

T. D. n° 2

Statistiques descriptives avec le logiciel R

Introduction : Ce T.D. a pour but de vous faire assimiler les représentations numériques et graphiques des statistiques descriptives univariée et bivariée.

Objectif de ce T.D. :

- Manipuler les données.
- Faire des résumés numériques et/ou graphiques.

Consignes pour ce T.D. :

- Suivre pas à pas les étapes et voir ce qui se passe.
- Ne pas hésiter à utiliser l'aide en ligne de R.
- Vous ne comprendrez peut-être pas tous les détails mais la meilleure chose à faire est de taper le code et de voir le résultat produit. Soyez curieux et n'hésitez pas à le modifier pour voir « ce qu'il se passe ».

Quelques remarques :

- Le symbole # signifie le début d'un commentaire.
- Lorsque vous travaillez sous R, il peut être intéressant de conserver les résultats et les graphiques de vos analyses. Le plus simple, dans un premier temps, est de les enregistrer dans un document word à l'aide du copier / coller. Pour ce faire, aller dans le menu « File » ou « Fichier », sélectionner « Copy to the clipboard » « as a Bitmap ». Noter que les graphes peuvent être réduits ou agrandis sans déformation.
- Parfois le signe + peut apparaître en début de ligne de commande de R. Ne le tapez pas svp. Il est là pour rappeler qu'une ligne a été coupée et que nous en somme en début de ligne.

Statistique descriptive univariée

Exercice 1 Fichier de données : iris.

Le logiciel R est un ensemble de bibliothèques de fonctions appelées « packages ». Chaque bibliothèque contient des jeux de données.

Pour connaître par exemple les jeux de données contenus dans le package `base`, écrire l'instruction suivante :

```
> data(package = "base").
```

Le résultat apparaît dans une fenêtre R data sets. En voici un extrait :

```
Data sets in package 'datasets':
```

```
AirPassengers.....Monthly Airline Passenger Numbers 1949-1960  
BJsales.....Sales Data with Leading Indicator  
BJsales.lead (BJsales)..Sales Data with Leading Indicator  
BOD.....Biochemical Oxygen Demand
```

...

`iris.....Edgar Anderson's Iris Data`

1. Noter la présence du fichier `iris` dans la liste ci-dessus. Les données de ce fichier sont célèbres. Elles ont été collectées par Edgar Anderson¹. Le fichier donne les mesures en centimètres des variables suivantes :
 - longueur du sépale (`Sepal.Length`),
 - largeur du sépale (`Sepal.Width`),
 - longueur du pétale (`Petal.Length`),
 - largeur du pétale (`Petal.Width`)
 pour trois espèces d'iris qui sont les :

1. Iris setosa,
2. Iris versicolor et
3. Iris virginica.

Sir R.A. Fisher² a utilisé ces données pour construire des combinaisons linéaires des variables permettant de séparer au mieux les trois espèces d'iris.

2. Pour analyser le fichier `iris`, il faut le charger. Quelle est l'instruction qu'il faut taper pour charger ce fichier ?
3. Taper une à une chacune des instructions ci-dessous et noter le résultat obtenu si possible.

Attention : le logiciel R n'est pas indifférent aux majuscules et aux minuscules, comme nous l'avons déjà souligné dans le T.D. 1.

```
> iris
```

```
> dim(iris)
```

```
> names(iris)
```

Quelle(s) différence(s) faites-vous avec la commande `> str(iris)` ?

4. Taper les lignes de commande suivantes :

```
> iris$Petal.Length
```

```
> iris$Species
```

Qu'observez-vous ?

5. La dernière colonne du fichier `iris` contient le nom des espèces réparties en trois catégories : setosa, versicolor et virginica. Pour accéder à celle-ci, il faut utiliser l'instruction `iris$Species`, comme vous venez de le constater à la question précédente. Nous disons alors que la dernière colonne contient une variable qualitative à trois modalités ou à trois niveaux³ appelés « levels » par le logiciel R. La fonction `levels` appliquée à la colonne `iris$Species` donne les modalités de la variable. En effet, il suffit de taper :

```
> levels(iris$Species)
```

1. E. Anderson (1935) The irises of the Gaspe Peninsula, *Bulletin of the American Iris Society*, 59, 2-5

2. R.A. Fisher, (1936) The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7, Part II, 179-188

3. Nous reverrons ces définitions et ces notations lors de l'analyse de la variance.

Pour résumer l'information contenue dans cette variable, vous utiliserez l'instruction suivante :

```
> summary(iris$Species)
```

Quel est le résultat que vous obtenez ?

6. Cette dernière information peut être obtenue en construisant un tableau (`table`) comptabilisant le nombre d'individus par modalité. Pour ce faire, taper l'instruction suivante :

```
> table(iris$Species)
```

Comparer avec le résultat obtenu à la question précédente.

7. **R permet également de réaliser des résumés graphiques.** Lorsqu'une instruction graphique est lancée, une nouvelle fenêtre, `R Graphics:Device`, s'ouvre. Les représentations graphiques liées aux variables qualitatives sont :
— le diagramme circulaire ou encore le camembert, voire la fonction (`pie()`)
— la diagramme en bâtons, voire la fonction (`barplot()`).

Taper les lignes de commande suivantes :

```
> pie(table(iris$Species))
```

```
> barplot(table(iris$Species))
```

8. Il existe une fonction, la fonction `par()`, permettant de découper la fenêtre graphique de deux façons :

```
par(mfrow=c(nl,nc)) ou par(mfcol=c(nl,nc)),
```

où `nl` définit le nombre de graphiques en lignes, `nc` définit le nombre de graphiques en colonnes, `mfrow` signifie que l'ordre d'entrée des graphiques s'effectue selon les lignes et `mfcol` signifie que l'ordre d'entrée des graphiques s'effectue selon les colonnes.

Supposons que vous vouliez représenter six graphiques dans une seule fenêtre en deux lignes et trois colonnes. La première instruction conduit à entrer les graphiques selon l'ordre :

1	2	3
4	5	6

La seconde instruction conduit à entrer les graphiques selon l'ordre :

1	3	5
2	4	6

Deux botanistes se sont également intéressés aux iris et ont collecté les espèces suivantes :

```
> collection1<-rep(c("setosa","versicolor","virginica"),
+c(15,19,12))
```

```
> collection2<-rep(c("setosa","versicolor","virginica"),
+c(22,27,17))
```

En utilisant la fonction `par(mfrow=c(2,2))`,

1. Construire les camemberts de ces deux nouvelles distributions. Commenter.

2. Construire les diagrammes en bâtons de ces deux nouvelles distributions. Commenter.
 3. Discuter des avantages et des inconvénients de ces deux types de représentations.
9. **Revenons au résumé numérique.** La troisième colonne du fichier `iris` contient la longueur du pétale. Il s'agit d'une variable mesurée qualifiée alors de variable quantitative. Pour résumer l'information contenue dans cette variable, nous utilisons la fonction `summary()`. Taper la ligne de commande suivante :

```
> summary(iris$Petal.Length)
```

Le résultat obtenu est le suivant :

Min.	1stQu.	Median	Mean	3rdQu.	Max.
1.000	1.600	4.350	3.758	5.100	6.900

La plus petite (Min.) longueur de pétale est égale à 1,000 *cm* tandis que la plus grande (Max.) est égale à 6,900 *cm*. La moyenne (Mean) représente la somme des valeurs de la distribution divisée par le nombre total d'iris. Elle est égale à 3,758 *cm*. Si l'ensemble des 150 longueurs de pétale est classé par ordre croissant, 1stQu., Median et 3rdQu. sont les trois valeurs qui permettent de couper la distribution en quatre parties égales. Nous les appelons respectivement premier quartile, médiane (ou deuxième quartile) et troisième quartile.

Essayons de retrouver ces six valeurs de tendance centrale. Taper les lignes de commande suivantes :

```
> min(iris$Petal.Length)
> max(iris$Petal.Length)
> mean(iris$Petal.Length)
```

Remarque : pour calculer la moyenne nous aurions pu procéder autrement. Taper les lignes de commande suivantes :

```
> sum(iris$Petal.Length)
> length(iris$Petal.Length)
> sum(iris$Petal.Length)/length(iris$Petal.Length)
```

Obtenez-vous le même résultat que précédemment ?

Maintenant occupons-nous de retrouver les valeurs des trois quartiles. Pour cela, taper la ligne de commande suivante :

```
> sort(iris$Petal.Length)
```

Que se passe-t-il ?

Puis continuer par taper les lignes de commandes suivantes :

```
> ordLpetal <- sort(iris$Petal.Length)
> ordLpetal # commenter le résultat
> sum(ordLpetal)/length(ordLpetal)
> ordLpetal[38]
> (ordLpetal[75]+ordLpetal[76])/2
> ordLpetal[113]
```

Auriez-vous pu faire plus simple en utilisant une autre fonction ? Si oui, laquelle ?

10. **Retour au résumé graphique.** Une des représentations adéquates est l'histogramme. Regarder l'aide de `hist()`. Puis, taper la ligne de commande suivante :

```
> hist(iris$Petal.Length,col=grey(0.6),main="Histogramme")
```

Remarque : `main` est l'option qui permet d'afficher un titre dans un graphique.

11. Réaliser le même type d'analyse sur chacune des trois autres variables quantitatives : largeur du pétale, longueur du sépale et largeur du sépale. Notez que vous n'avez pas toutes les instructions à réécrire en utilisant le système de flèches du clavier. \uparrow et \downarrow vous permettent de retrouver les fonctions que vous avez utilisées. \leftarrow et \rightarrow vous permettent de vous déplacer dans la fonction et donc, dans changer certains paramètres.

Exercice 2 Fichier de données : Europe.

Nous vous demandons dans cet exercice de faire des résumés numériques mais aussi de tracer une boîte à moustaches sur le jeu de données `Europe`.

Pour cela, il faut d'abord que vous installiez le package `BioStatR`. La marche à suivre pour installer un package est toujours la même.

- Dans la barre du haut, vous cliquez sur `Packages`,
- vous allez sur l'onglet `Installer le(s) package(s)`,
- vous choisissez un `mirror`, le plus proche du lieu où vous vous trouvez.
- Une nouvelle fenêtre va s'ouvrir qui contient une liste de packages. Vous cliquez sur le package `BioStatR` (ou sur un autre quand ce sera un autre qu'il faudra utiliser).

Vous allez voir des lignes qui s'écrivent dans la fenêtre `R Console` ce qui signifie que tout s'installe bien. N'oubliez pas à l'issue de ces installations de taper la ligne de commande suivante :

```
> library(BioStatR)
```

Si tout s'est bien passé, vous devez lire :

Message d'avis :

```
le package "BioStatR" a été compilé avec la version R 3.0.3
```

```
# Description du jeu de données et statistiques descriptives
```

```
> head(Europe)
```

```
> str(Europe)
```

```
> summary(Europe)
```

```
> range(Europe$Duree)
```

```
> sd(Europe$Duree)
```

```
# Boîte à moustaches
```

```
> boxplot(Europe$Duree,ylab="Durée (heures)")
```

```
> points(1,mean(Europe$Duree),pch=2)
```

Remarque : `pch` est une option graphique qui définit le symbole qui représente les observations.

Exercice 3 Données brutes ou groupement en classes.

Parfois, lorsque nous étudions une série statistique sur un caractère quantitatif qui comporte un grand nombre de valeurs, nous préférons alors regrouper par classes puis ensuite remplacer chaque classe par son milieu. Mais les résultats en sont légèrement modifiés, ce que vous pouvez imaginer. D'ailleurs certains auteurs suggèrent des corrections par exemple en ce qui concerne la variance, comme la correction de Sheppard, comme le déclarent Couty, Debord et Fredon dans leur livre « Mini manuel de probabilités et statistique », Dunod. D'ailleurs de ce livre, nous allons extraire le jeu de données qui va nous permettre de faire cet exercice.

Nous considérons une série statistique de 60 taux d'hémoglobine dans le sang (g/L) mesurés chez des adultes présumés en bonne santé :

Femmes	105	110	112	112	118	119	120	120	125	126
	127	128	130	132	133	134	135	138	138	138
	138	142	145	148	148	150	151	154	154	158
Hommes	141	144	146	148	149	150	150	151	153	153
	153	154	155	156	156	160	160	160	163	164
	164	165	166	168	168	170	172	172	176	179

1. Nous considérons le groupement en classes suivant :

$$]104; 114];]114; 124];]124; 134];]134; 144];]144; 154];]154; 164]; \\]164; 174];]174; 184].$$

- Pour chacune des deux séries : femmes et hommes, déterminer les effectifs et les fréquences de chaque classe.
2. Effectuer une représentation graphique adaptée des deux distributions groupées en classe de la question 1.
 3. Calculer les moyennes des trois distributions initiales : ensemble, femmes, hommes.
 4. Calculer les moyennes des trois distributions (ensemble, femmes, hommes) après le regroupement en classes de la question 1., en remplaçant chaque classe par son milieu.
 5. Calculer les médianes des trois distributions initiales : ensemble, femmes, hommes.
 6. Calculer l'écart interquartile pour chacune des trois distributions initiales : ensemble, femmes, hommes.
 7. Calculer les variances et les écart-types des trois distributions initiales : ensemble, femmes, hommes.
 8. Calculer les variances et les écart-types des trois distributions après le regroupement en classes de la question 1., en remplaçant chaque classe par son milieu.
 9. Pour la distribution des femmes, calculer les moments jusqu'à l'ordre 4 puis déduire les moments centrés jusqu'à l'ordre 4 puis les paramètres de caractéristique de forme de Fisher.

Statistique descriptive bivariée

Exercice 4 Suite de l'exercice 1.

Nous allons reprendre les données de l'exercice 1. Voici la suite des questions.

12. Une fois réalisées les graphiques pour chaque variable prise séparément, l'étude peut porter sur la relation entre deux variables. Nous parlons alors de croisement de deux variables ou d'étude bivariée. La représentation graphique liant deux variables quantitatives est le nuage de points.

Représentons par exemple la longueur et la largeur du pétale pour les 150 iris contenus dans le fichier de données. Pour cela, exécuter la ligne de commentaire suivante :

```
>plot(iris$Petal.Length, iris$Petal.Width,
+ xlab="Longueur du pétale", ylab="Largeur du pétale",
+ main="Nuage de points", pch=20)
```

Faire un commentaire.

Dans cette représentation graphique, plusieurs individus peuvent être situés sur un même point. La fonction `sunflowerplot` permet de visualiser ces superpositions. Taper la ligne de commande suivante :

```
> sunflowerplot(iris$Petal.Length, iris$Petal.Width,
+ xlab="Longueur du pétale", ylab="Largeur du pétale",
+ main="Nuage de points", pch=20)
```

13. Réaliser l'étude du croisement de deux variables quantitatives de votre choix. Il est clair que le sens biologique de l'étude ne doit pas être négligé.
14. La représentation graphique permettant de lier une variable qualitative et une variable quantitative est la boîte à moustaches (`boxplot`). Représentons par exemple la longueur des pétales en fonction de l'espèce. Pour cela, taper la ligne de commande suivante :

```
> boxplot(iris$Petal.Length ~ iris$Species, col=grey(0.6))
```

Commenter.

15. Choisir une autre variable quantitative, croiser-la avec la variable espèce d'iris et commenter.

16. Le nuage de points comme les boîtes à moustaches montrent que les données morphologiques des iris semblent liées à l'espèce. Il pourrait donc être intéressant de réaliser des graphiques différents pour chacune des modalités Iris setosa, Iris Versicolor et Iris Virginica ou de superposer l'information espèce dans le graphique des nuages de points. Nous vous proposons ici quelques développements. Libre à vous, de les refaire ou d'en trouver d'autres...

Taper alors les lignes de commande qui vont suivre :

```
# Tracé des histogrammes des longueurs des pétales de l'ensemble
des iris, des iris setosa, des iris versicolor et
des iris virginica
> par(mfrow=c(2,2))
```

```

> br0=seq(0,8,le=20)
> hist(iris$Petal.Length, main="Ensemble des 150 iris",
+ xlab="Longueur du pétale", br=br0)
> hist(iris$Petal.Length[iris$Species=="setosa"], main="Setosa",
+ xlab="Longueur du pétale", br=br0)
> hist(iris$Petal.Length[iris$Species=="versicolor"],
+ main="Versicolor", xlab="Longueur du pétale", br=br0)
> hist(iris$Petal.Length[iris$Species=="virginica"],
+ main="Virginica", xlab="Longueur du pétale", br=br0)

# Tracé des nuages des points de la largeur du pétale en fonc-
tion de la longueur du pétales de l'ensemble des iris,
des iris setosa, des iris versicolor et des iris virginica
> par(mfrow=c(2,2))
> plot(iris$Petal.Length, iris$Petal.Width,
+ xlab="Longueur du pétale", ylab="Largeur du pétale",
+ main="Nuage de points", pch=20)
> plot(iris$Petal.Length[iris$Species=="setosa"],
+ iris$Petal.Width[iris$Species=="setosa"],
+ xlim=c(1,6.9), ylim=c(0.1,2.5), xlab="",ylab="",
+ main="iris setosa", pch=20)
> plot(iris$Petal.Length[iris$Species=="versicolor"],
+ iris$Petal.Width[iris$Species=="versicolor"],
+ xlim=c(1,6.9), ylim=c(0.1,2.5), xlab="", ylab="",
+ main="iris versicolor", pch=20)
> plot(iris$Petal.Length[iris$Species=="virginica"],
+ iris$Petal.Width[iris$Species=="virginica"],
+ xlim=c(1,6.9), ylim=c(0.1,2.5), xlab="", ylab="",
+ main="iris virginica", pch=20)
17. Et pour finir...Taper la ligne de commande suivante :
> pairs(iris[1:4], main = "Anderson's Iris Data - 3 species",
pch = 21, bg = c("red", "green3", "blue")[unclass(iris$Species)])
# Représentation graphique de toutes les possibilités de
variables par variables

```

Qu'observez-vous ?

La fonction `pairs()` reproduit tous les graphiques variables par variables possibles sur une seule fenêtre graphique et `bg` est une option graphique pour définir la couleur.