

## TP 3 : Graphiques avec R

*\*\*\* Manipulation des objets matrix et data.frame \*\*\**

### Table des matières

<b>1</b>	<b>Les graphiques avec R</b>	<b>2</b>
1.1	Les fonctions graphiques de “haut niveau” . . . . .	2
1.2	Fonctions graphiques de “bas niveau” . . . . .	3
1.3	La gestion des graphiques . . . . .	3
1.4	Une illustration graphique . . . . .	4
<b>2</b>	<b>Zoom sur certaines commandes graphiques</b>	<b>5</b>
2.1	Diagramme en bâtons . . . . .	5
2.2	Diagramme circulaire . . . . .	7
2.3	Boîtes à moustaches . . . . .	7
2.4	Histogramme . . . . .	8
<b>3</b>	<b>Exercices d’application</b>	<b>9</b>
3.1	Exercice 1 . . . . .	9
3.2	Exercice 2 . . . . .	10

# 1 Les graphiques avec R

R est un logiciel qui offre une multitude de possibilités graphiques qu'il est impossible de détailler dans ce document (`demo(graphics)` permet d'obtenir un avant-goût de ces possibilités). Cependant, les paragraphes suivants vous permettront d'acquérir les bases nécessaires à l'utilisation des graphiques en R.

R fonctionne au moyen de deux types de fonctions graphiques :

- les fonctions graphiques principales ou de “haut niveau” qui créent un nouveau graphique,
- les fonctions graphiques secondaires ou de “bas niveau” qui permettent d'ajouter de nouveaux éléments sur un graphique déjà produit ou de prédéfinir des paramètres graphiques d'un graphique à produire.

## 1.1 Les fonctions graphiques de “haut niveau”

Le tableau TAB. 1 fournit une liste non exhaustive des fonctions graphiques principales.

Fonction graphique	Description
<code>plot(X)</code>	Graphe des valeurs de X en fonction de leur indice
<code>plot(X,Y)</code>	Graphe des valeurs de Y en fonction des valeurs de X
<code>matplot(X,Y)</code>	Graphe de la 1ère colonne de Y en fonction de la 1ère colonne de X, graphe de la 2ème colonne de Y en fonction de la 2ème colonne de X ...
<code>pairs(X)</code>	Si X est une matrice ou un data.frame, <code>pairs</code> crée tous les graphes bivariés entre les colonnes de X
<code>contour(X,Y,Z)</code>	Courbes de niveau (X et Y sont les vecteurs de localisation dans le plan, Z est la matrice des niveaux de dimension <code>c(length(X),length(Y))</code> )
<code>pie(X)</code>	Diagramme circulaire de X
<code>boxplot(X)</code>	Boîte à moustaches de X
<code>dotchart(X)</code>	Graphe de Cleveland de X
<code>hist(X)</code>	Histogramme des fréquences de X
<code>barplot(X)</code>	Diagramme en bâtons de X
<code>qqplot(X,Y)</code>	Graphe quantile-quantile des échantillons X et Y

TABLE 1 – Fonctions graphiques de haut-niveau

Chaque fonction possède des options graphiques dont certaines sont communes à plusieurs fonctions. Le tableau TAB. 2 liste un certain nombre d'options graphiques.

Option graphique (valeur par défaut)	Description
<code>add=FALSE</code>	Si <code>add=TRUE</code> , superpose le graphe au graphe précédent
<code>axes=TRUE</code>	Si <code>axes=FALSE</code> , ne trace pas les axes ni le cadre
<code>type="p"</code>	Type de graphique (“p” points, “l” lignes, “b” lignes et points, “h” lignes verticales, “s” et “S” escaliers...)
<code>xlim, ylim</code>	Limites inférieure et supérieure des axes
<code>xlab, ylab</code>	Annotations des axes
<code>main, sub</code>	Titre et sous-titre du graphique

TABLE 2 – Options graphiques

## 1.2 Fonctions graphiques de “bas niveau”

Plusieurs fonctions graphiques secondaires de R permettent d’ajouter de nouveaux éléments à un graphique déjà existant. Le tableau TAB 3. recense un grand nombre de ces fonctions.

Fonction graphique secondaire	Description
<code>points(X,Y)</code>	Ajoute les points de coordonnées définies par les vecteurs <b>X</b> et <b>Y</b>
<code>lines(X,Y)</code>	Analogue de <code>points(X,Y)</code> avec des lignes
<code>text(x,y,labels, ...)</code> <code>text(locator(1),labels, ...)</code>	Ajoute le texte défini par <code>labels</code> au point <code>(x,y)</code> Analogue avec un choix de la localisation à l’aide de la souris
<code>legend(x,y,legend, ...)</code>	Inscrit la légende donnée par <code>legend</code> au point <code>(x,y)</code>
<code>title(main,sub)</code>	Ajoute le titre <code>main</code> , le sous-titre <code>sub</code>
<code>abline(b,a)</code>	Trace une ligne d’équation $y=ax+b$
<code>abline(h=b)</code>	Trace une ligne d’équation $y=b$
<code>abline(v=k)</code>	Trace une ligne d’équation $x=k$
<code>rug(x)</code>	Dessine des traits verticaux aux points d’abscisses données par <code>x</code>

TABLE 3 – Fonctions graphiques secondaires permettant l’ajout d’éléments graphiques

*Remarque.* Il est possible de rajouter sur un graphique, au point de coordonnées  $(x,y)$ , une expression mathématique au moyen de la fonction `text` :

```
text(x,y,expression(...))
```

Si l’on souhaite inclure dans une expression mathématique la valeur numérique d’une variable `v`, on peut faire appel aux fonctions `substitute` et `as.expression` sous la forme suivante :

```
text(x,y,as.expression(substitute(label==value,list(value=v))))
```

Aux fonctions graphiques secondaires s’ajoutent les paramètres graphiques. Ils peuvent être utilisés sous forme d’options graphiques pour certains ou bien à l’aide de la fonction `par()`. La liste de ces paramètres graphiques est obtenue par la commande `help(par)`. Le tableau TAB. 4 dresse la liste des paramètres que nous serons amenés à utiliser en TP.

Paramètre	Description
<code>bg</code>	Spécifie la couleur du fond graphique (ex : <code>bg='lightyellow'</code> )
<code>font</code>	Contrôle le style du texte (1 : normal, 2 : italique, 3 : gras, 4 : gras italique)
<code>pch</code>	Contrôle le type de symbole (ex : <code>pch=25</code> )
<code>col</code>	Contrôle la couleur des symboles ou des lignes
<code>lwd</code>	Gère la taille des lignes
<code>xaxt, yaxt</code>	Si <code>xaxt='n'</code> ( <code>yaxt='n'</code> ) l’axe des abscisses (ordonnées) n’est pas tracé
<code>new</code>	Si <code>new=TRUE</code> , le nouveau graphique sera tracé sur l’ancien graphique

TABLE 4 – Paramètres graphiques

## 1.3 La gestion des graphiques

Un dispositif graphique s’ouvre lors de l’exécution de la première fonction graphique. L’appel à une autre fonction graphique remplace, dans ce même dispositif graphique, l’ancien graphique

par le graphique courant. Pour visualiser simultanément les deux représentations graphiques, trois approches sont possibles :

- La commande `par(new=TRUE)`, intercalée entre les deux commandes graphiques, permet de superposer le nouveau graphique à l'ancien dans le dispositif graphique.
- Il est aussi possible d'ouvrir un nouveau dispositif graphique grâce à la commande `windows()` sous Windows (`x11()` sous Linux, `quartz()` sous OS X ou `dev.new()` le cas échéant) avant de taper la deuxième commande graphique. On visualise alors deux graphiques chacun dans une fenêtre graphique différente.
- Une autre possibilité consiste à afficher plusieurs graphiques dans un même dispositif graphique au moyen de la commande `layout` qui partitionne le dispositif en plusieurs parties où s'affichent successivement les différents graphiques. Par exemple,

```
> layout(matrix(1:6,3,2)) # Création d'un dispositif graphique à 6 fenêtres
> layout.show(6)         # Affichage de l'organisation du dispositif
```

Enfin, pour sauvegarder un graphique dans un fichier *Plot.eps* (Encapsulated PostScript) ou *Plot.pdf* (Portable Document Format), il convient de taper la commande `dev.copy2eps(file="Plot.eps")` ou `dev.copy2pdf(file="Plot.pdf")`.

## 1.4 Une illustration graphique

Le code suivant permet de réaliser différents graphiques dont celui de la FIG 1.

```
> par(mfrow=c(1,2))
> # Tracé de l'objet Nile de type ts
> plot(Nile)
> title(main="Un petit titre", sub="Un sous-titre")
> # Tracé de l'objet AirPassengers de type ts dans le même dispositif
> plot(AirPassengers,lwd=3,type="b",col="skyblue2",col.main=3,main="Joli !",font.main=2)
> # Tracé de l'objet volcano de type matrix
> par(mfrow=c(1,2))
> x=10*1:87
> y=10*1:61
> image(x,y,volcano,col=terrain.colors(20))
> contour(x,y,volcano,col=terrain.colors(15))
> windows()
> persp(x,y,volcano,phi=40,theta=30,expand=0.75,col="lightgreen")
> # Etude de l'objet cars de type data.frame
> names(cars) # Affiche le nom des variables du data.frame
> row.names(cars) # Affiche le nom des lignes du data.frame
> cars$speed # Affiche les données de la première variable
> cars[,1] # Affichage identique par syntaxe matricielle
> attach(cars) # Accès aux variables par les noms speed et dist
> dev.new() # Ouverture d'un nouveau dispositif graphique
> pairs(cars) # Tracé des nuages de points
> matplot(cars,type="l") # Tracé des courbes sur un même graphique
> windows() # Ouverture d'un nouveau dispositif
> layout(matrix(1:4,2,2),width=c(4,1),height=c(2,2))
> # Graphiques de Cleveland et boîtes à moustaches des données cars
> dotchart(speed, main="Cleveland speed")
> text(5,20,as.expression(substitute(min==value1,list(value1=min(speed)))))
> text(5,16,as.expression(substitute(MAX==value2,list(value2=max(speed)))))
```

## Volcan 3D

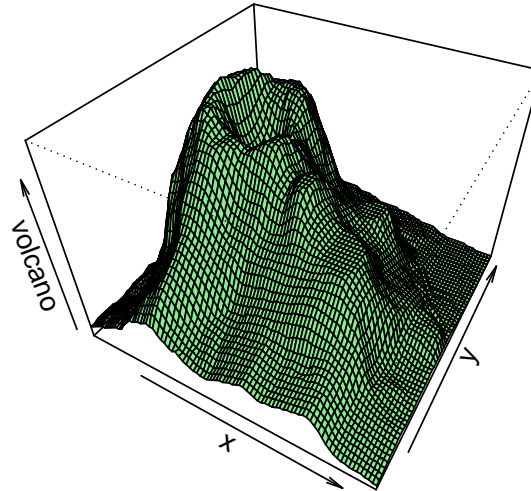


FIGURE 1 – Graphique 3D des données volcano

```
> dotchart(dist, main="Cleveland dist")
> boxplot(speed,main="Boxplot speed")
> boxplot(sort(dist), main="Boxplot dist")
> # Histogrammes des données speed
> x11()
> par(bg="lightgreen",mfrow=c(1,2))
> barplot(table(speed), main="Diagramme en bâtons",col=rainbow(10))
> hist(speed, main="Histogramme",col="tomato")
> rug(speed)
> dev.copy2pdf(file="Histogramme.pdf") # Sauvegarde du dernier graphique
> detach(cars)                        # Détachons les données cars
```

## 2 Zoom sur certaines commandes graphiques

### 2.1 Diagramme en bâtons

Pour une série statistique d'observations  $(x_j)_{j=1\dots n}$  (généralement  $x_j$  est l'effectif de la  $j$ ème valeur d'un caractère discret), le diagramme en bâtons est formé de  $n$  bâtons (ou de barres si l'on donne une largeur aux bâtons) où la hauteur du  $j$ ème bâton est égale à  $x_j$ .

La commande R pour obtenir un diagramme en bâtons est `barplot(height)` où `height` est le

vecteur ou la matrice de données. Sont disponibles, en plus des options classiques `main`, `xlim`, `xlab`, `ylab` ..., les options suivantes :

- `width` : vecteur de largeurs de barres.
- `space` : espacement inter-barres (ou inter-blocs si `height` est une matrice) exprimé sous forme de fraction de la largeur moyenne des barres.
- `names.arg` : vecteur des noms à inscrire sous les barres.
- `beside` : si `beside=FALSE` les barres sont espacées, si `beside=TRUE` les barres sont juxtaposées
- `col` : un vecteur de couleurs pour les barres.
- `horiz` : si `horiz=FALSE` les barres sont verticales, si `horiz=TRUE` les barres sont horizontales.

Par exemple, les données `VADeaths` de type `matrix`, représentant le nombre de morts pour 1000 personnes dans plusieurs sous-populations de Virginie en 1940, peuvent être représentées par diagrammes en bâtons.

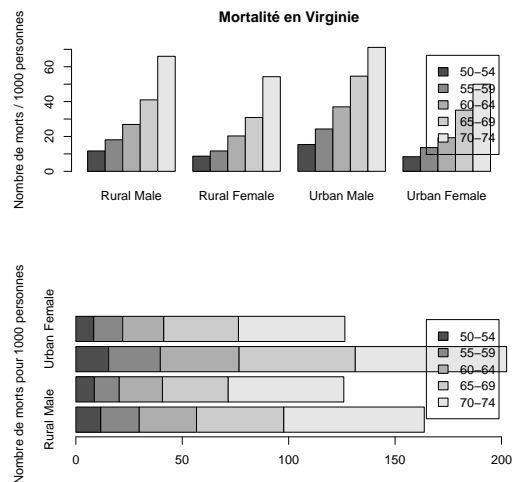


FIGURE 2 – Barplots des données `VADeaths`

Le code produisant la figure FIG. 2 est le suivant :

```
> VADeaths
      Rural Male Rural Female Urban Male Urban Female
50-54      11.7         8.7      15.4         8.4
55-59      18.1        11.7      24.3        13.6
60-64      26.9        20.3      37.0        19.3
65-69      41.0        30.9      54.6        35.1
70-74      66.0        54.3      71.1        50.0
> layout(c(1,2))
> barplot(VADeaths, beside=TRUE, legend=TRUE, main="Mortalité en
+ virginie", ylab="Nombre de morts / 1000 personnes")
> barplot(VADeaths, beside=FALSE, horiz=TRUE, legend=TRUE, ylab=
+ "Nombre de morts pour 1000 personnes")
```

## 2.2 Diagramme circulaire

Pour une série statistique d'observations  $(x_j)_{j=1\dots n}$  (généralement  $x_j$  est l'effectif de la  $j$ ème valeur d'un caractère discret), le diagramme circulaire est un disque divisé en  $n$  secteurs où l'angle (et donc la surface) du  $j$ ème secteur est proportionnel à  $x_j$ .

Si une population est subdivisée en sous-groupes dont la répartition en pourcentages est `prop` et les noms des sous-groupes sont `noms`, la commande R `pie(prop,noms)` permet de tracer le camembert de la population. Hormis les options classiques, on peut préciser l'angle (en degrés) de départ avec l'option `init.angle` et la couleur de chaque secteur avec l'option `col`.

Exemple de répartition de notes

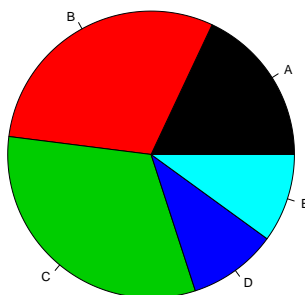


FIGURE 3 – Exemple de diagramme circulaire

La figure FIG. 3 est le résultats des commandes suivantes :

```
> prop=c(18,30,32,10,10)
> noms=c("A","B","C","D","E","F")
> pie(prop,noms,col=1:6, main="Répartition des notes d'une classe
+ d'étudiants en statistique)
```

## 2.3 Boîtes à moustaches

Pour des données brutes  $(x_j)_{j=1\dots n}$ , observations d'un caractère continu, la boîte à moustaches est une alternative à l'histogramme permettant de se faire rapidement une idée de la répartition d'un ensemble de données :

- La partie principale d'une boîte à moustaches est un rectangle horizontal dont la borne inférieure correspond au premier quartile  $Q_1$  des données et la partie supérieure au troisième quartile  $Q_3$ . Ainsi 50% des données sont situées dans la boîte et autant à l'extérieur.
- La ligne verticale située à l'intérieur du rectangle est la médiane des données.
- Les pattes de la boîtes sont fixées à la plus petite donnée supérieure à  $Q_1 - 1.5(Q_3 - Q_1)$  pour la patte inférieure et la plus grande donnée inférieure à  $Q_3 + 1.5(Q_3 - Q_1)$  pour la patte supérieure.
- Les cercles au-delà des pattes correspondent aux données restantes souvent qualifiées d'aberrantes.

La commande R qui trace la boîte à moustaches est `boxplot`. Les données `iris` de type `data.frame`, correspondant aux mesures effectuées sur 50 fleurs pour trois espèces d'iris, fournissent un bon exemple d'application.

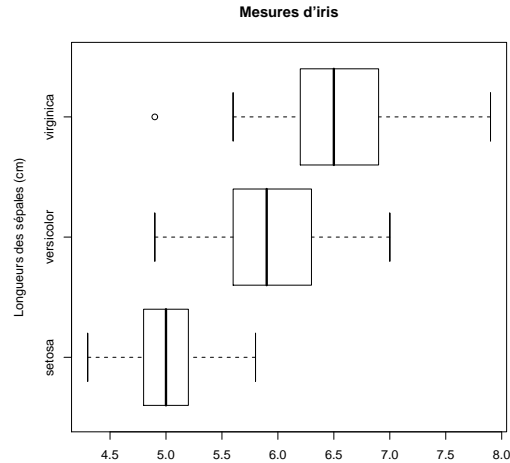


FIGURE 4 – Boîtes à moustaches

La figure FIG. 4 est le résultat des commandes suivantes :

```
> attach(iris)
> boxplot(Sepal.Length ~ Species, horizontal=TRUE, main="Mesures
+ d'iris", ylab="Longueur des sépales (cm)")
> detach(iris)
```

## 2.4 Histogramme

Pour une série statistique d'observations  $(x_j)_{j=1\dots n}$  (d'un caractère généralement continu) regroupées en  $k$  classes  $[a_i, a_{i+1}[$  d'effectif  $n_i$ , l'histogramme est un bon moyen de visualiser la répartition des données. Un histogramme est formé de rectangles accolés de bases  $[a_i, a_{i+1}[$  et de hauteurs  $\omega_i$  proportionnelle à la fréquence  $f_i = \frac{n_i}{n}$  si les classes ont la même amplitude et de hauteur proportionnelle à  $f_i/(a_{i+1} - a_i)$  sinon. On remarque qu'en choisissant,

$$w_i = \frac{f_i}{a_{i+1} - a_i} \quad \forall i = 1 \dots k \quad (1)$$

la somme des aires des rectangles est égale à 1 :

$$\sum_{i=1}^k \frac{f_i}{a_{i+1} - a_i} \times (a_{i+1} - a_i) = \sum_{i=1}^k f_i = 1.$$

En R, la commande `hist(x)` permet de tracer l'histogramme des données `x` avec les options classiques `main`, `xlim`, `xlab`, `ylab` ... et les options spécifiques :

- `breaks` : vecteur des délimitations entre classes (par défaut toutes les classes ont même amplitude).



- `freq` : si `freq=FALSE`, l'aire sous le graphe est égale à 1 (par défaut `freq=TRUE` et la hauteur d'un rectangle est égale à l'effectif de la classe).
- `nclass` : nombre de classes.
- `labels` : permet d'associer un nom à chaque rectangle.

Par ailleurs, la commande `h=hist()` permet aussi de conserver les caractéristiques de l'histogramme dans la liste `h`. Précisons certaines de ces caractéristiques :

- `breaks` :  $n + 1$  points de frontière entre classes.
- `counts` : vecteur des effectifs  $n_i$ .
- `mids` : vecteur des centres des classes.
- `density` : hauteur des rectangles pour l'estimation de densité.

Traçons, à titre d'exemple, l'histogramme d'un échantillon de taille 100 d'une loi  $\mathcal{N}(2,1)$ . Afin d'évaluer l'estimation de la densité fournie par l'histogramme, traçons aussi la densité d'une loi  $\mathcal{N}(2,1)$  en chacun des centres de classes. Le graphique obtenu à la figure FIG. 5 a été produit par les commandes suivantes :

```
> x=rnorm(100,2,1)
> h=hist(x,freq=FALSE,col="lightgrey",main="Histogramme d'un
+ échantillon gaussien N(2,1)")
> lines(h$mids,dnorm(h$mids,2,1),lwd=2)
```

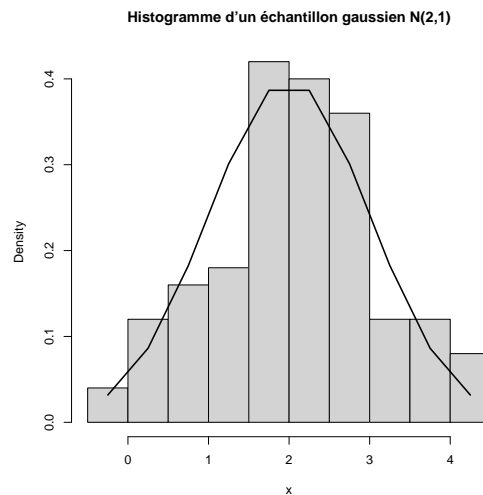


FIGURE 5 – Histogramme et densité d'un échantillon gaussien

## 3 Exercices d'application

### 3.1 Exercice 1

Les données `islands` de type `named vector` de R sont constituées des surfaces (en milliers de miles carrés) de 48 continents ou îles de plus de 10 000 miles carrés.

1. Tracez un histogramme des données.
2. Prenez le logarithme des données `islands` pour former un nouveau vecteur appelé `log-islands`. Tracez l'histogramme des données `log-islands`. Commentez.

3. Ouvrez un dispositif graphique à deux fenêtres. Puis, pour les données `log-islands`, tracez dans la première fenêtre un histogramme avec l'option `breaks='Sturges'` et dans l'autre fenêtre un histogramme avec l'option `breaks='Scott'`.
4. Construisez la boîte à moustaches des données `log-islands`.
5. Utilisez la commande graphique `dotchart` pour représenter les données `islands`. Faites la même chose pour les données par ordre croissant (à l'aide de la commande `sort`).
6. Quel(s) graphique(s) vous semble(nt) le(s) plus approprié(s) pour étudier la répartition des données ?

## 3.2 Exercice 2

Le `data.frame` `stackloss` de R rassemble 21 observations de 4 variables mesurées dans une entreprise transformant de l'ammoniac en acide nitrique.

1. Affichez les noms des 4 variables de `stackloss`.
2. A l'aide de la fonction `plot`, tracez tous les couples possibles de 2 variables.
3. Utilisez la commande `matplot` pour tracer sur un même graphique les 4 courbes des données `stackloss`. Ajoutez la légende à l'aide de la commande
 

```
> legend("topleft",legend=names(stackloss),lty=1:4,col=1:4)
```
4. Utilisez maintenant la commande `pairs` sur les données `stackloss`. Identifiez le(s) couple(s) de variables qui semblent reliées linéairement ou non linéairement.
5. Effectuez la régression linéaire entre les variables `Water.Temp` et `stack.loss` et stockez les informations relatives à la régression dans une liste nommée `info`. En R, la régression linéaire entre deux variables `var1` et `var2` s'effectue au moyen de la commande `lm` comme suit :

```
> # Calculs des éléments de la régression
> info=lm(var2~var1)
> # Affichage des résultats
> summary(info)
> # Représentation graphique
> plot(var1,var2)
> # Droite de régression
> abline(info$coefficients)
```

6. Testez la normalité des résidus `info$residuals` de la régression en utilisant les commandes `qqnorm`, `qqline`, `ks.test` et `shapiro.test`.