

Analyses statistiques et représentations graphiques sous R

Sami Mestiri

Faculté des sciences économiques et de gestion Mahdia
mestirisami2007 @ gmail.com



Contents

Chapitre 1 Initiation à l'utilisation de R
1 Installation de logiciel R.....
2 Création et manipulation des données.....
2.1 Les vecteurs.....	9
2.2 Les matrices.....	11
2.3 Les facteurs.....	13
2.4 Les listes d'objets.....	13
2.5 Les data.frame.....	14
3 Quelques fonctions utiles.....	15
4.1 Manipulation de vecteurs.....	16
4.2 Recherche d'éléments dans un vecteur.....	16
4.3 Arrondissement.....	17
4.4 Statistiques descriptives.....	17
4.5 Opérations sur les matrices.....	18
4.6 Produit extérieur.....	19
4 Les graphiques.....	20
5 Les fonctions.....	21
6 Les structures de contrôle.....	22
6.1 Exécution conditionnelle.....	22
6.2 Boucles.....	24
Chapitre 2 Les distributions statistiques à un caractère	25
1 La terminologie statistique.....	26
2 Distribution statistique d'un caractère qualitatif.....	27
2.1 Tableau statistique.....	27
2.2 Représentation graphique.....	29
2.2.1 Diagramme en barres.....	29
2.2.2 Diagramme en secteurs.....	30
3 Distribution statistique d'un caractère quantitatif discret.....	31
3.1 Tableau statistique.....	31
3.2 Représentation graphique : Digramme en bâtons.....	33

4	Distribution statistique d'un caractère quantitatif continu.....	34
4.1	Représentation graphique.....	37
4.1.1	a) Cas d'amplitudes égales.....	38
4.1.2	b) Cas d'amplitudes inégales.....	39
5	Fonction de répartition.....	40
5.1	Cas d'un caractère quantitatif discret.....	41
5.2	Cas d'un caractère quantitatif continu.....	42
Chapitre 3 Les caractéristiques d'une distribution statistique.....		44
1	Les caractéristiques de position ou de tendance centrale.....	44
1.1	Le Mode.....	44
1.1.1	Cas d'une variable discrète.....	44
1.1.2	Cas d'une variable continue.....	45
1.2	La Médiane.....	47
1.2.1	Cas d'une variable discrète.....	47
1.2.2	Cas d'une variable continue.....	50
1.3	Les quantiles.....	51
1.3.1	Les quartiles.....	51
1.3.2	Les déciles.....	52
1.3.3	Les centiles.....	52
1.4	La moyenne arithmétique.....	53
1.4.1	Cas d'une série statistique de n observations.....	53
1.4.2	Cas d'une variable statistique discrète.....	53
1.4.3	Cas d'une variable statistique continue.....	54
1.5	Moyenne géométrique.....	54
1.6	Moyenne quadratique.....	56
1.7	Moyenne harmonique.....	56
2	Les caractéristiques de dispersion.....	57
2.1	L'étendue.....	58
2.2	L'intervalle interquartile.....	58
2.3	L'écart absolu moyen.....	58
2.4	La variance.....	58
2.5	L'écart-type.....	59
2.6	Le coefficient de variation.....	59
3	Les caractéristiques de forme.....	61
3.1	Les coefficients d'asymétrie.....	61
3.1.1	Repérage graphique de l'asymétrie.....	61
3.1.2	Mesure quantitative de l'asymétrie.....	62
3.1.3	Mesure quantitative d'aplatissement.....	63

Chapitre 4 La concentration.....	64
1 L'écart Médiale-Médiane.....	66
2 La courbe de concentration.....	68
3 Indice de concentration de Gini.....	70
Chapitre 5 Les distributions statistiques à deux caractères.....	72
1 La distribution à deux dimensions.....	72
2 Distributions marginales	73
3 Distributions conditionnelles.....	74
3.1 Distribution conditionnelle relative à X	74
3.2 Distribution conditionnelle relative à Y	74
4 Paramètres d'une distribution à deux dimensions.....	74
4.1 Moyennes et variances marginales.....	74
4.2 Moyennes et variances conditionnelles.....	75
5 Application.....	75
Chapitre 6 Analyse statistique bivariée	77
1 Deux variables quantitatives	77
1.1 La représentation graphique.....	77
1.2 Le coefficient de corrélation linéaire de Pearson	78
1.3 Droite de régression.....	80
1.4 L'ajustement linéaire par la méthode des moindres carrés.....	82
2 Deux variables qualitatives	84
2.1 Tableau de contingence.....	84
2.2 Tableau des fréquences	85
2.3 Le test de khi-deux.....	85
3 L'analyse d'une variable quantitative et qualitative.....	90
3.1 Applications sur l'ANOVA à un seul facteur	91

Chapitre 7 Les séries chronologiques	93
1 Définition.....	93
1.1 Les composantes d'une série chronologique	95
2 Estimation de la tendance.....	95
2.1 Moyennes mobiles.....	96
2.2 Ajustement analytique.....	97
3 Estimation des mouvements saisonniers.....	98
3.1 Modèle additif- multiplicatif.....	98
3.1.1 Méthode de profil.....	98
3.1.2 Méthode de la bande.....	99
3.1.3 Méthode analytique.....	100
4 Correction des variations saisonnières.....	101
4.1 Modèle additif- Données désaisonnalisées.....	101
4.2 Modèle multiplicatif- Données désaisonnalisées.....	105

Avant propos

La statistique est devenu de plus en plus une science de pleine actualité. Ses champs d'application se sont multipliés : Agronomie, Démographie, Économie, Marketing, Gestion de l'entreprise, Finance, Médecine, Sciences politiques, etc. Il est donc indispensable de connaître les principes et les indicateurs fondamentaux de la statistique.

L'objet de ce livre est de fournir aux lecteurs des outils de traitements statistiques des données avec logiciel R. Chaque chapitre de comporte une partie théorique et des exercices d'applications. Les applications sont programmées sous le logiciel R.

Dans certains cas on est amené à étudier plusieurs caractères sur une même population. Dans ce cas là, l'intérêt peut porter non seulement sur chaque caractère, mais également sur les liens qui peuvent exister entre les variables. En fait, certaines études visent à étudier conjointement deux variables mesurés sur un même individu. Il s'agit par exemple de la répartition des employés d'une entreprise selon leurs salaires et ancienneté, ou bien la répartition des ménages selon le nombre d'enfants à charge et leurs revenus. Nous nous intéresserons donc dans ce livre à l'étude d'une distribution statistique à deux dimensions.

Dans ce livre nous focalisant notre intérêt sur deux points particuliers. Le premier est la présentation des outils d'analyse bivariée, le but étant d'initier les étudiants à l'élaboration d'une étude utile pour l'exploitation des relations entre deux variables. Le second point sera une introduction à l'étude des séries chronologiques pour la détermination du cycle économique.

Chapitre1 : Initiation à l'utilisation de R

R est un logiciel libre distribué par "GNU Public Licence" (qui est une licence qui fixe les conditions légales de distribution des logiciels libres du projet GNU) et dérive du langage S (le logiciel S-PLUS). Il présente des caractéristiques remarquables comme la possibilité d'effectuer du calcul matriciel et d'autres opérations complexes, du stockage et de manipulation des données. Il contient des nombreuses fonctions pour les analyses statistiques et des outils graphiques flexibles. Il s'adresse à un large public formé de spécialistes et non-spécialistes en informatique.

Il a été initialement créé, en 1996, par Robert Gentleman et Ross Ihaka du département de statistique de l'Université d' Auckland en Nouvelle Zélande. Depuis 1997, il s'est formé une équipe "R Core Team" qui développe R. Il est conçu pour pouvoir être utilisé avec les systèmes d'exploitation Unix, Linux, Windows et MacOS.

R est devenu un logiciel libre et gratuit en 1995. R est à la fois un langage de programmation et un progiciel de fonctions statistiques. La version de base de R contient déjà un grand nombre de fonctions statistiques et graphiques permettant, par exemple, de calculer une moyenne ou une variance ou de tracer un histogramme.

Le logiciel R disponible gratuitement sur le site CRAN, à l'adresse suivante : cran.r-project.org. Ce logiciel est comparable à Matlab à certains égards, mais se révèle beaucoup plus puissant dans le domaine des traitements statistiques. De nombreux chercheurs ont développé au cours des années des fonctions plus avancées qui sont disponibles à tous les utilisateurs de R. Ces fonctions sont regroupées en bibliothèques qui sont disponibles pour téléchargement sur le site du projet R : <http://www.r-project.org/>.

Le " Comprehensive R Archive Network" (CRAN Réseau d'archives de R globales) est un ensemble de sites qui fournit ce qui est nécessaire à la distribution de R, ses extensions sa documentation, ses fichiers sources et ses fichiers binaires. Le site maître du CRAN est situé en Autriche à Vienne, nous pouvons y accéder par L'URL : <http://cran.r-project.org/>

1 Installation de logiciel R

Le logiciel est téléchargé sur le site web officiel de R (<http://www.r-project.org/>), il faut ensuite se diriger dans download CRAN. Choisissez un site miroir proche de chez vous (par exemple : en France), les téléchargements seront probablement plus rapides, vous trouvez ensuite un encadrement légendé download and install R.

Pour Windows, cliquez sur Windows puis base. Cliquez ensuite sur le fichier "exist" (par exemple : R-2.14.1-win.exe). Le programme d'installation est alors téléchargé sur votre ordinateur. Il suffit de cliquer dessus et de suivre les instructions. Un dossier portant le nom de la version de R téléchargé (R-2.14.1-win.exe dans ce cas) est créé. Il est situé, à partir du disque dur C, dans la série de dossiers suivante : program files /R.

Dans ce dossier se trouve le dossier library qui comprend les packages base de R. Un autre élément utile doit être localisé : le fichier .R data. Celui-ci n'est pas apparent au début. Il contiendra tous les objets que vous créez et sauvegardez dans R. Sur mon ordinateur, il apparaît, à partir du disque dur C.

2 Création et manipulation des données

On peut interagir directement avec R en ligne de commande, ou bien éditer un fichier texte (avec l'extension .R) sous son éditeur préféré.

Les données manipulées sous R sont stockées dans un espace de travail. Le répertoire de travail actuel peut être visualisé en tapant `getwd()`, tandis que l'accès à des répertoires se fait grâce à la commande `setwd()`.

On peut à n'importe quel moment de la session visualiser le contenu de cet espace de travail à l'aide de la commande `ls()`, et supprimer une variable à l'aide de la commande `rm()`, en lui passant comme argument le nom de la variable. Pour supprimer l'ensemble des variables contenues dans l'espace de travail, on utilise la commande: `rm(list=ls())`

Pour quitter l'environnement, il suffit de taper `q()`. Comme on vient de le dire, le travail sous R s'effectue via une session, que l'on peut à tout moment sauver à l'aide de la commande `save.image()`.

En tapant `ls()`, on constatera que nos variables sont toujours présentes dans l'espace de travail (celui-ci a en fait été sauve dans un fichier `.RData`, dans le répertoire de travail). La commande `history()` permet de lister les dernières commandes. Enfin, on peut sauver toutes les commandes ayant servi à analyser des données dans un fichier script avec l'extension `.R`, et taper `source('script.R',echo=T)`, depuis l'invite de commande, pour effectuer l'analyse.

R est un langage en ligne de commande (interprète), comme Matlab, et peut manipuler de nombreux types de données à l'aide de commandes prédéfinies. On va s'intéresser aux principaux types d'objets, à savoir :

- les vecteurs
- les matrices
- les facteurs
- les listes d'objets
- les dataframe

et voir quels sont les moyens de créer, manipuler et traiter de tels objets.

Notons que, comme dans tout langage, certains mots-clés sont réservés et ne peuvent être utilisés comme noms de variable ou de fonction :

FALSE Inf NA NaN NULL TRUE
break else for function if in next repeat while

2.1 Les vecteurs

Le vecteur est le type de base de R. Un nombre est simplement un vecteur à un seul élément. Notons dans un premier temps que R peut être utilisé comme un simple calculateur :

```
2+2
[1]4
2*pi
[1]6.283185
```

Certaines constantes (e.g. `pi`) et fonctions (`exp`, `cos`, etc.) sont connues de R et peuvent être appelées directement. On peut également assigner une valeur à une variable à l'aide de l'opérateur `←` (ou `<-`) :

```
x ← 10
```

La variable peut ensuite être utilisée comme dans n'importe quel langage interprété :

```
x+2
[1]12
y ← x+4
y
[1]14
```

Bien que nous n'ayons utilisé que des variables de type numérique, il existe d'autres modes de représentation des données (ou classes) dans un vecteur :

character, integer, logical, complex, list.

Enfin, les types NULL et NA jouent un rôle particulier puisqu'ils désignent respectivement l'absence de valeur (i.e. un ensemble vide) et le codage par défaut d'une valeur manquante.

Différentes fonctions permettent de générer facilement des vecteurs :c, seq, rep. Par exemple, on peut créer le vecteur $X = [1\ 2\ 3\ 4\ 5]$ de différentes façons :

```
X ← c(1,2,3,4,5)
X
[1]12345
rm(X)
```

La fonction seq() est utile pour générer des suites d'entiers, et possède plusieurs arguments pour spécifier la séquence recherchée (voir ?seq).

```
Y ← seq(1:5)
Y
[1]12345
```

On peut dupliquer un vecteur à l'aide de la commande rep() :

```
Z ← rep(X,2)
Z
[1]1234512345
```

La fonction rev() permet quant à elle de renverser l'ordre de la séquence (1 2 3 devient 3 2 1).

```
W ← rev(X)
W
[1]54321
```

Les opérations arithmétiques s'appliquent aussi sur des vecteurs. Par exemple, on peut additionner deux vecteurs a et b, et effectuer les produits scalaires ou membre à membre classiques :

```
a ← 1:5
b ← 5:9
a × 2
[1]246810
a+b
[1]68101214
a × b
[1]512213245
```

De nombreuses fonctions permettent également de classer les éléments (ou les indices) d'un vecteur, de les sommer, etc.

```
u ← c(i)
u
[1]13274
sort(u) # ordonner les éléments du vecteur u
[1]12347
```

`order(u)` # donner le rang des éléments du vecteur u
[1]13254

`sum(u)` # donner la somme des éléments du vecteur u
[1]17

`length(u)` donner la longueur du vecteur u
[1]5

Enfin, une particularité de R est que les éléments d'un vecteur peuvent avoir des noms. La fonction `names()` permet en effet d'associer une étiquette à chacun des éléments d'un vecteur :

```
x ← 1:5  
names(x) ← c( a , b , c , d , e )
```

```
x  
a b c d e  
1 2 3 4 5
```

2.2 Les matrices

Pour les matrices, on peut soit générer un vecteur de taille n, et le réarranger pour créer une matrice de taille (l,m) , soit directement créer une matrice à l'aide de la fonction `matrix()`, en spécifiant en arguments le nombre de lignes et de colonnes, par exemple :

```
M ← c(1,2,3,4,5,6,7,8,9,10), nrow=2, ncol=5
```

```
           [,1]      [,2]      [,3]      [,4]      [,5]  
x         1.00      3.00      5.00      7.00      9.00  
y         2.00      4.00      6.00      8.00     10.00
```

Les fonctions `cbind` et `rbind` permettent de manipuler des vecteurs de manière à former une matrice par concaténation sur les colonnes ou sur les lignes.

```
x ← c(1,3,5,7,9)  
y ← c(2,4,6,8,10)  
A ← cbind(x,y)
```

```
           x         y  
[1]      1.00      2.00  
[2]      3.00      4.00  
[3]      5.00      6.00  
[4]      7.00      8.00  
[5]      9.00     10.00
```