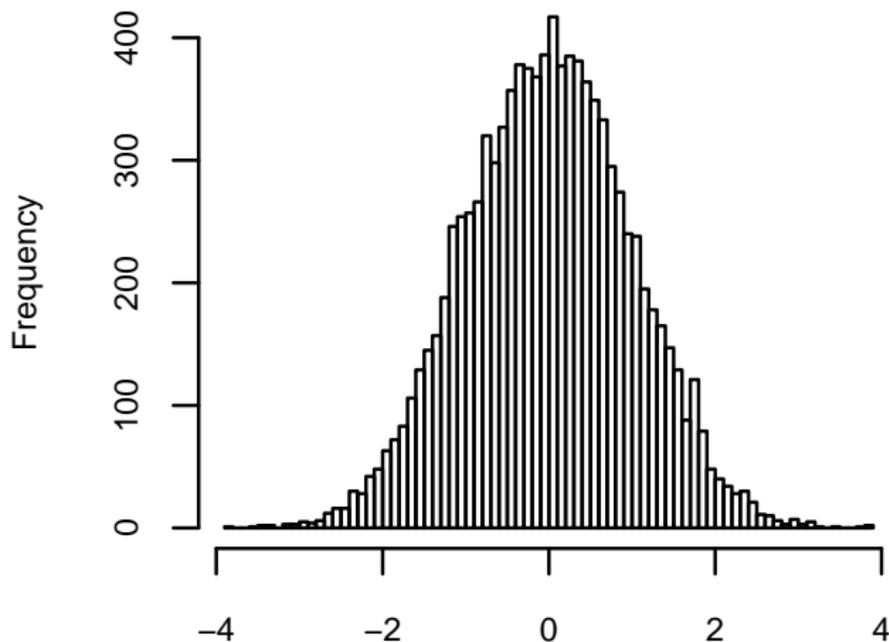
taille de l'échantillon = 200



Exemple avec la loi normale : histogramme

Avec $n = 10000$

taille de l'échantillon = 10000



Définir une distribution discrète

La fonction `sample(x, size, replace=FALSE, prob=NULL)` permet d'échantillonner les éléments de `x` : le tirage est de taille `size`, avec ou sans remise. Si `prob` est vide, chaque élément est équiprobable.

```
> sample(1:5)
```

```
[1] 5 4 2 3 1
```

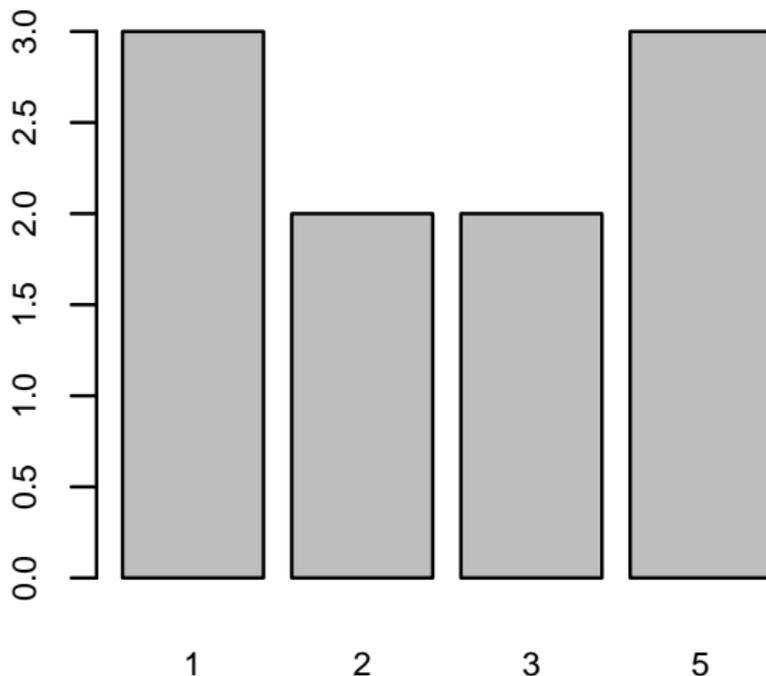
```
> sample(1:5,10,replace=TRUE)
```

```
[1] 4 3 3 1 5 4 4 4 3 4
```

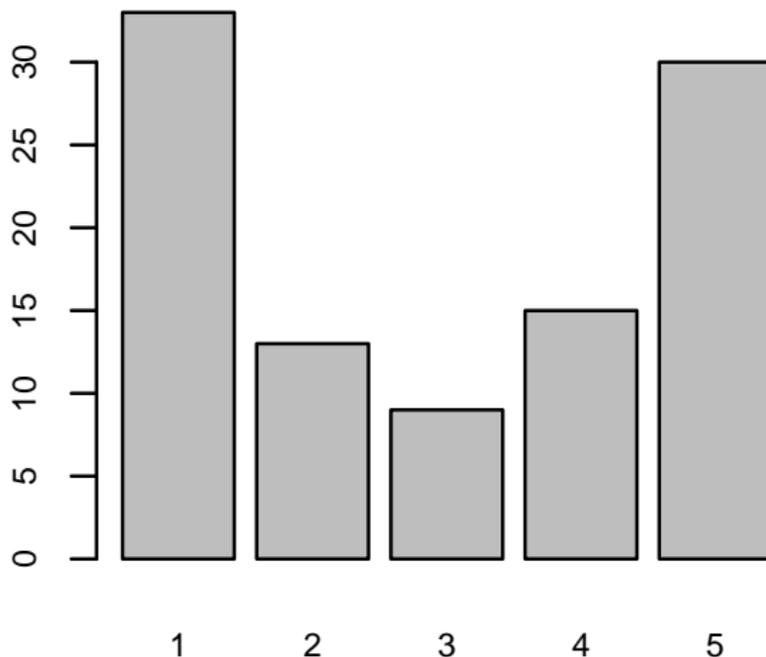
```
> sample(1:5,10,replace=TRUE,prob=c(.35,.1,.1,.1,0.35))
```

```
[1] 1 5 5 2 1 5 5 5 1 4
```

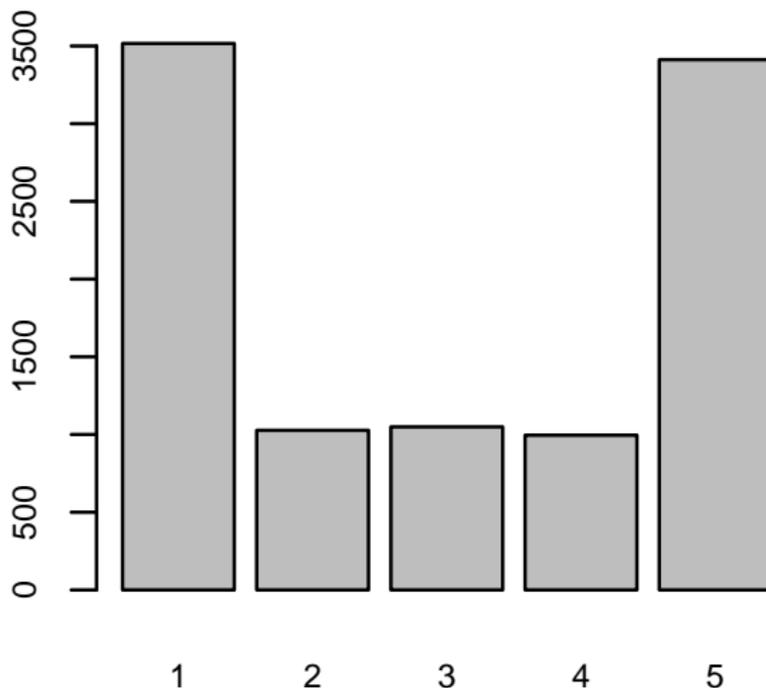
taille de l'échantillon = 10



taille de l'échantillon = 100



taille de l'échantillon = 10000



Forme générique : `p+distrib(x, ...)`

`p` pour « probability distribution function » : donne $\mathbb{P}(X \leq x)$, où X est une variable aléatoire de loi `distrib`.

```
> pnorm(0.5)
```

```
[1] 0.6914625
```

```
> pnorm(0.5, mean=2, sd=3)
```

```
[1] 0.3085375
```

```
> pnorm((0.5-2)/3)
```

```
[1] 0.3085375
```

```
> pbinom(5, 10, .25)
```

```
[1] 0.9802723
```

Forme générique : `d+distrib(x,...)`

`d` pour « density » : donne la densité pour une variable aléatoire continue et $\mathbb{P}(X = \mathbf{x})$ pour X une variable aléatoire discrète.

```
> dnorm(0.5)
```

```
[1] 0.3520653
```

```
> dexp(3,1/8)
```

```
[1] 0.08591116
```

```
> dbinom(5,10,.25)
```

```
[1] 0.0583992
```

```
> dpois(4,2)
```

```
[1] 0.09022352
```

Forme générique : `q+distrib(alpha, ...)`

`q` pour « quantile » : donne la valeur de x définie par

$$\mathbb{P}(X \leq x) = \alpha,$$

où X est une variable aléatoire de loi `distrib`.

```
> qnorm(0.95)
```

```
[1] 1.644854
```

```
> qt(0.4, df=28)
```

```
[1] -0.2557675
```

```
> qchisq(0.05, df=6)
```

```
[1] 1.635383
```

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Reprenons l'exemple de la vigne

Renommons les variables et considérons uniquement les données de 2008 pour une manipulation plus agréable. J'enlève également la colonne « variété », car je ne vois pas à quoi elle sert.

```
> vigne <- read.delim("mesures_baie_raisin_2008-2009.txt",header=TRUE)
> vigne <- vigne[-c(2,6,7)]
> colnames(vigne) <- c("pop","pepin.08","poids.08","volcm3.08")
> head(vigne)
```

```
  pop pepin.08 poids.08 volcm3.08
1  CE      1.0     0.89      7.70
2  CE      1.0     1.14      8.82
3  CE      1.2     1.26     10.20
4  CE      1.2     0.66      NA
5  CE      1.2     0.83      NA
6  CE      1.3     0.54      4.61
```

```
> attach(vigne)
```

Le résumé numérique s'adapte selon la nature des variables (univariée, multivariée, factorielle)

```
> summary(pepin.08)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.000	1.400	1.800	1.858	2.400	3.200	27.000

```
> summary(pop)
```

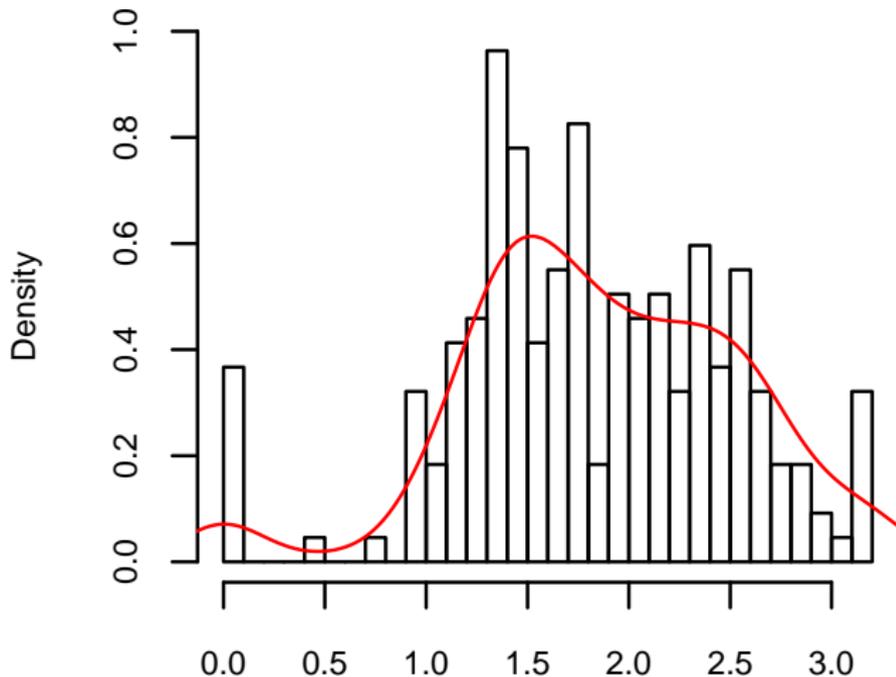
```
CE CO TE  
84 89 72
```

```
> summary(vigne)
```

pop	pepin.08	poids.08	volcm3.08
CE:84	Min. : 0.000	Min. : 0.390	Min. : 3.44
CO:89	1st Qu.: 1.400	1st Qu.: 0.820	1st Qu.: 7.14
TE:72	Median : 1.800	Median : 1.060	Median : 8.91
	Mean : 1.858	Mean : 1.212	Mean : 10.47
	3rd Qu.: 2.400	3rd Qu.: 1.360	3rd Qu.: 11.90
	Max. : 3.200	Max. : 3.750	Max. : 33.80
	NA's : 27.000	NA's : 28.000	NA's : 44.00

Histogramme

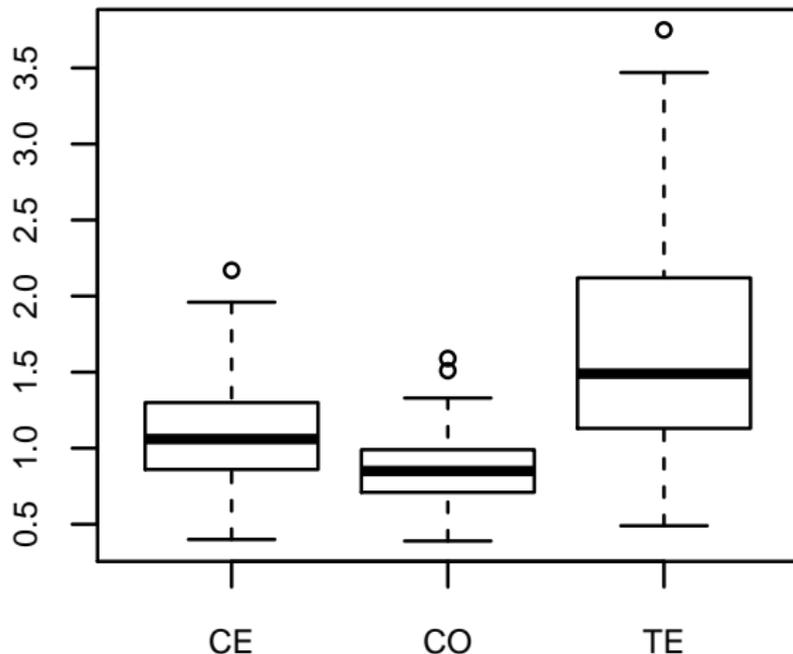
```
> hist(pepin.08,nclass=25,prob=TRUE)  
> lines(density(pepin.08[!is.na(pepin.08)]))
```



Boîtes à moustaches

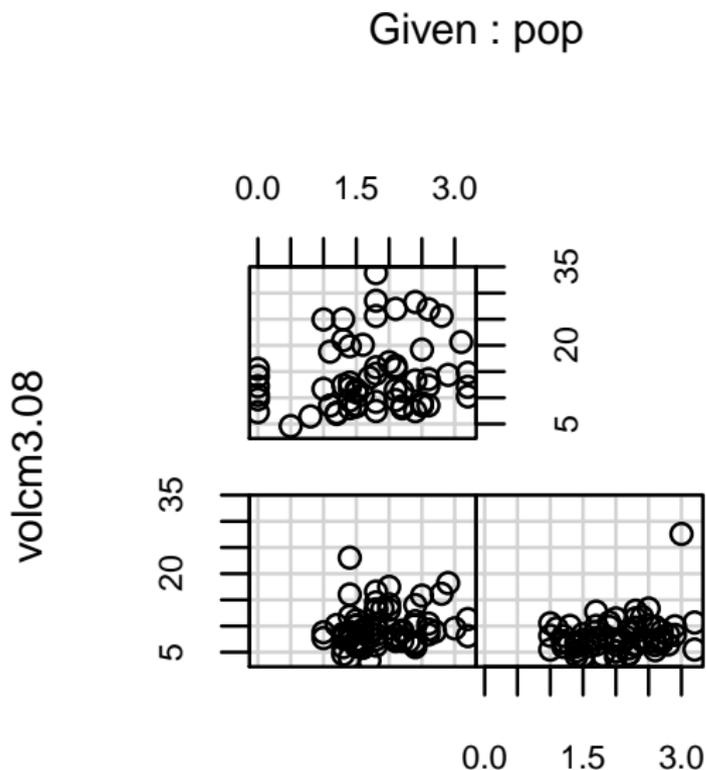
La boîte à moustache permet de visualiser les grands traits caractéristiques d'une distribution.

```
> boxplot(poids.08~pop)
```



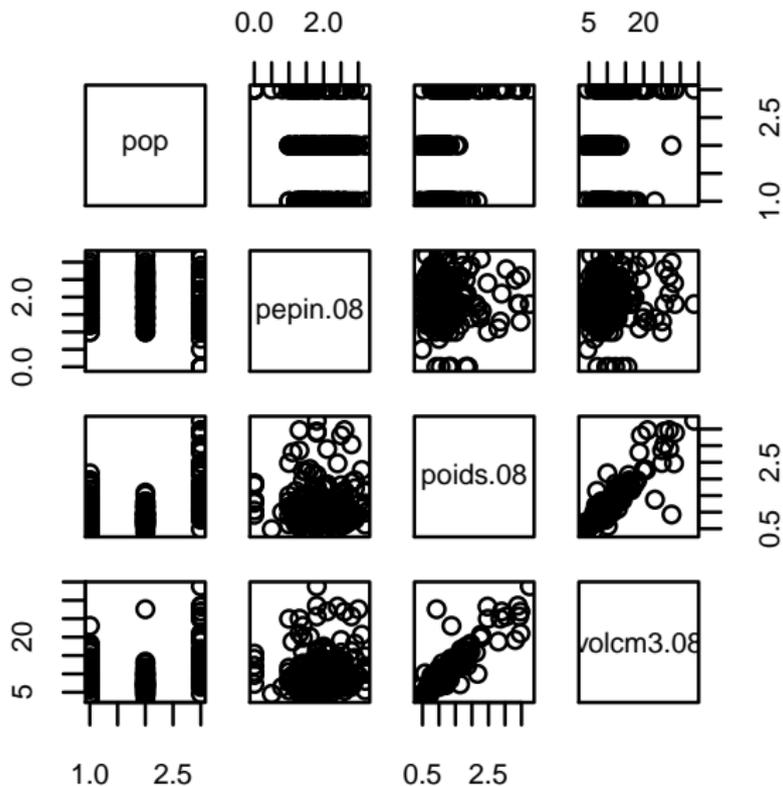
Graphe conditionné par une variable

```
> coplot(volcm3.08 ~ pepin.08 | pop, show.given=FALSE)
```



Graphes pair à pair

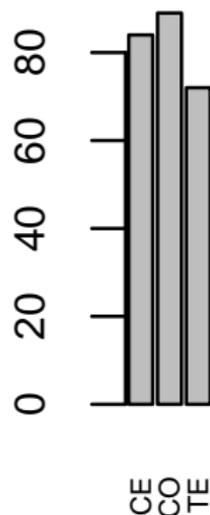
```
> pairs(vigne)
```



Camembert et diagramme en barres

Couplés à la commande `table` Le diagramme en barres et le graphe en camembert permettent de visualiser le découpage d'une population en donnée catégorielle.

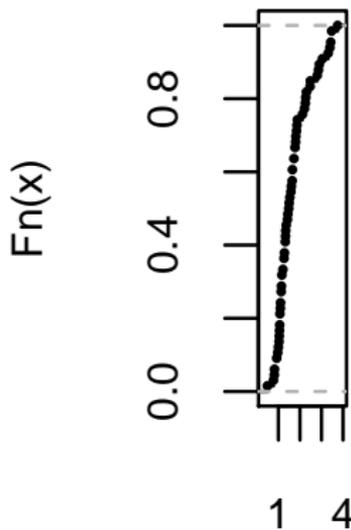
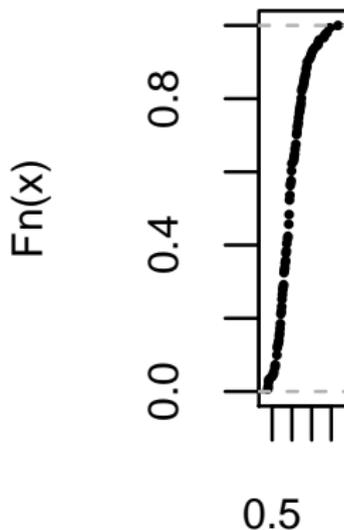
```
> par(mfrow=c(1,2))  
> pie(table(pop))  
> barplot(table(pop),las=3)
```



Fonction de répartition empirique

`ecdf` crée un objet qui peut être tracé avec `plot`.

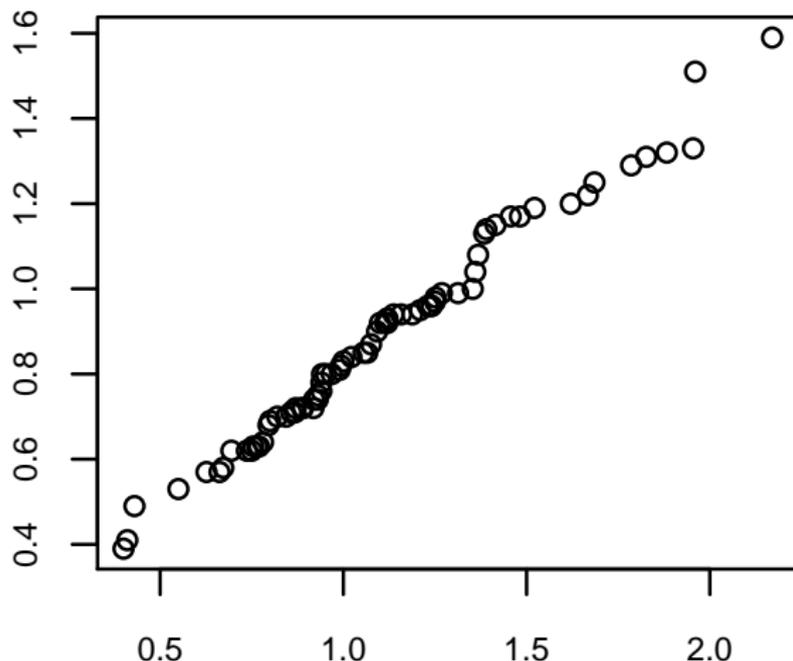
```
> par(mfrow=c(1,2))  
> plot(ecdf(poids.08[pop!="TE"]))  
> plot(ecdf(poids.08[pop=="TE"]))
```



Comparaison de distribution

Pour comparer visuellement deux distributions, la manière la plus efficace est le graphe quantile/quantile (qui doivent correspondre si les distributions sont proches.)

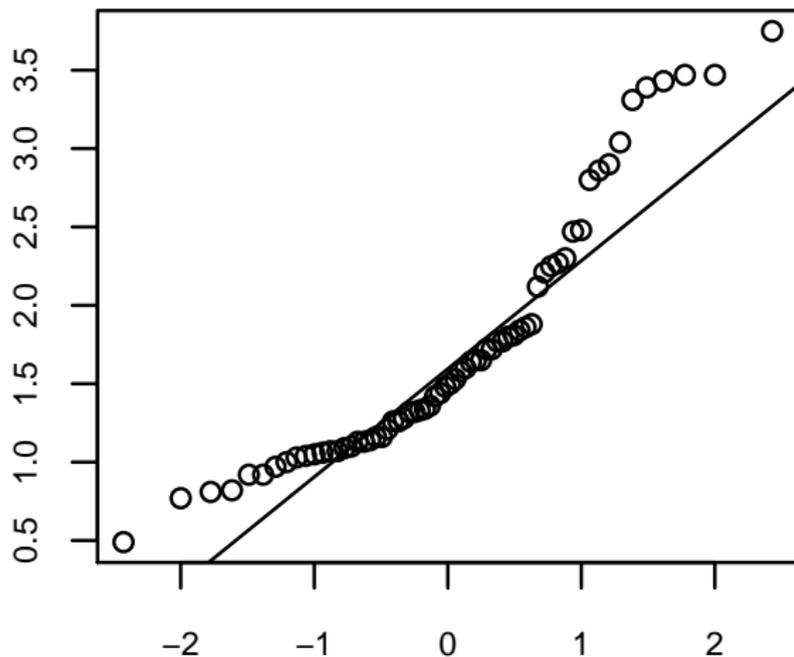
```
> qqplot(poids.08[pop=="CE"],poids.08[pop=="CO"])
```



Normalité d'une distribution

Une distribution est-elle normale ? Avant d'y répondre par un test, les commandes `qqnorm/qqline` donne une indication.

```
> qqnorm(poids.08[pop=="TE"])  
> qqline(poids.08[pop=="TE"])
```



Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Soit X un caractère d'étude dans une population \mathcal{P} .

Les ingrédients

- ▶ $X \sim \mathcal{F}$ (le **modèle**),
- ▶ n copies indépendantes (X_1, \dots, X_n) , (un **échantillon**),
- ▶ $X \sim \mathcal{F}_0$ (l'**hypothèse nulle**).

Le test

Décider, au vu d'une réalisation de l'échantillon, laquelle hypothèse est la plus possible parmi

$$\begin{cases} \mathcal{H}_0 : \mathcal{F}^* = \mathcal{F}_0, \\ \mathcal{H}_1 : \mathcal{F}^* = \mathcal{F}_1 \neq \mathcal{F}_0. \end{cases}$$

Ingrédients

- ▶ **erreurs** de type 1 et 2 : décider $\mathcal{H}_0|\mathcal{H}_1 / \mathcal{H}_1|\mathcal{H}_0$,
- ▶ **risques** α et β (probabilités d'erreur),
- ▶ **statistique de test** T_n (différente sous \mathcal{H}_0 et \mathcal{H}_1).

Stratégie

- ▶ **Contrôler** les faux positifs en fixant $\alpha = \alpha^*$,
- ▶ Construire une **règle de décision** dépendant de T_n et α^*

Rejet de \mathcal{H}_0 si $T_n \in \text{zone de rejet}(\alpha^*)$.

Degré de significativité

C'est le plus petit niveau α tel que le test est rejeté. Pour un test unilatéral, par exemple

$$\text{p-value} = \mathbb{P}_{\mathcal{H}_0}(T_n > t_{\text{obs}}),$$

c'est donc le risque commis en décidant \mathcal{H}_1 d'après les observations.

La puissance

C'est la probabilité de rejeter \mathcal{H}_0 à raison

$$\pi = 1 - \beta = \mathbb{P}_{\mathcal{H}_1}(\text{décider } \mathcal{H}_1),$$

c'est-à-dire les vrais positifs.

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités**

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Definition

On appelle loi du Khi-deux χ_n^2 à n degrés de libertés (ddl) la v.a. Y définie par

$$Y = U_1^2 + \dots + U_n^2, \quad \text{où } U_i \sim \mathcal{N}(0, 1).$$

Definition

On appelle loi de Student \mathcal{T}_n à n ddl la v.a. Z vérifiant

$$Z = \frac{U}{\sqrt{Y/n}}, \quad \text{où } U \sim \mathcal{N}(0, 1) \text{ et } Y \sim \chi_n^2.$$

Definition

On appelle loi de Fisher $\mathcal{F}_{n,m}$ à n et m ddl la v.a. F définie par

$$F = \frac{X/n}{Y/m}, \quad \text{où } X \sim \chi_n^2, Y \sim \chi_m^2.$$

Théorème

Soit $(X_i)_{i=1,\dots,n}$ une suite de variables aléatoires indépendantes de même loi et soit S_n la somme de ces variables. Alors,

$$\frac{S_n}{n} \xrightarrow[n \rightarrow \infty]{p.s.} \mathbb{E}[X].$$

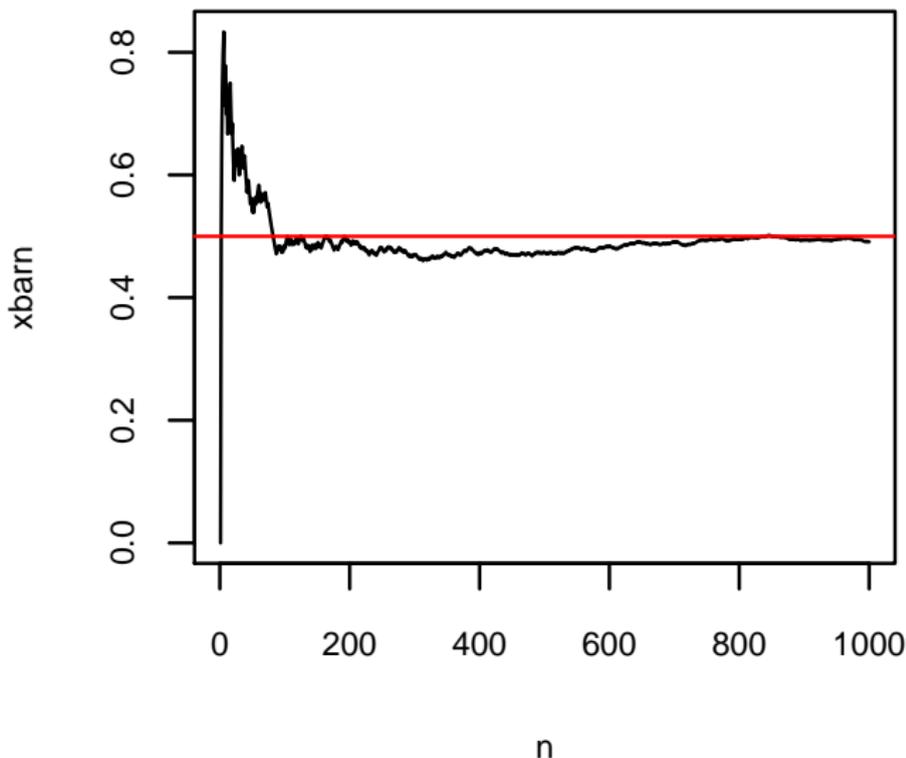
Conséquence théorique

La moyenne empirique $\bar{X} = S_n/n$ est un estimateur fortement convergent de l'espérance d'une loi.

Conséquence pratique

C'est le **fondement** de la plupart des simulations numériques en statistiques.

L'exemple passe-partout (mais parlant) : simuler l'issue d'un tirage de pile ou face, et observer l'évolution.



Théorème

Soit $(X_i)_{i=1,\dots,n}$ une suite de variable aléatoire indépendante de même loi, d'espérance μ et de variance σ^2 . Alors,

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

Conséquences

- ▶ L'écart à la moyenne suit une loi normale lorsque l'on observe beaucoup d'individus, quelque soit la loi de la variable observée.
- ▶ C'est un résultat portant sur un indice **global** de la population **pas** sur les individus!!!

Protocole

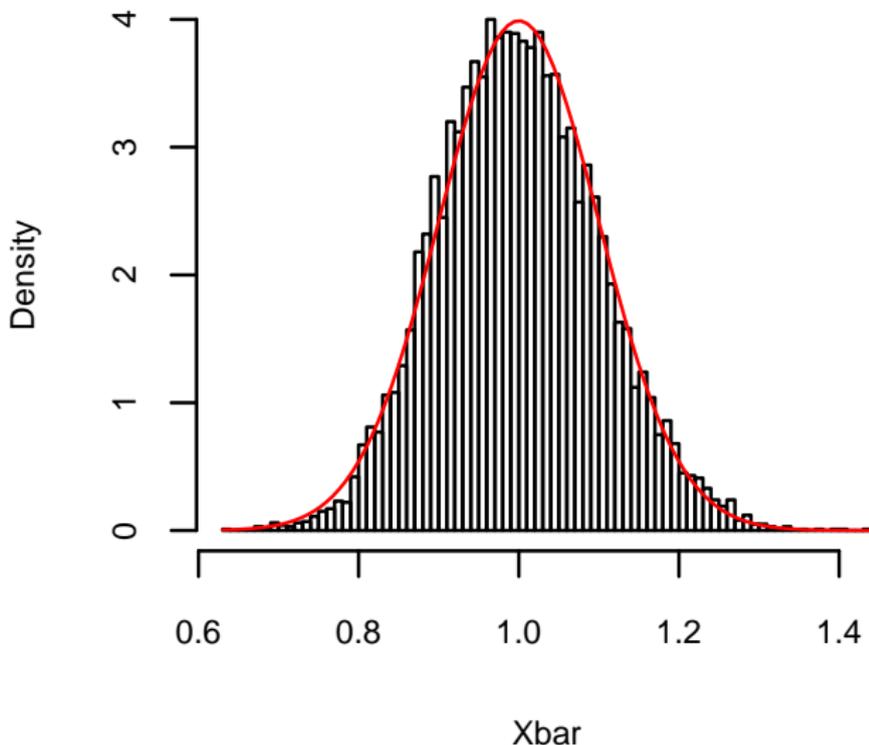
- ▶ On simule n réalisations d'une variable exponentielle,
- ▶ On calcule la valeur observés de la moyenne empirique,
- ▶ On répète m fois cette expérience,
- ▶ On trace l'histogramme des m observations de \bar{X}_n .

En R (compact !), ca donne :

```
> n <- 100  
> m <- 10000  
> Xi <- matrix(rexp(n*m),n,m)  
> Xbar <- colMeans(Xi)  
> hist(Xbar,nclass=100)
```

Illustration : graphique

```
> hist(Xbar, nclass=100, prob=TRUE, main="")  
> curve(dnorm(x, mean=1, sd=1/sqrt(n)), add=TRUE, col="red")
```



Théorème

Soit $(X_i)_{i=1,\dots,n}$ une série de variables aléatoires indépendantes, de même loi, d'espérance μ et de variance σ^2 . Alors, lorsque $n \rightarrow \infty$,

$$\sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \sigma^2 \chi_{n-1}^2,$$

ce qu'on écrit également

$$n\hat{\sigma}^2 \sim \sigma^2 \chi_{n-1}^2,$$

où $\hat{\sigma}^2$ est la variance empirique.

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student**

- Test d'égalité des variances

Introduction au modèle linéaire

Problème

Soit X un caractère d'intérêt d'une population \mathcal{P} de loi, d'espérance μ et de variance σ^2 inconnue. On propose une valeur μ_0 de l'espérance, et on teste sa validité en observant un échantillon i.i.d $(X_i)_{i=1,\dots,n}$:

$$\begin{cases} H_0 : \mu = \mu_0, \\ H_1 : \mu = \mu_1, \text{ avec } \mu_1 \neq \mu_0. \end{cases}$$

Par le théorème central limite et le théorème de Cochran, on montre que

$$\frac{\bar{X}_n - \mu}{\hat{\sigma}/\sqrt{n}} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{T}(n-1),$$

où $\mathcal{T}(n-1)$ est une loi de Student à $n-1$ degrés de liberté.

Si on teste

$$\begin{cases} H_0 : \mu = \mu_0, \\ H_1 : \mu = \mu_1, \text{ avec } \mu_1 > \mu_0, \end{cases}$$

alors, au niveau de confiance α ,

$$\text{on rejette } H_0 \text{ si } \bar{x}_n > \mu_0 + \frac{\hat{\sigma}}{\sqrt{n}} u_{1-\alpha}.$$

Définition

La p -valeur ou degré de significativité donne le risque que l'on prend à choisir H_1 plutôt que H_0 . Si elle est supérieure à α , on conserve H_0 et on rejette sinon.

Exemple

(Sirop contre la toux) La notice d'un sirop contre la toux indique comme valeur de référence pour la moyenne m_0 de l'agent actif 40g/litre. Le contrôleur de la fabrication décidera d'arrêter provisoirement la production si la moyenne m inconnue est strictement inférieure à cette valeur de référence.

Le contrôleur de la fabrication prélève de manière indépendantes 9 bouteilles au hasard dans la production et mesure la quantité d'agent actif. Conclusion ?

```
> x <- c(38.7, 39.6, 37.9, 40.6, 40.5, 37.7, 41.2, 37.5, 39.1)
> t.test(x,mu=40,alternative="less")
```

One Sample t-test

```
data: x
t = -1.7586, df = 8, p-value = 0.05835
alternative hypothesis: true mean is less than 40
95 percent confidence interval:
 -Inf 40.04593
sample estimates:
mean of x
 39.2
```

Problème

Soit X et Y deux caractères d'intérêt d'une population \mathcal{P} de loi, d'espérance μ_1 et μ_2 et de variance commune σ^2 , toutes inconnues. On propose de tester l'égalité des espérances en observant deux échantillons i.i.d $(X_i)_{i=1,\dots,n}$ et $(Y_i)_{i=1,\dots,m}$:

$$\begin{cases} H_0 : \mu_1 = \mu_2, \\ H_1 : \mu_1 \neq \mu_2, \end{cases}$$

Par le théorème central limite et le théorème de Cochran, on montre que

$$\frac{\bar{X}_n + \bar{Y}_m - (\mu_1 + \mu_2)}{s^* \sqrt{n^{-1} + m^{-1}}} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{T}(n + m - 2),$$

$$\text{où } s^{*2} = \frac{(n-1)S_X^{*2} + (m-1)S_Y^{*2}}{n+m-2}$$

Testons l'égalité des poids entre espèces (CE,CO) et TE chez la vigne :

```
> t.test(poids.08[pop!="TE"],poids.08[pop=="TE"])
```

```
Welch Two Sample t-test
```

```
data: poids.08[pop != "TE"] and poids.08[pop == "TE"]  
t = -7.1015, df = 75.698, p-value = 5.751e-10  
alternative hypothesis: true difference in means is not equal to 0  
95 percent confidence interval:  
 -0.9142257 -0.5137193  
sample estimates:  
mean of x mean of y  
0.9949669 1.7089394
```

↪ Existe en version apparié sous R, avec ou sans variances égales (avec `paired=TRUE`, `var.equal=TRUE`).

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Problème

Soit X et Y deux caractères d'intérêt d'une population \mathcal{P} de lois et de variances σ_1^2 et σ_2^2 inconnues. On propose de tester l'égalité des variances en observant deux échantillons i.i.d $(X_i)_{i=1,\dots,n}$ et $(Y_i)_{i=1,\dots,m}$:

$$\begin{cases} H_0 : \sigma_1 = \sigma_2, \\ H_1 : \sigma_1 \neq \sigma_2, \end{cases}$$

Par le théorème central limite et le théorème de Cochran, on montre que

$$\frac{S_X^{*2}/\sigma_1^2}{S_Y^{*2}/\sigma_2^2} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{F}(n-1, m-1).$$

Testons l'égalité des variances des poids entre espèces (CE,CO) et TE chez la vigne :

```
> var.test(poids.08[pop!="TE"],poids.08[pop=="TE"])
```

```
      F test to compare two variances
```

```
data:  poids.08[pop != "TE"] and poids.08[pop == "TE"]
```

```
F = 0.1846, num df = 150, denom df = 65, p-value < 2.2e-16
```

```
alternative hypothesis: true ratio of variances is not equal
```

```
95 percent confidence interval:
```

```
 0.1198802 0.2746481
```

```
sample estimates:
```

```
ratio of variances
```

```
 0.1845923
```

- ▶ `shapiro.test()` : gaussianité d'une population ;
- ▶ `wilcoxon.test` : test sur la médiane ou de comparaison de médianes (populations non-gaussienne - existe en version apparié) ;
- ▶ `prop.test` : test de proportion sur le paramètre d'une loi de Bernoulli,
- ▶ `chisq.test()` : test du χ^2 d'adéquation ou d'indépendance ;
- ▶ `ks.test()` : test de Kolmogorov-Smirnov d'adéquation d'une fonction de répartition ;
- ▶ `kruskal.test()` : test de Kruskal-Wallis d'effet d'un facteur pour une population non-gaussienne ;
- ▶ et beaucoup d'autres `help.search("test")`.

Statistiques descriptives

- Tables statistiques

- Analyse élémentaire d'une population

Les tests d'hypothèses

- Principe du test d'hypothèse

- Rappels de probabilités

- Tests de Student

- Test d'égalité des variances

Introduction au modèle linéaire

Formulation générale

$$y = \beta_0 + \sum_{i=1}^p X_i \beta_i + \varepsilon,$$

où

- ▶ y est la variable à expliquer ou *réponse*,
- ▶ $(X_i)_{i=1, \dots, p}$ sont les variable explicatives ou *prédicteurs*,
- ▶ β_0 est le biais (à estimer)
- ▶ $(\beta_i)_{i=0, \dots, p}$ sont les paramètres à estimer,
- ▶ $\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$ sont les résidus (ce que le modèle n'explique pas)

Remarque

le modèle est linéaire au sens *d'une combinaison linéaire des paramètres* ! Les modèles suivants sont des modèles linéaires :

1. Polynomial

$$y = \beta_0 + X_1\beta_1 + X_2^2\beta_2 + \varepsilon$$

2. Cobb-Douglas

$$\ln y = \beta_0 + \beta_1 \ln X_1 + \varepsilon$$

3. Logistique

$$\ln \frac{y}{1-y} = \beta_0 + \beta_1 X_1 + \varepsilon$$

↪ Un des modèles les plus puissants et les plus utilisés de la statistique

La fonction `lm`

Pour « linear model » : s'utilise avec le concept de *formule* :

- ▶ 1 variable explicative avec intercept (régression linéaire)

$$\text{lm}(y \sim X + 1),$$

- ▶ 2 variables explicatives sans intercept

$$\text{lm}(y \sim X1 + X2 - 1),$$

- ▶ 1 variable factorielle (anova à 1 facteur)

$$\text{lm}(y \sim A),$$

- ▶ 2 variables factorielles (anova à 2 facteurs sans interactions)

$$\text{lm}(y \sim A+B),$$

- ▶ 1 variable factorielles (anova à 2 facteurs avec interactions)

$$\text{lm}(y \sim A*B).$$

Exemple

On veut expliquer la variable de poids par la population d'origine et le volume de la baie. Plusieurs modèles sont possibles, par exemple

1. sans interactions

```
> m1 <- lm(poids.08 ~ volcm3.08 + pop)
> anova(m1)
```

Analysis of Variance Table

Response: poids.08

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
volcm3.08	1	62.069	62.069	903.80	< 2.2e-16	***
pop	2	1.766	0.883	12.86	5.632e-06	***
Residuals	197	13.529	0.069			

Signif. codes: 0

2. avec interactions

```
> m2 <- lm(poids.08 ~ volcm3.08 * pop)
> anova(m2)
```

Analysis of Variance Table

Response: poids.08

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
volcm3.08	1	62.069	62.069	1088.990	< 2.2e-16	***
pop	2	1.766	0.883	15.495	5.684e-07	***
volcm3.08:pop	2	2.415	1.207	21.183	4.733e-09	***
Residuals	195	11.114	0.057			

Signif. codes: 0

Comparer plusieurs modèles

La fonction `anova(modele1, modele2, ...)`.

```
> anova(m1, m2)
```

Analysis of Variance Table

Model 1: `poids.08 ~ volcm3.08 + pop`

Model 2: `poids.08 ~ volcm3.08 * pop`

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	197	13.529				
2	195	11.114	2	2.4147	21.183	4.733e-09 ***

Signif. codes: 0

↪ le modèle avec interactions est plus intéressant