



# R POUR L'ANALYSE DE DONNÉES

## PRISE EN MAIN

### Ressources et cours consultés

- Udemy : *Programmer avec R en 5h en partant de zéro + Projet (2025)* de Mikail Altundas
- Udemy : *Programmer en R en 8h pour la Data Science de A à Z* d'Amandine Velt
- Wikipédia R
- Guide data.gouv R
- Bookdown — Introduction à R
- CRAN (site officiel)

11 septembre 2025

# Table des matières

<b>1</b>	<b>Installation de R et RStudio</b>	<b>5</b>
<b>2</b>	<b>Présentation de l'interface RStudio</b>	<b>6</b>
<b>3</b>	<b>Créer et exécuter un script R dans RStudio</b>	<b>6</b>
<b>4</b>	<b>Opérations Script.R</b>	<b>7</b>
<b>5</b>	<b>Vecteurs</b>	<b>9</b>
5.1	Syntaxe . . . . .	9
5.2	Opérations sur les vecteurs . . . . .	9
5.2.1	Création par répétition . . . . .	9
5.2.2	Création par suite ou séquence . . . . .	9
5.2.3	Création élément par élément . . . . .	9
5.2.4	Création avec nombres aléatoires . . . . .	10
5.2.5	Opérations arithmétiques . . . . .	10
5.2.6	Opérations logiques . . . . .	10
5.2.7	Modification . . . . .	10
5.2.8	Combinaison de vecteurs . . . . .	10
<b>6</b>	<b>Listes</b>	<b>11</b>
6.1	Syntaxe . . . . .	11
6.2	Opérations sur les listes . . . . .	11
6.2.1	Accès aux éléments . . . . .	11
6.2.2	Modification des éléments . . . . .	11
6.2.3	Listes imbriquées . . . . .	11

6.2.4	Fonctions utiles . . . . .	12
<b>7</b>	<b>Matrices</b>	<b>12</b>
7.1	Syntaxe . . . . .	12
7.2	Opérations sur les matrices . . . . .	12
7.2.1	Accès aux éléments . . . . .	12
7.2.2	Assemblage de colonnes ou de lignes . . . . .	13
7.2.3	Création d'une matrice diagonale . . . . .	13
7.2.4	Fonctions . . . . .	13
<b>8</b>	<b>Data Frames</b>	<b>13</b>
8.1	Syntaxe . . . . .	13
8.2	Opérations sur les data frames . . . . .	13
8.2.1	Accès aux éléments . . . . .	13
8.2.2	Modification du contenu . . . . .	14
8.2.3	Ajout ou suppression de colonnes / lignes . . . . .	14
8.2.4	Nommer les lignes et colonnes . . . . .	14
<b>9</b>	<b>Conditions, boucles et fonctions</b>	<b>14</b>
9.1	Structure conditionnelle : <code>if</code> . . . . .	14
9.2	La boucle <code>for</code> . . . . .	15
9.3	La boucle <code>while</code> . . . . .	16
9.4	La boucle <code>repeat</code> . . . . .	16
9.5	Créer ses propres fonctions . . . . .	16
<b>10</b>	<b>Fonctions pour la manipulation de données en R</b>	<b>17</b>
10.1	1. Fonctions de base sur les objets R . . . . .	17

10.2	2. Sélection et extraction . . . . .	18
10.3	3. Fusion . . . . .	18
10.4	4. Valeurs manquantes . . . . .	18
10.5	5. Application de fonctions (matrice, dataframe) . . . . .	19
10.6	6. Transformation et création . . . . .	19
10.7	7. Réorganisation . . . . .	19
10.8	8. Fonctions utiles de <code>dplyr</code> . . . . .	19
<b>11</b>	<b>Exemple pratique</b>	<b>20</b>
<b>12</b>	<b>Visualisation</b>	<b>21</b>
12.1	Graphiques statistiques classiques . . . . .	21
12.2	Graphiques multivariés et 3D . . . . .	22
12.3	Visualisation avec <code>ggplot2</code> . . . . .	22
12.4	Packages spécialisés pour la visualisation . . . . .	22
12.5	Personnalisation des axes . . . . .	22
12.6	Grilles et repères . . . . .	23
12.7	Légendes et textes . . . . .	23
12.8	Fenêtres graphiques multiples . . . . .	23
<b>13</b>	<b>Exemple pratique</b>	<b>24</b>
13.1	La fonction <code>plot()</code> . . . . .	24
13.2	La fonction <code>lines()</code> . . . . .	26
13.3	Ajout de repères et annotations graphiques . . . . .	28
13.4	Histogramme avec <code>hist()</code> . . . . .	30
13.5	Diagramme en barres avec <code>barplot()</code> . . . . .	32
13.6	Diagramme en boîte avec <code>boxplot()</code> . . . . .	33

13.7	Diagramme de dispersion avec <code>stripchart()</code> . . . . .	34
13.8	Autres graphiques . . . . .	36
13.9	Visualisation avancée en 2D et 3D . . . . .	38
<b>14</b>	<b>Quelques lois usuelles :</b>	<b>40</b>
14.1	Normale $\mathcal{N}(\mu, \sigma^2)$ . . . . .	41
14.2	Student $t_\nu$ (ddl $\nu$ ) . . . . .	41
14.3	Khi-deux $\chi_\nu^2$ . . . . .	41
14.4	Binomiale $\mathcal{B}(n, p)$ . . . . .	41
14.5	Poisson $\mathcal{P}(\lambda)$ . . . . .	41
14.6	Exponentielle $\text{Exp}(\lambda)$ . . . . .	42
14.7	Uniforme $\mathcal{U}(a, b)$ . . . . .	42
<b>15</b>	<b>Loi normale <math>\mathcal{N}(\mu, \sigma^2)</math></b>	<b>43</b>
15.1	Définition et paramètres . . . . .	43
15.2	Quand l'utiliser ? . . . . .	43
15.3	Exemples : . . . . .	43

# Introduction

L'**analyse de données** est l'ensemble des méthodes et techniques qui permettent de transformer des données brutes en informations compréhensibles et exploitables pour orienter la prise de décision.

Parmi les langages utilisés, **R** occupe une place centrale pour le traitement, l'analyse et la visualisation des données. Il est particulièrement apprécié pour ses capacités statistiques, sa richesse en bibliothèques spécialisées et ses outils graphiques puissants.

Ce support a pour objectif de présenter les bases du langage R, utiles pour manipuler des données, produire des analyses, et construire des visualisations claires et reproductibles.

## 1 Installation de R et RStudio

Pour utiliser le langage **R**, deux outils sont nécessaires :

- **R** : le moteur de calcul statistique.
- **RStudio** : une interface de développement (IDE) plus conviviale pour écrire et exécuter du code R.

### 1. Installer R

1. Aller sur le site officiel : <https://cran.r-project.org>
2. Choisir le système d'exploitation (*Windows*, *macOS* ou *Linux*).
3. Télécharger la dernière version stable de R.
4. Suivre les étapes d'installation par défaut.

### 2. Installer RStudio

1. Aller sur le site : <https://posit.co/download/rstudio-desktop>
2. Télécharger la version gratuite de **RStudio Desktop Open Source**.
3. Installer le logiciel comme n'importe quelle application.

### 3. Lancer RStudio

Une fois les deux outils installés, il suffit de lancer **RStudio**. Celui-ci utilise automatiquement l'interpréteur R installé en arrière-plan.

*Remarque* : Il est recommandé d'utiliser RStudio plutôt que l'interface de base de R, car il offre un environnement plus complet, avec éditeur de code, console, visualisation des données, historique des commandes, etc.

## 2 Présentation de l'interface RStudio

Une fois lancé, **RStudio** se présente sous la forme d'une interface divisée en quatre volets principaux :

- **Script (en haut à gauche)** : permet d'écrire et enregistrer du code R dans des fichiers.
- **Console (en bas à gauche)** : exécute les commandes R en temps réel.
- **Environnement / Historique (en haut à droite)** : affiche les objets en mémoire, les jeux de données, et l'historique des commandes.
- **Plots / Packages / Fichiers (en bas à droite)** : permet d'afficher les graphiques, gérer les packages, explorer les fichiers, etc.

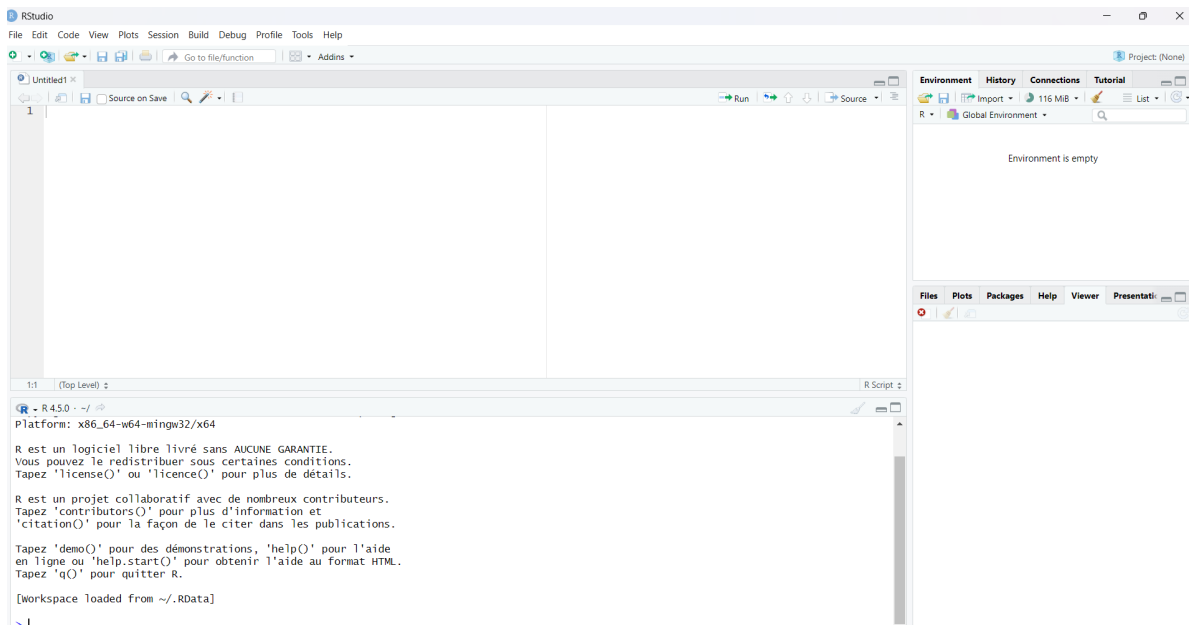


FIGURE 1 – Interface RStudio au démarrage

## 3 Créer et exécuter un script R dans RStudio

Dans RStudio, l'écriture du code se fait dans un fichier appelé **script R**, qui peut être sauvegardé et réutilisé.

## 1. Créer un nouveau fichier R

Pour créer un nouveau script R :

- Aller dans le menu **File > New File > R Script**.
- Un nouvel onglet s'ouvre dans la partie supérieure gauche de RStudio.
- Ce fichier peut être sauvegardé avec l'extension **.R**.

## 2. Types de fichiers courants

RStudio permet de créer différents types de fichiers selon les besoins :

- **R Script (.R)** : pour écrire du code R classique.
- **R Markdown (.Rmd)** : pour combiner code, texte et rendu dynamique.
- **Shiny App** : pour créer des applications interactives en R.
- **Quarto (.qmd)** : alternative moderne à R Markdown pour des rapports interactifs.

## 3. Exécuter du code R

Une fois un script ouvert, il est possible d'exécuter une ou plusieurs lignes de code de deux manières :

- **Ctrl + Entrée** (ou **Cmd + Entrée** sur Mac) : exécute la ligne sélectionnée.
- Cliquer sur le bouton **Run** en haut à droite de l'éditeur.

Le résultat de l'exécution s'affiche dans la console, en bas à gauche de l'interface.

## 4 Opérations Script.R

R peut être utilisé pour effectuer des opérations algébriques élémentaires. On peut également enregistrer des valeurs dans des variables et travailler ensuite avec celles-ci.

### 1. Calculs simples

Voici quelques exemples de calculs réalisés directement dans la console R :

```
1 3 + 2
2 3 - 2
3 3 * 2
```



```

4 3 / 2
5 3^2
6
7 # Plusieurs instructions sur une seule ligne
8 3 + 2 ; 3 - 2 ; 3 * 2 ; 3^2

```

## 2. Affectations et expressions plus complexes

On peut stocker des valeurs dans des variables (ou scalaires) avec le symbole = ou <-. Ensuite, ces variables peuvent être combinées dans des expressions :

```

1 x = 2
2 a = 2
3 b = 3
4 c = 4
5
6 x + c
7 b * x + c
8 a * x^2 + b * x + c
9 b / (x + c)^a
10 x ^ (a + b + c)

```

## 3. Opérations logiqueset de comparaisons

R permet également de manipuler des valeurs logiques :

```

1 TRUE
2 FALSE
3 Exemple
4 3 > 2      # TRUE
5 4 == 5     # FALSE
6 5 != 2     # TRUE
7
8 == : égalité
9 != : différence
10 > : supérieur
11 < : inférieur
12 & : et logique
13 | : ou logique

```

## 5 Vecteurs

Les **vecteurs** sont les structures de base de R. Ce sont des séquences ordonnées d'éléments *de même type* : numériques, chaînes de caractères ou logiques.

### 5.1 Syntaxe

```
nom_vecteur <- c(valeur1, valeur2, ..., valeurN)
```

Exemple :

```
1 vec1 = c(2.8, 2.4, 2.1, 3.6, 2.8)      # vecteur numérique
2 vec2 = c("rouge", "vert", "vert", "jaune") # vecteur de chaînes de caractères
3 vec3 = c(TRUE, TRUE, FALSE, FALSE, FALSE) # vecteur logique (booléen)
```

### 5.2 Opérations sur les vecteurs

#### 5.2.1 Création par répétition

```
1 rep(4, 3)          # répète la valeur 4 trois fois -> 4 4 4
2 vec4 = rep(vec1, 2) # répète tout le vecteur vec1 deux fois
3 vec5 = rep(vec1, c(2, 1, 3, 3, 2)) # répète chaque élément de vec1
```

#### 5.2.2 Création par suite ou séquence

```
1 vec6 = 1:10          # crée une séquence d'entiers de 1 à 10
2 vec7 = seq(3, 5, by = 0.2) # crée une séquence de 3 à 5 avec un pas de 0.2
```

#### 5.2.3 Création élément par élément

```
1 vec8 = numeric()      # crée un vecteur numérique vide
2 vec8[1] = 41.8         # affecte la valeur 41.8 à la 1ère position
3 vec8[2] = -0.3         # affecte la valeur -0.3 à la 2ème position
4 vec8[3] = 92           # affecte la valeur 92 à la 3ème position
5 vec8                   # affiche le vecteur rempli
```

### 5.2.4 Création avec nombres aléatoires

```
1 set.seed(123) # fixe le hasard
2 moyenne_de_la_classe <- sample(1:20, 20) # 20 nombres tirés entre 1 et 20
3 moyenne_de_la_classe # affiche le résultat
```

### 5.2.5 Opérations arithmétiques

```
1 vec1 = c(1, 2, 3) # vecteur 1
2 vec2 = c(4, 5, 6) # vecteur 2
3 vec1 + vec2 # addition élément par élément
4 vec1 - vec2 # soustraction élément par élément
5 vec1 * vec2 # multiplication élément par élément
6 vec2 / vec1 # division élément par élément
7 vec1 ^ 2 # puissance (carré) de chaque élément
```

### 5.2.6 Opérations logiques

```
1 v = c(4, 7, 9, 3, 6) # vecteur numérique
2 v > 5 # test supérieur à 5 (renvoie TRUE/FALSE)
3 v == 3 # test égal à 3
4 v != 7 # test différent de 7
```

### 5.2.7 Modification

```
1 v = c(10, 20, 30, 40, 50) # vecteur initial
2 v[2] = 99 # remplace la 2ème valeur par 99
3 v[c(1, 3)] = 0 # remplace les 1ère et 3ème valeurs par 0
4 v # affiche le vecteur modifié
```

### 5.2.8 Combinaison de vecteurs

```
1 a = c(1, 2, 3) # vecteur a
2 b = c(4, 5) # vecteur b
3 c1 = c(a, b) # concatène a et b
4 c2 = a + 10 # ajoute 10 à chaque élément de a
5 c3 = a * b # multiplication élément par élément
```

## 6 Listes

Une **liste** en R est une structure de données flexible qui peut contenir des objets de types variés : vecteurs, matrices, fonctions, data frames ou même d'autres listes.

### 6.1 Syntaxe

```
nom_liste <- list(nom1 = objet1, nom2 = objet2, ...)
```

Exemple :

```
1 ma_liste <- list(  
2   nom = "Alice",           # élément texte (chaîne de caractères)  
3   notes = c(15, 18, 13),    # vecteur numérique  
4   moyenne = mean(c(15, 18, 13)), # valeur calculée (moyenne des notes)  
5   validation = TRUE        # valeur logique  
6 )  
7 ma_liste                   # affiche la liste complète
```

### 6.2 Opérations sur les listes

#### 6.2.1 Accès aux éléments

```
1 ma_liste$notes      # Accès par le nom  
2 ma_liste[[2]]       # Accès par position  
3 ma_liste[2]         # Renvoie une sous-liste  
4 names(ma_liste)     # Noms des composantes
```

#### 6.2.2 Modification des éléments

```
1 ma_liste$remarque <- "Très bon travail" # Ajout  
2 ma_liste$validation <- FALSE           # Modification  
3 ma_liste$nom <- NULL                    # Suppression
```

#### 6.2.3 Listes imbriquées

```
1 liste2 <- list(  
2   nom = "Projet",          # élément texte  
3   contenu = list()         # sous-liste
```

```

4     data = c(1, 2, 3),           # vecteur numérique
5     titre = "Analyse 2025"      # texte
6 )
7 )
8 liste2$contenu$data            # accès via le nom
9 liste2[[2]][["data"]]          # accès via les indices de liste

```

### 6.2.4 Fonctions utiles

```

1 length(ma_liste)               # Nombre d'éléments
2 str(ma_liste)                  # Structure
3 names(ma_liste)                # Noms des composantes
4 is.list(ma_liste)              # Est-ce une liste ?

```

## 7 Matrices

Une **matrice** en R est une structure bidimensionnelle (lignes  $\times$  colonnes) contenant uniquement des éléments du même type. Elle est utile pour les calculs statistiques, algébriques, ou pour représenter des tableaux numériques.

### 7.1 Syntaxe

```
nom_matrice <- matrix(vecteur, nrow = x, ncol = y)
```

Exemple :

```

1 mat1 <- matrix(1:10, ncol = 5)  # matrice avec 10 nombres sur 5 colonnes
2 mat1                            # affiche la matrice

```

### 7.2 Opérations sur les matrices

#### 7.2.1 Accès aux éléments

```

1 mat1[, 3]                       # colonne 3
2 mat1[2, ]                       # ligne 2
3 mat1[1, c(2, 4)]                # éléments aux colonnes 2 et 4 de la ligne 1

```

### 7.2.2 Assemblage de colonnes ou de lignes

```
1 cbind(1:3, 4:6)           # concaténation colonne
2 rbind(c(1,2,3), c(4,5,6)) # concaténation ligne
```

### 7.2.3 Création d'une matrice diagonale

```
1 diag(c(2, 1, 5))  # matrice diagonale (2,1,5 sur la diagonale)
```

### 7.2.4 Fonctions

```
1 rowSums(mat)  # somme par ligne
2 colMeans(mat) # moyenne par colonne
3 t(mat)        # transposée
4 dim(mat)      # dimensions (nb lignes, nb colonnes)
```

## 8 Data Frames

Un **data frame** est une table à deux dimensions, chaque colonne pouvant avoir un type différent (numérique, caractère, logique, etc.). C'est la structure la plus utilisée pour manipuler les jeux de données en R.

### 8.1 Syntaxe

```
df <- data.frame(colonne1 = ..., colonne2 = ..., ...)
```

Exemple :

```
1 df <- data.frame(
2   nom = c("Alice", "Bob", "Chloe"),  # colonne texte
3   age = c(25, 30, 22),               # colonne numérique
4   etudiant = c(TRUE, FALSE, TRUE)   # colonne logique
5 )
```

### 8.2 Opérations sur les data frames

#### 8.2.1 Accès aux éléments

```

1 df$age           # accès à la colonne "age"
2 df[1, ]         # sélection de la 1ère ligne
3 df[, "etudiant"] # accès à la colonne "etudiant"

```

### 8.2.2 Modification du contenu

```

1 df$age[1] <- 26      # modifie la 1ère valeur de la colonne "age"
2 df$etudiant[3] <- FALSE # change la 3ème valeur de la colonne "etudiant"

```

### 8.2.3 Ajout ou suppression de colonnes / lignes

```

1 df$ville <- c("Paris", "Lyon", "Caen")      # ajoute une colonne "ville"
2 df[4, ] <- list("David", 28, TRUE, "Nice")  # ajoute une 4ème ligne complète
3 df$etudiant <- NULL                        # supprime la colonne "etudiant"

```

### 8.2.4 Nommer les lignes et colonnes

```

1 colnames(df) <- c("Nom", "ge", "Étudiant", "Ville") # renomme les colonnes
2 rownames(df) <- c("A", "B", "C", "D")              # renomme les lignes

```

## 9 Conditions, boucles et fonctions

Dans cette partie, nous allons découvrir les bases de la logique en programmation avec R : tester des conditions, répéter des actions, et créer des fonctions réutilisables.

Ce sont des outils essentiels pour rendre vos analyses automatiques, puissantes et intelligentes.

### 9.1 Structure conditionnelle : if

Les instructions conditionnelles permettent d'exécuter du code seulement si une condition est vraie.

**Syntaxe :**

```

1 if (condition) {
2   # instructions si vrai
3 } else {

```

```
4   # instructions si faux
5 }
```

### Exemples :

```
1 x <- 10
2 if (x > 5) {
3   print("x est supérieur à 5")      # cas où x est > 5
4 } else {
5   print("x est inférieur ou égal à 5") # sinon
6 }
```

### Avec else if :

```
1 x <- 7
2 if (x < 5) {
3   print("Petit")    # si x est inférieur à 5
4 } else if (x == 5) {
5   print("Moyen")    # si x est égal à 5
6 } else {
7   print("Grand")    # sinon (x > 5)
8 }
```

## 9.2 La boucle for

La boucle `for` est utilisée pour répéter une instruction un certain nombre de fois, pour chaque élément d'une séquence.

### Syntaxe :

```
1 for (variable in séquence) {
2   # instructions à répéter
3 }
```

### Exemples :

```
1 for (i in 1:5) {
2   print(i^2)  # affiche le carré de i
3 }
4
5 for (jour in c("lundi", "mardi", "mercredi")) {
6   print(paste("Aujourd'hui, c'est", jour)) # affiche le jour
7 }
```



## 9.3 La boucle while

La boucle `while` répète une instruction tant qu'une condition reste vraie.

**Syntaxe :**

```
1 while (condition) {  
2   # instructions à répéter  
3 }
```

**Exemple :**

```
1 x <- 1  
2 while (x < 10) {  
3   print(x)      # affiche la valeur de x  
4   x <- x + 2    # incrémente x de 2  
5 }
```

## 9.4 La boucle repeat

La boucle `repeat` répète un bloc d'instructions indéfiniment, jusqu'à rencontrer une instruction `break`.

**Syntaxe :**

```
1 repeat {  
2   # instructions  
3   if (condition) break  
4 }
```

**Exemple :**

```
1 x <- 1  
2 repeat {  
3   print(x)      # affiche la valeur de x  
4   x <- x + 1    # incrémente x  
5   if (x > 5) break # stoppe la boucle si x > 5  
6 }
```

## 9.5 Créer ses propres fonctions

Une fonction est un bloc de code qui peut être réutilisé plusieurs fois. Elle prend des paramètres en entrée et retourne un résultat.

Syntaxe :

```
1 nom_fonction <- function(argument1, argument2, ...) {  
2   # instructions  
3   return(résultat)  
4 }
```

Exemple simple :

```
1 somme_carre <- function(x) {  
2   return(x^2 + x)    # renvoie x**2 + x  
3 }  
4 somme_carre(3)        # appel de la fonction -> 12
```

## 10 Fonctions pour la manipulation de données en R

Ce chapitre regroupe les fonctions incontournables pour explorer, transformer, trier, résumer ou nettoyer des données en R et les fonctions applicables aux vecteurs.

### 10.1 1. Fonctions de base sur les objets R

```
1 #vecteurs  
2 v = c(10, 20, 30, 40, 50)  
3 sum(v)                # la somme de toutes les valeurs  
4 mean(v)               # la moyenne des valeurs  
5 min(v)                # le minimum des valeurs  
6 max(v)                # le maximum des valeurs  
7 length(v)             # le nombre de valeur  
8 sort(v)               # trier par croissant  
9 rev(v)                # reverser la liste  
10 range(v)              # retourner le minimum et le maximum  
11 unique(v)             # les vaeurs uniques  
12 duplicated(v)         # TRUE pour les valeurs duppliquées  
13 which(v > 30)          # retourner les valeurs supérieurs à 30  
14 any(v > 60)           # est ce qu'il y'a au moins un vrai  
15 all(v > 5)            # est ce que c'est vrai pour toutes les valeurs  
16 cumsum(v)             # somme cumulée des valeurs  
17 cumprod(v)            # le produit cumulé des valeurs  
18 diff(v)               # la difference entre les valeurs  
19 rep(v, times = 3)     # Répéter v trois fois  
20 seq(1, 10, by = 2)    # Séquence  
21 which.max(v)          # Position du max  
22 which(v > 5)          # Positions vérifiant une condition
```

```

23 gl(2, 5) # Générer un facteur
24 replicate(3, mean(rnorm(10))) # Répétitions
25 append(vec, 100, after = 2) # Ajouter élément
26
27 # dataframe
28 head(df) # Premières lignes
29 tail(df) # Dernières lignes
30 dim(df) # Dimensions (lignes, colonnes)
31 nrow(df) # Nombre de lignes
32 ncol(df) # Nombre de colonnes
33 length(vec) # Longueur d'un vecteur
34 names(df) # Noms des colonnes
35 colnames(df) # Colonnes
36 rownames(df) # Lignes
37 class(df) # Classe de l'objet
38 str(df) # Structure (type et aperçu)
39 summary(df) # Résumé statistique

```

## 10.2 2. Sélection et extraction

```

1 df$colonne # Extraire une colonne
2 df[1, ] # Extraire la première ligne
3 df[, 2] # Extraire la deuxième colonne
4 df[1:3, c("Age", "Sexe")] # Sous-table
5 subset(df, Age > 25) # Filtrer les lignes

```

## 10.3 3. Fusion

```

1 cbind(v1, v2) # Coller des vecteurs en colonnes
2 rbind(v1, v2) # Coller des vecteurs en lignes
3 merge(df1, df2, by = "ID") # Jointure par colonne

```

## 10.4 4. Valeurs manquantes

```

1 is.na(x) # vérifier, cellule par cellule, si la valeur NA
2 complete.cases(df) # repérer les lignes complètes (sans NA)
3 anyNA(df) # vérifier existe au moins un NA
4 df = df[complete.cases(df),] # garder uniquement les lignes complètes

```

## 10.5 5. Application de fonctions (matrice, dataframe)

```
1 apply(mat, 1, mean)           # Moyenne par ligne
2 apply(mat, 2, sd)             # Écart-type par colonne
3 lapply(df, class)             # Classe de chaque colonne
4 sapply(df, mean)              # Moyenne par colonne
5 vapply(df, is.numeric, TRUE) # type numerique?
6 tapply(df$Age, df$Sexe, mean) # Moyenne d'age par sexe
7
8 by(df, df$Sexe, summary)      # Résumé par groupe
```

## 10.6 6. Transformation et création

```
1 transform(df, Age2 = Age + 10) # Nouvelle colonne
2 within(df, { Note = Age * 2 }) # Modification interne
3 cut(df$Age, breaks = 3)        # Discrétiser
4 rank(c(10, 5, 3))              # Rang
```

## 10.7 7. Réorganisation

```
1 stack(df)                      # Transforme ton data frame en format long (values + ind)
2 unstack(stack(df))             # rend le même df de départ
3 split(df$Age, df$Sexe)         # Découper par groupe
```

## 10.8 8. Fonctions utiles de dplyr

```
1 library(dplyr)
2
3 select(df, Age, Sexe)          # Sélection de colonnes
4 filter(df, Age > 20)           # Filtrage de lignes
5 mutate(df, Age2 = Age + 2)     # Nouvelle colonne
6 arrange(df, desc(Age))         # Tri décroissant
7 rename(df, NouveauNom = Age)   # Renommer une colonne
8 group_by(df, Sexe)             # Grouper
9 distinct(df)                  # Supprimer doublons
10 count(df, Sexe)               # Compter par groupe
11 relocate(df, Sexe, .before = Age) # Réorganiser
```

## 11 Exemple pratique

Voici un exemple de manipulation de données.

```
1 # --- 0. Création dun data frame ---
2 employes <- data.frame(
3   Nom      = c("Alice", "Bob", "Claire", "David", "Eva", "François"),
4   Age      = c(25, 34, 28, NA, 45, 28),
5   Sexe     = c("F", "H", "F", "H", "F", "H"),
6   Salaire  = c(2000, 2500, 2200, 2700, NA, 2500),
7   Département = c("RH", "IT", "IT", "Comptabilité", "RH", "IT")
8 )
9
10 # --- 1. Exploration de base ---
11 head(employes)                # 6 premières lignes
12 str(employes)                 # structure du data frame
13 summary(employes)            # résumé statistique
14 names(employes)              # noms des colonnes
15 nrow(employes)               # nombre de lignes
16 ncol(employes)               # nombre de colonnes
17
18 # --- 2. Sélection & extraction ---
19 employes$Nom                  # extraire une colonne
20 employes[1:3, c("Nom", "Salaire")] # lignes 1-3 et colonnes choisies
21 subset(employes, Age > 30)    # employés > 30 ans
22 employes[is.na(employes$Age), ] # lignes avec Age manquant
23
24 # --- 3. Nettoyage ---
25 employes <- employes[complete.cases(employes), ] # supprime les NA
26 unique(employes$Age)                # âges distincts
27 duplicated(employes)                # détecte les doublons
28
29 # --- 4. Fusion de data frames ---
30 infos_suppl <- data.frame(
31   Nom      = c("Alice", "Bob", "Claire"),
32   Anciennete = c(3, 5, 2)
33 )
34 employes <- merge(employes, infos_suppl, by = "Nom", all.x = TRUE) #fusion
35
36 # --- 5. Nouvelles colonnes ---
37 employes$Prime <- employes$Salaire * 0.1 # prime = 10% du salaire
38 employes <- transform(employes, Bonus = Salaire * 0.05) # autre méthode
39
40 # --- 6. Fonctions apply ---
41 apply(employes[, c("Salaire", "Prime", "Bonus")], 2, mean) # moyenne par colonne
42
43 # --- 7. Statistiques par groupes ---
```

```

44 tapply(employees$Salaire, employees$Sexe, mean)           # moyenne par sexe
45 aggregate(Salaire ~ Département, data = employees, mean) # moyenne par département
46 by(employees$Salaire, employees$Département, summary)    # résumé par département
47
48 # --- 8. Tri des données ---
49 employees[order(employees$Age), ]                         # tri croissant par âge
50 employees[order(-employees$Salaire), ]                   # tri décroissant par salaire
51
52 # --- 9. Manipulations avec dplyr ---
53 library(dplyr)
54 employees %>%
55   select(Nom, Salaire) %>%                                # sélectionner colonnes
56   filter(Salaire > 2300) %>%                               # filtrer
57   mutate(Net = Salaire - 150) %>%                         # nouvelle colonne Net
58   arrange(desc(Net)) %>%                                  # tri décroissant
59   group_by(Sexe) %>%                                       # regroupement par sexe
60   summarise(Salaire_moyen = mean(Salaire), Max = max(Salaire)) # stats par sexe
61
62 # --- 10. Pipe avancé ---
63 employees %>%
64   group_by(Département) %>%                                # regroupement par dép.
65   filter(Salaire > 2000) %>%                               # filtrer
66   summarise(Salaire_moyen = mean(Salaire))                # moyenne par groupe
67
68 # --- 11. Fonctions complémentaires ---
69 which(employees$Salaire == max(employees$Salaire)) # indice du salaire max
70 rank(employees$Salaire)                             # rang des salaires
71 seq(18, 60, by = 6)                                 # séquence
72 rep(employees$Nom[1], times = 3)                     # répète un nom
73 append(employees$Nom, "Zara", after = 2)             # insère "Zara"

```

## 12 Visualisation

R propose une grande variété de fonctions pour représenter visuellement les données. Ces fonctions peuvent être regroupées selon leur origine (base R ou packages) et leur usage (statistique, multivarié, interactif, etc.).

### 12.1 Graphiques statistiques classiques

- `plot()` : nuages de points, courbes...
- `hist()` : histogramme pour la distribution des valeurs.
- `boxplot()` : visualisation des quartiles et des valeurs extrêmes.

- `barplot()` : diagramme en barres.
- `pie()` : camembert.
- `dotchart()`, `stripchart()` : valeur individuelles en ligne ou en trait.
- `stem()` : graphique tige et feuille.
- `density()` : estimation de densité.
- `rug()` : ajoute des tics pour les observations.
- `pairs()` : nuages de points pour toutes les paires de variables.
- `coplot()` : nuages conditionnels.
- `mosaicplot()`, `assocplot()`, `fourfoldplot()` : pour les tableaux de contingence.

## 12.2 Graphiques multivariés et 3D

- `image()`, `heatmap()`, `filled.contour()` : représentation de matrices de valeurs.
- `contour()` : courbes de niveaux.
- `persp()` : représentation 3D.
- `symbols()` : graphiques à bulles, cercles proportionnels.

## 12.3 Visualisation avec ggplot2

- `ggplot()` : initialise le graphe.
- `aes()` : définit les esthétiques (x, y, couleur...).
- `geom_point()`, `geom_line()`, `geom_bar()`, `geom_col()`, `geom_boxplot()`, `geom_histogram()`, `geom_violin()` : types de graphes.
- `geom_density()`, `geom_rug()`, `geom_smooth()`, `geom_text()` : couches supplémentaires.
- `facet_wrap()`, `facet_grid()` : créer plusieurs graphes selon une variable.
- `labs()`, `theme()`, `coord_flip()`, `coord_polar()` : personnalisation.

## 12.4 Packages spécialisés pour la visualisation

- `lattice` : `xyplot()`, `bwplot()`, pour visualisation conditionnelle.
- `plotly`, `highcharter`, `echarts4r` : visualisations interactives.
- `corrplot`, `GGally`, `ggbridges`, `ggraph` : pour corrélations, graphes, densités.
- `RColorBrewer`, `viridis` : palettes de couleurs.
- `patchwork`, `cowplot`, `gridExtra` : mise en page multiple de graphes `ggplot2`.

## 12.5 Personnalisation des axes

- `xaxt = "n"`, `yaxt = "n"` : Empêche l’affichage automatique des axes (X ou Y).

- **axis(side, at, labels)** : Ajoute un axe personnalisé.
  - **side** = 1 (bas), 2 (gauche), 3 (haut), 4 (droite)
  - **at** : vecteur des positions des graduations
  - **labels** : étiquettes à afficher
  - **tick** : afficher les ticks (TRUE/FALSE)
  - **col.axis**, **cex.axis**, **font.axis** : couleur, taille, police
- **las** : Orientation du texte des axes (0 = parallèle, 1 = horizontal, 2 = perpendiculaire, 3 = vertical).
- **tck** : Longueur des ticks (0 = invisibles, négatif = vers l'intérieur).
- **mgp = c(titre, label, ligne)** : Contrôle l'espacement entre titre, label et ligne des axes.

## 12.6 Grilles et repères

- **grid(nx, ny)** : Affiche une grille par défaut dans le fond du graphique.
  - **nx**, **ny** : nombre de lignes verticales et horizontales
  - **col**, **lty**, **lwd** : couleur, type, épaisseur de ligne
- **abline()** : Ajoute des lignes droites de repère :
  - **abline(h = valeur)** : ligne horizontale
  - **abline(v = valeur)** : ligne verticale
  - **abline(a, b)** : droite de régression  $y = a + bx$
- **segments(x0, y0, x1, y1)** : Trace un segment.
- **arrows(x0, y0, x1, y1)** : Trace une flèche.

## 12.7 Légendes et textes

- **legend(x, y, legend)** : Ajoute une légende dans le graphique.
  - **x**, **y** : position ou mots-clés : "topright", "bottomleft", etc.
  - **legend** : vecteur des intitulés
  - **pch**, **lty**, **col**, **cex** : symboles et apparence
  - **bty** = "n" : supprime la boîte
- **text(x, y, labels)** : Affiche un texte à la position donnée.
- **mtext()** : Texte dans la marge.
- **title()** : Ajoute ou modifie le titre ou les axes.

## 12.8 Fenêtres graphiques multiples

- **par(mfrow = c(n, m))** : Affiche plusieurs graphiques dans une fenêtre, ligne par ligne.
- **par(mfcol = c(n, m))** : Affiche plusieurs graphiques colonne par colonne.



- `layout(mat)` : Organisation personnalisée des graphiques.
- `par(mar = c(b, l, t, r))` : Marges internes (bottom, left, top, right).
- `par(oma = c(b, l, t, r))` : Marges externes.
- `par(bg = "color")` : Couleur de fond de la fenêtre graphique.

## 13 Exemple pratique

### 13.1 La fonction `plot()`

La fonction `plot()` est la fonction graphique générique de R. Elle permet de tracer un nuage de points, une courbe ou un graphique personnalisé à partir de deux vecteurs ou d'un objet graphique.

#### Syntaxe

`plot(x, y = NULL, type = "p", ...)`

- `x, y` : vecteurs numériques ou objets graphiques.
- `type` : type de tracé :
  - `"p"` : points (valeur par défaut)
  - `"l"` : lignes
  - `"b"` : points + lignes
  - `"o"` : points + lignes superposés
  - `"h"` : segments verticaux (type histogramme)
  - `"s", "S"` : marches (lignes en escalier)
  - `"n"` : ne trace rien (utile pour personnaliser à la main)
- `xlab, ylab` : noms des axes.
- `xlim, ylim` : limites des axes.
- `main, sub` : titre principal et sous-titre.
- `col, pch, lty, lwd` : couleur, symbole, type et épaisseur de ligne.
- `cex, cex.axis, cex.lab, cex.main` : tailles des éléments.
- `axes = FALSE, xaxt = "n", yaxt = "n"` : suppression des axes.
- `asp` : rapport d'aspect (utile pour les cercles).
- `frame.plot = FALSE, bty = "n"` : désactivation de la boîte autour du graphe.

Des fonctions complémentaires permettent d'enrichir le graphique après un appel à `plot()` :

- `lines()`, `points()` : ajouter des courbes ou points supplémentaires.
- `abline()`, `segments()`, `arrows()` : ajouter des repères ou flèches.
- `text()`, `legend()`, `mtext()`, `title()` : annotation, titres et légendes.
- `grid()`, `box()`, `axis()` : personnalisation des axes et du fond.
- `identify()`, `locator()` : interaction avec la souris.

## Exemple : Visualisation avec plot()

```
1 # Génère une suite de 10 valeurs de 1 à 10
2 x <- 1:10
3
4 # Génère 10 valeurs aléatoires selon une loi normale
5 y <- rnorm(10)
6
7 # Simple tracé
8 plot(x, y)
9
10 # Tracé avec personnalisation avancée
11 plot(x, y,
12       type = "b",
13       col = "darkgreen",
14       pch = 17,
15       lwd = 2,
16       main = "Évolution aléatoire",
17       xlab = "Indice",
18       ylab = "Valeur",
19       xlim = c(0, 12),
20       ylim = c(-3, 3),
21       cex = 1.5,
22       cex.main = 1.2,
23       las = 1,
24       bty = "l",
25       xaxt = "n")
26
27 # Axe X personnalisé
28 axis(1,
29       at = 1:10,
30       labels = paste("V", 1:10),
31       col.axis = "blue",
32       cex.axis = 0.8)
33
34 # Grille et ligne horizontale
35 grid()
36 abline(h = 0, col = "gray", lty = 2)
37
38 # Légende
39 legend("topright",
40       legend = "Série",
41       col = "darkgreen",
42       pch = 17,
43       bty = "n")
```

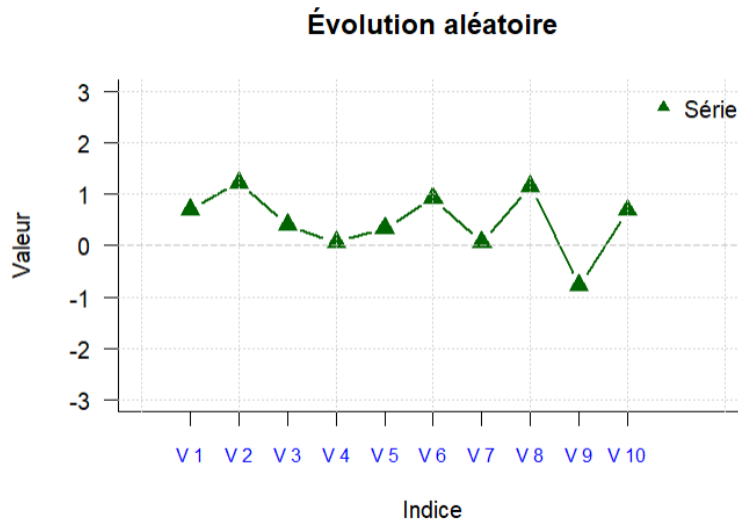


FIGURE 2 – Visualisation enrichie avec `plot()`

Ce graphique combine des points reliés, une personnalisation des axes, une grille, un repère horizontal, et une légende pour une visualisation claire et structurée.

## 13.2 La fonction `lines()`

La fonction `lines()` permet d'ajouter des lignes sur un graphique existant, souvent pour superposer plusieurs séries de données ou enrichir la courbe d'un tracé.

### Syntaxe

`lines(x, y, ...)`

- `x, y` : vecteurs numériques représentant les coordonnées des points à relier.
- `col` : couleur de la ligne.
- `lty` : type de ligne (1 = continue, 2 = pointillée, etc.).
- `lwd` : épaisseur de la ligne.
- `type` : permet aussi d'ajouter des symboles avec "b", "o", etc.

**Important** : `lines()` ne crée pas de graphique seule — elle doit suivre un `plot()`.

### Exemple complet

```
1 x <- 1:10
```

```

2 y <- rnorm(10)
3
4 # Génère une série secondaire à superposer
5 y2 <- rnorm(10, mean = 0.5)
6
7 # Tracé du graphique principal
8 plot(x, y,
9       type = "b",
10      col = "darkgreen",
11      pch = 17,
12      lwd = 2,
13      main = "Comparaison de deux séries",
14      xlab = "Indice",
15      ylab = "Valeur",
16      xlim = c(0, 12),
17      ylim = range(c(y, y2)),
18      cex = 1.5,
19      cex.main = 1.2,
20      las = 1,
21      bty = "l",
22      xaxt = "n")
23
24 # Axe X personnalisé
25 axis(1, at = 1:10, labels = paste("V", 1:10), col.axis = "blue", cex.axis = 0.8)
26
27 # Ajoute une grille et un repère
28 grid()
29 abline(h = 0, col = "gray", lty = 2)
30
31 # Ajout de la deuxième série avec lines()
32 lines(x, y2,           # Coordonnées des points à relier
33       type = "b",      # "b" = both points + lignes
34       col = "orange",  # Couleur orange pour différencier la série
35       pch = 19,        # Forme du point : rond plein
36       lty = 3,         # Type de ligne : pointillée
37       lwd = 2)         # Épaisseur de la ligne
38
39 # Ajoute une légende pour les deux séries
40 legend("topright",
41       legend = c("Série 1", "Série 2"),
42       col = c("darkgreen", "orange"),
43       pch = c(17, 19),
44       lty = c(1, 3),
45       bty = "n")
46 \end{

```

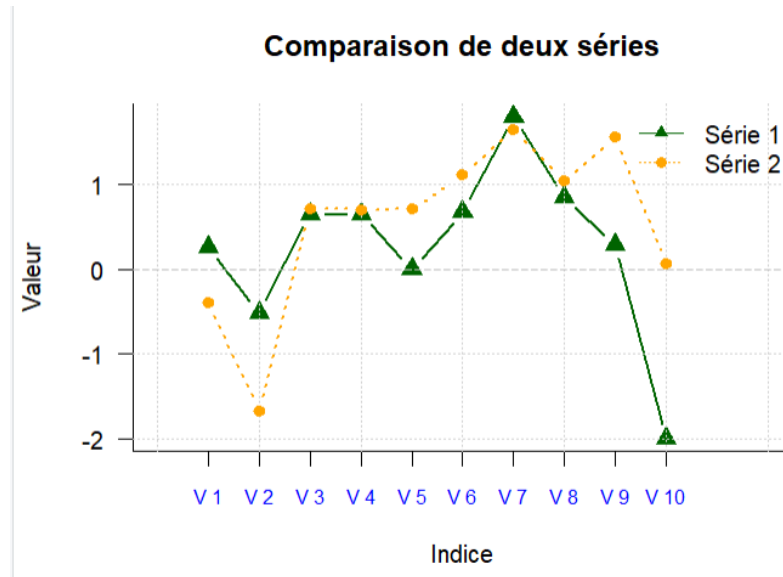


FIGURE 3 – Superposition de deux séries avec `plot()` et `lines()`

### 13.3 Ajout de repères et annotations graphiques

En plus de tracer des courbes, **R** permet d'ajouter manuellement des repères, des annotations, des flèches ou des points pour enrichir la lisibilité d'un graphique. Voici les principales fonctions utilisées à cet effet :

- `abline()` : ajoute une droite (horizontale, verticale ou oblique).
- `points()` : ajoute manuellement un ou plusieurs points.
- `segments()` : trace un segment entre deux points.
- `arrows()` : ajoute une flèche entre deux coordonnées (utile pour pointer ou annoter).

#### Exemple combiné

```

1 # Génération de données
2 x <- 1:10
3 y <- rnorm(10)
4
5 # Tracé principal avec courbe et points
6 plot(x, y,
7       type = "b",           # both = points + lignes
8       col = "blue",         # couleur bleue
9       pch = 16,             # points ronds pleins
10      main = "Exemples de repères et annotations",
11      xlab = "Index",
12      ylab = "Valeur")
13

```

```

14 # Ajout d'une droite horizontale
15 abline(h = 0,
16        col = "gray",          # couleur grise
17        lty = 2)              # ligne pointillée
18
19 # Ajout dun point spécifique
20 points(5, 1.5,
21        pch = 17,              # triangle plein
22        col = "darkgreen",
23        cex = 1.5)
24
25 # Ajout dun segment
26 segments(x0 = 2, y0 = -1, x1 = 8, y1 = 1.5,
27          col = "orange",
28          lwd = 2)
29
30 # Ajout d'une flèche
31 arrows(x0 = 5, y0 = -1, x1 = 5, y1 = 1.5,
32        col = "brown",
33        length = 0.1)
34
35 # Légende
36 legend("topright",
37        legend = c("Série (plot)",
38                  "Repère horizontal (abline)",
39                  "Point manuel (points)",
40                  "Segment (segments)",
41                  "Flèche (arrows)"),
42        col = c("blue", "gray", "darkgreen", "orange", "brown"),
43        pch = c(16, NA, 17, NA, NA),
44        lty = c(1, 2, NA, 1, 1),
45        bty = "n")

```

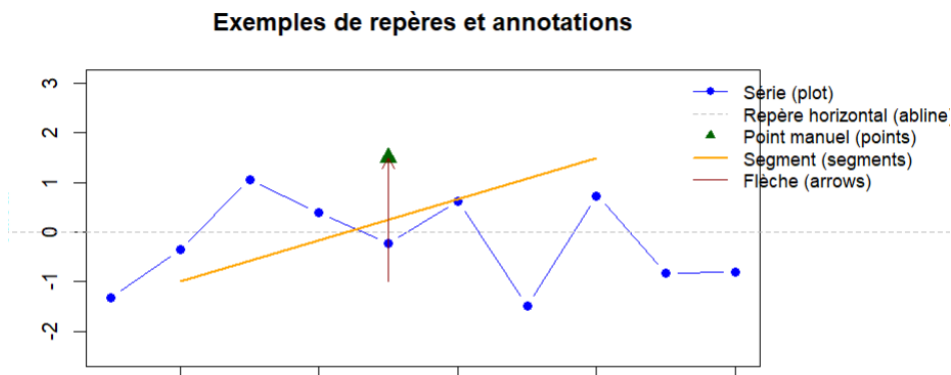


FIGURE 4 – Superposition d'éléments graphiques avec `abline()`, `points()`, `segments()` et `arrows()`

### 13.4 Histogramme avec `hist()`

L'histogramme permet de représenter la distribution d'une variable numérique continue en la divisant en intervalles (ou "classes") et en comptant le nombre d'observations dans chaque intervalle.

#### Syntaxe

```
hist(x,
     breaks = "Sturges",
     col = NULL,
     main = NULL,
     xlab = NULL,
     ylab = NULL,
     freq = TRUE,
     ...)
```

- `x` : vecteur numérique (variable à représenter)
- `breaks` : méthode ou nombre de classes (ex. "Sturges", 10, `seq(0, 100, 10)`)
- `col` : couleur des barres
- `main` : titre du graphique
- `xlab`, `ylab` : titres des axes
- `freq = TRUE` : affiche les effectifs (ou `FALSE` pour les densités)

#### Exemple : Visualisation avec `hist()`

```

1 data(airquality)
2
3 # Histogramme simple
4 hist(airquality$Ozone)
5
6 # Histogramme personnalisé
7 hist(airquality$Ozone,
8     breaks = 10,
9     col = "skyblue",
10    border = "white",
11    main = "Distribution de l'Ozone",
12    xlab = "Ozone (ppb)",
13    ylab = "Fréquence",
14    freq = T,
15    xlim = c(0,150),
16    ylim = c(0,40))

```

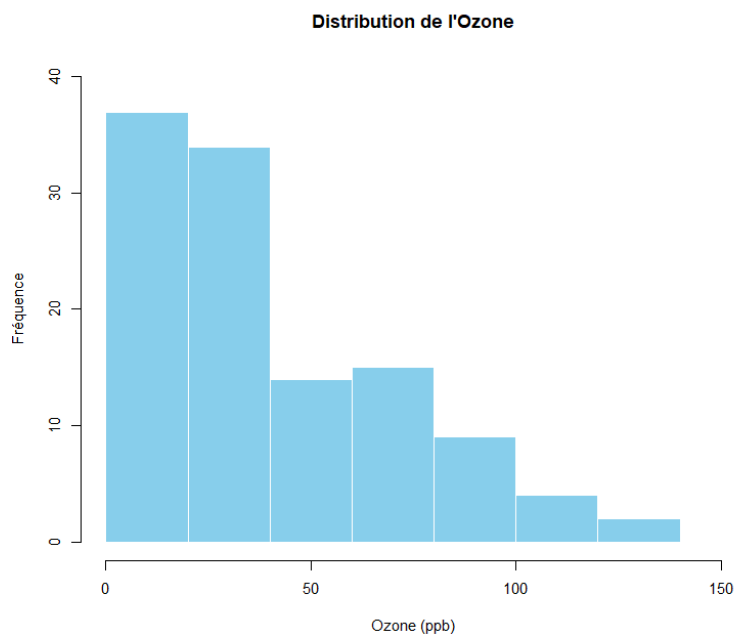


FIGURE 5 – Histogramme de la concentration d’ozone dans l’air

Cet histogramme met en évidence la distribution de la variable `Ozone`, avec un découpage en 10 classes et des barres colorées pour une meilleure lisibilité.



## 13.5 Diagramme en barres avec `barplot()`

Le diagramme en barres permet de représenter des **\*\*valeurs catégorielles\*\*** (ou des fréquences) sous forme de barres verticales ou horizontales. Il est souvent utilisé pour comparer les effectifs ou les moyennes entre plusieurs groupes.

### Syntaxe

```
barplot(height,  
        names.arg = NULL,  
        horiz = FALSE,  
        col = NULL,  
        main = NULL,  
        xlab = NULL,  
        ylab = NULL,  
        ...)
```

- `height` : vecteur numérique ou tableau (valeurs à représenter)
- `names.arg` : noms des catégories (axe X ou Y selon l'orientation)
- `horiz` : si TRUE, barres horizontales
- `col` : couleur des barres
- `main`, `xlab`, `ylab` : titres

### Exemple : Visualisation avec `barplot()`

```
1 genre <- c("Homme", "Femme", "Autre")  
2 effectifs <- c(45, 52, 3)  
3 barplot(effectifs,  
4         names.arg = genre,  
5         col = "steelblue",  
6         main = "Répartition par genre",  
7         ylab = "Effectifs",  
8         xlab = "Genre")
```

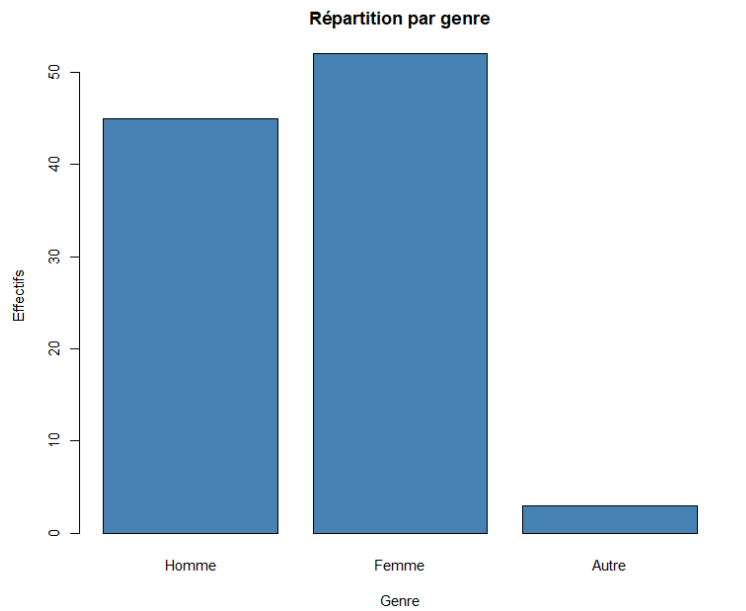


FIGURE 6 – Diagramme en barres représentant la répartition par genre

Ce graphique permet une comparaison rapide des effectifs pour chaque modalité de la variable **genre**. L'utilisation de couleurs rend la visualisation plus intuitive.

### 13.6 Diagramme en boîte avec `boxplot()`

Le **boxplot**, ou diagramme en boîte, est un graphique qui permet de résumer la distribution d'une variable quantitative : médiane, quartiles, minimum, maximum et valeurs aberrantes.

#### Syntaxe

```
boxplot(x,
        horizontal = FALSE,
        col = NULL,
        main = NULL,
        xlab = NULL,
        ylab = NULL,
        ...)
```

- **x** : vecteur numérique ou liste / formule (ex : `note ~ classe`)
- **horizontal** : si `TRUE`, boîte horizontale
- **col** : couleur
- **main**, **xlab**, **ylab** : titres

### Exemple : Visualisation avec boxplot()

```
1
2 note <- c(12, 15, 13, 18, 14, 9, 10, 16, 12, 20, 8, 17)
3 classe <- c(rep("A", 6), rep("B", 6))
4
5 # Boxplot comparatif entre deux classes
6 boxplot(note ~ classe,
7         col = c("tomato", "skyblue"),
8         main = "Distribution des notes par classe",
9         xlab = "Classe",
```

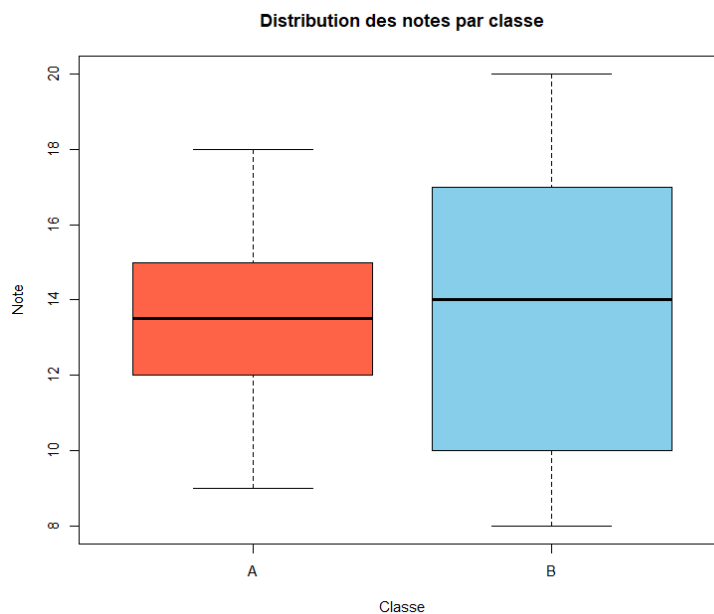


FIGURE 7 – Diagramme en boîte comparant la distribution des notes par classe

Ce graphique montre la répartition des notes dans chaque classe, en mettant en évidence la médiane, l'étendue interquartile, les extrêmes et les éventuels outliers.

### 13.7 Diagramme de dispersion avec stripchart()

La fonction `stripchart()` permet de visualiser la distribution d'une variable numérique à l'aide d'un nuage de points unidimensionnel. C'est une alternative compacte à l'histogramme ou au boxplot.

## Syntaxe

```
stripchart(x, method = "overplot", vertical = FALSE, ...)
```

- **x** : vecteur numérique ou liste/groupe.
- **method** : façon de gérer les superpositions :
  - "overplot" (valeur par défaut) : les points peuvent se superposer.
  - "jitter" : ajout d'un bruit aléatoire pour éviter la superposition.
  - "stack" : empilement vertical ou horizontal.
- **vertical** : TRUE pour un affichage vertical, FALSE pour horizontal.
- **group.names** : noms des groupes (si liste ou facteur).

## Exemple simple avec bruit (jitter)

```
1 # Tracé horizontal avec jitter
2 stripchart(x,
3           method = "jitter",      # évite les superpositions
4           pch = 19,              # point plein
5           col = "darkblue",      # couleur
6           main = "Diagramme de dispersion",
7           xlab = "Valeur")
```

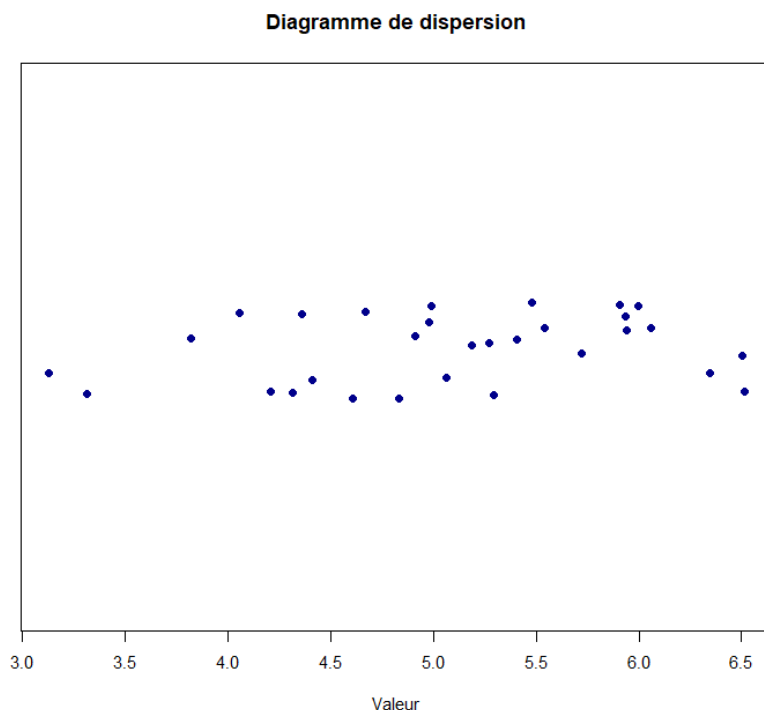


FIGURE 8 – Diagramme de dispersion horizontal avec `stripchart()` et bruit aléatoire

Ce type de graphique est très utile pour observer rapidement la densité locale et les valeurs aberrantes dans un petit échantillon.

## 13.8 Autres graphiques

`dotchart()`

```
1 valeurs <- c(23, 17, 35, 29, 12)
2 noms <- c("Produit A", "Produit B", "Produit C", "Produit D", "Produit E")
3
4 # Dotchart simple
5 dotchart(valeurs,
6           labels = noms,
7           main = "Diagramme en points",
8           xlab = "Quantité vendue",
9           col = "darkblue",
10          pch = 19)
```

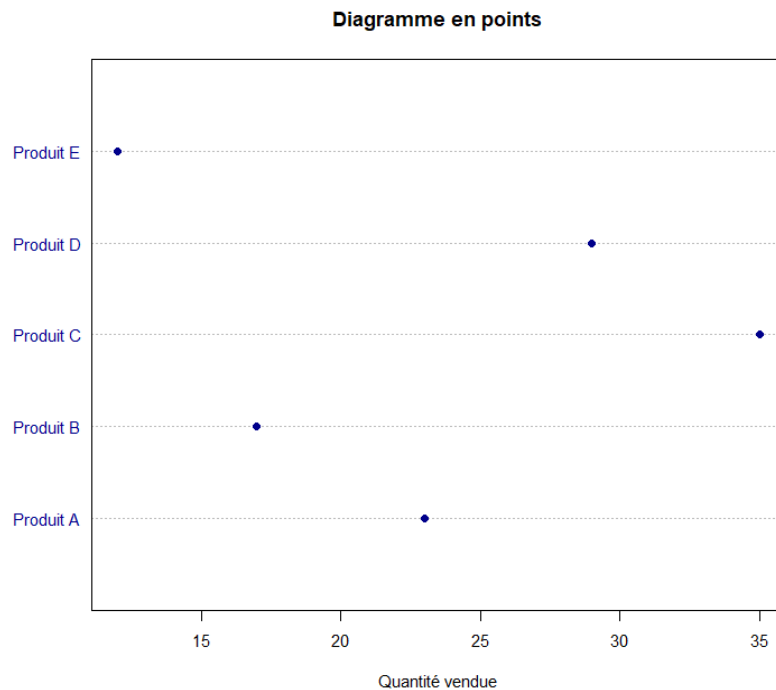


FIGURE 9 – Diagramme en points avec `dotchart()`

Matrice de scatterplots — `pairs()`

```

1 # Exemple sur le jeu de données iris
2 data(iris)
3
4 # Affiche la matrice de scatterplots pour les variables numériques
5 pairs(iris[, 1:4],
6       main = "Matrice de scatterplots - iris",
7       pch = 21,
8       bg = c("red", "green3", "blue")[iris$Species])

```

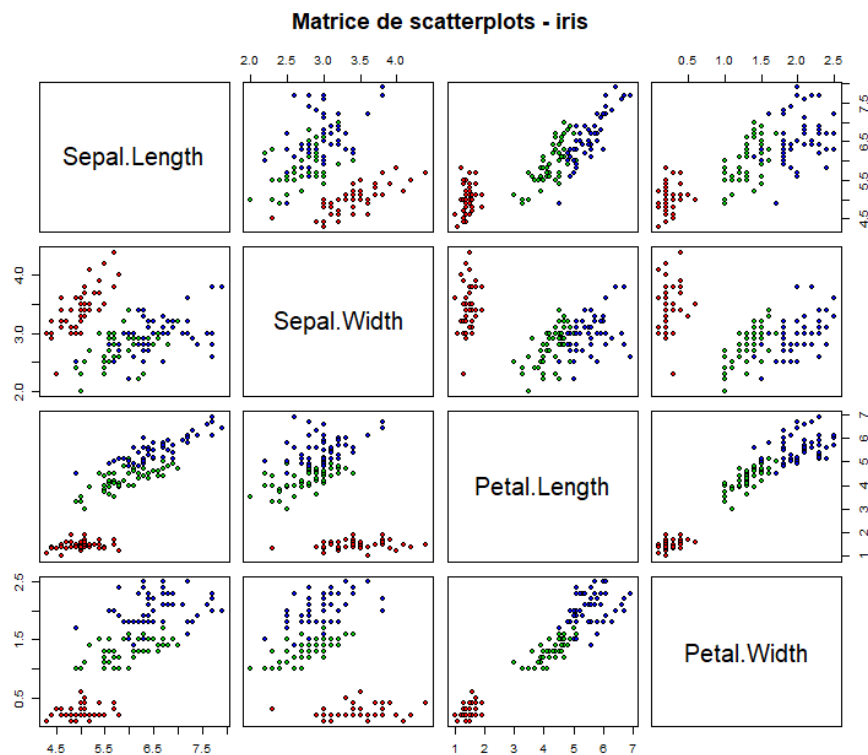


FIGURE 10 – Matrice de scatterplots avec `pairs()` sur le jeu iris

### Tracés multiples — `matplot()`

```

1 # Création de données
2 x <- 1:10
3 y <- cbind(sin(x), cos(x), sin(x / 2))
4
5 # Tracés multiples
6 matplot(x, y,
7         type = "l",           # type ligne
8         lty = 1,             # ligne continue
9         col = c("red", "blue", "darkgreen"),
10        lwd = 2,

```

```

11     main = "Tracés superposés",
12     xlab = "Temps",
13     ylab = "Valeur")
14
15 legend("topright",
16       legend = c("sin(x)", "cos(x)", "sin(x/2)"),
17       col = c("red", "blue", "darkgreen"),
18       lty = 1,
19       bty = "n")

```

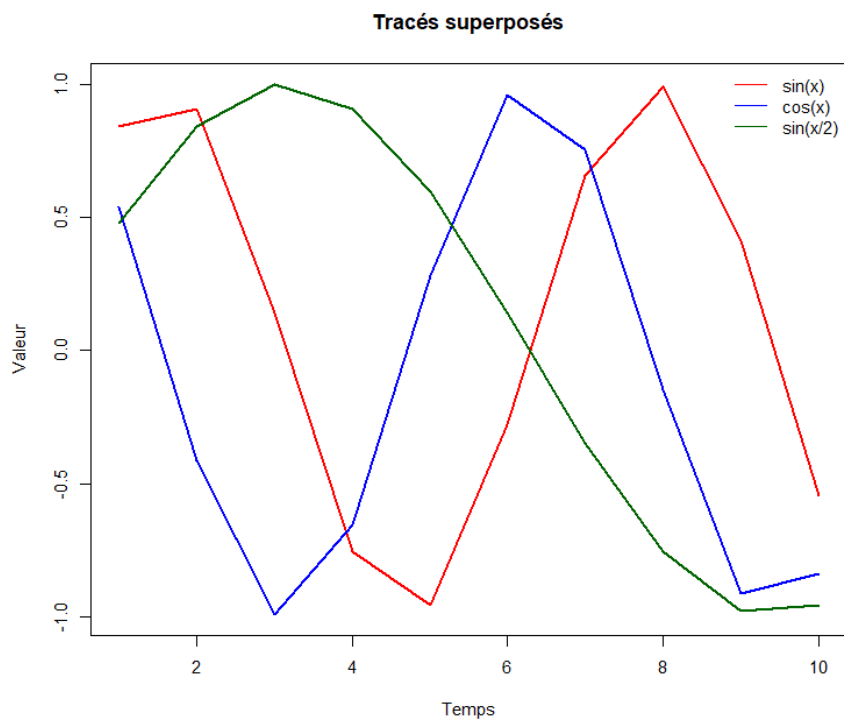


FIGURE 11 – Tracés multiples avec `matplot()`

## 13.9 Visualisation avancée en 2D et 3D

### 1. Affichage d'une image matricielle — `image()`

```

1 # Création d'une matrice de données
2 z <- outer(1:20, 1:20, function(x, y) sin(x / 3) + cos(y / 5))
3
4 # Affichage matriciel
5 image(z,
6       main = "Visualisation matricielle",
7       col = terrain.colors(20))

```

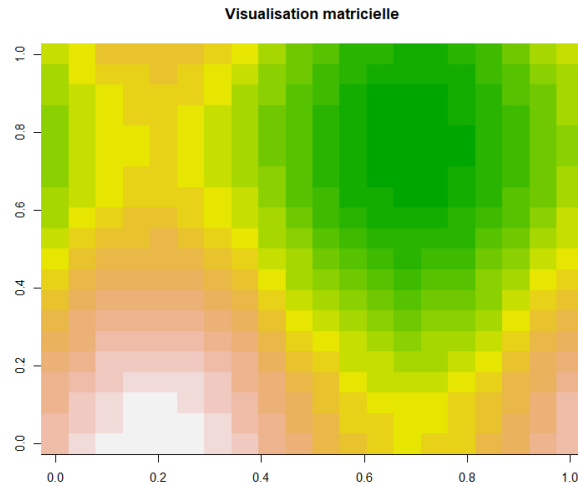


FIGURE 12 – Affichage matriciel avec `image()`

## 2. Courbes de niveau — `contour()`

```

1 # Même matrice que précédemment
2 z <- outer(1:20, 1:20, function(x, y) sin(x / 3) + cos(y / 5))
3
4 # Courbes de niveau
5 contour(z,
6         main = "Courbes de niveau",
7         col = "blue")

```

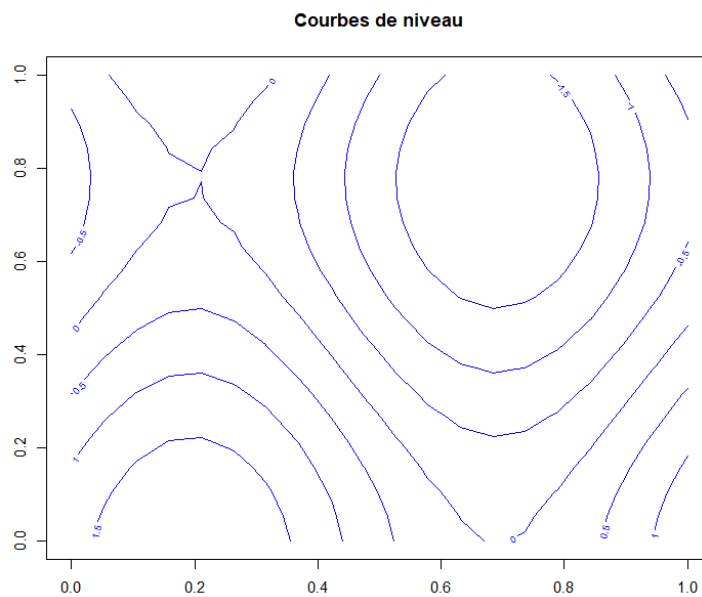


FIGURE 13 – Courbes de niveau avec `contour()`



### 3. Représentation 3D — persp()

```
1 x <- seq(-10, 10, length = 30)
2 y <- seq(-10, 10, length = 30)
3 z <- outer(x, y, function(x, y) cos(sqrt(x^2 + y^2)))
4
5 # Affichage 3D
6 persp(x, y, z,
7       theta = 30, phi = 30,
8       col = "lightblue",
9       border = NA,
10      shade = 0.5,
11      main = "Surface 3D")
```

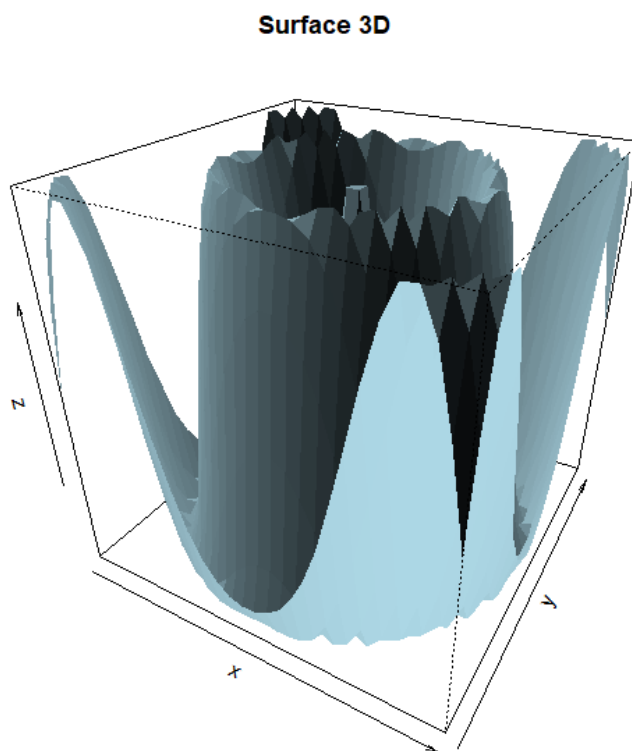


FIGURE 14 – Surface 3D avec persp()

## 14 Quelques lois usuelles :

Pour une loi  $L$  : dL (densité/proba), pL (répartition), qL (quantiles), rL (tirage).

Ex. : dnorm(x,mean,sd), pbinom(k,n,p), rt(n,df).

## 14.1 Normale $\mathcal{N}(\mu, \sigma^2)$

**Usage** : somme de petits effets additifs, erreurs de mesure, phénomènes “centrés”.

- Paramètres : moyenne  $\mu$ , écart-type  $\sigma$ .
- $E[X] = \mu$ ,  $\text{Var}(X) = \sigma^2$ .
- **Quand ?** Variable continue & symétrique ; théorème central limite ; données gaussiennes.
- R : `dnorm(x, mu, sd)`, `pnorm(x, mu, sd)`, `qnorm(p, mu, sd)`, `rnorm(n, mu, sd)`.

## 14.2 Student $t_\nu$ (ddl $\nu$ )

**Usage** : moyenne d'un petit échantillon quand  $\sigma$  inconnu, tests/IC avec  $n$  petit.

- Centre 0, queues plus épaisses que la normale (se rapproche de  $\mathcal{N}(0, 1)$  quand  $\nu \rightarrow \infty$ ).
- **Quand ?** Statistique  $T = \sqrt{n} \frac{\bar{X} - \mu_0}{S}$  ; modèles robustes aux valeurs extrêmes.
- R : `dt(x, df)`, `pt(x, df)`, `qt(p, df)`, `rt(n, df)`.

## 14.3 Khi-deux $\chi_\nu^2$

**Usage** : sommes de carrés normalisés, tests d'ajustement/indépendance.

- $\chi_\nu^2 = \sum_{i=1}^\nu Z_i^2$  avec  $Z_i \sim \mathcal{N}(0, 1)$  iid.
- $E[X] = \nu$ ,  $\text{Var}(X) = 2\nu$ .
- **Quand ?** Variance empirique  $\sim \chi^2$  ; tableaux de contingence (test  $\chi^2$ ).
- R : `dchisq(x, df)`, `pchisq(x, df)`, `qchisq(p, df)`, `rchisq(n, df)`.

## 14.4 Binomiale $\mathcal{B}(n, p)$

**Usage** :  $n$  essais indépendants, succès/échec (Bernoulli).

- $k = 0, \dots, n$ ,  $\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$ .
- $E[X] = np$ ,  $\text{Var}(X) = np(1 - p)$ .
- **Quand ?** Comptage de succès (clics oui/non, défaut oui/non) avec  $n$  fixé.
- R : `dbinom(k, n, p)`, `pbinom(k, n, p)`, `qbinom(p, n, p)`, `rbinom(n, size, p)`.

## 14.5 Poisson $\mathcal{P}(\lambda)$

**Usage** : nombre d'événements rares sur un intervalle (temps/espace), indépendants.

- $k = 0, 1, \dots$ ,  $\Pr(X = k) = e^{-\lambda} \lambda^k / k!$ .
- $E[X] = \lambda$ ,  $\text{Var}(X) = \lambda$ .
- **Quand ?** Comptage sans  $n$  fixé (arrivées par minute, défauts par mètre).

- R : `dpois(k, lambda)`, `ppois(k, lambda)`, `qpois(p, lambda)`, `rpois(n, lambda)`.

## 14.6 Exponentielle $\text{Exp}(\lambda)$

**Usage** : *temps d'attente* jusqu'au prochain événement de Poisson (mémoire nulle).

- $x \geq 0$ ,  $f(x) = \lambda e^{-\lambda x}$ ;  $E[X] = 1/\lambda$ ,  $\text{Var}(X) = 1/\lambda^2$ .
- **Quand ?** Durées entre événements, fiabilité (sans vieillissement).
- R : `dexp(x, rate)`, `pexp(x, rate)`, `qexp(p, rate)`, `rexp(n, rate)`.

## 14.7 Uniforme $\mathcal{U}(a, b)$

**Usage** : ignorance totale dans un intervalle ; générateurs aléatoires de base.

- $E[X] = (a + b)/2$ ,  $\text{Var}(X) = (b - a)^2/12$ .
- R : `dunif(x, min, max)`, `punif`, `qunif`, `runif`.

## Choisir la loi ?

1. **Type de variable** : *discrète* (comptages) ou *continue* (mesures) ?
2. **Comptage** :
  - $n$  **fixé**, succès/échec  $\Rightarrow$  **Binomiale**  $\mathcal{B}(n, p)$ .
  - $n$  **non fixé**, événements rares dans un intervalle  $\Rightarrow$  **Poisson**  $\mathcal{P}(\lambda)$ .
3. **Durée jusqu'à un événement** (processus de Poisson)  $\Rightarrow$  **Exponentielle**.
4. **Somme ou moyenne de mesures bruitées, symétriques**  $\Rightarrow$  **Normale**.
5. **Inférence sur la moyenne avec  $\sigma$  inconnu** :
  - $n$  petit  $\Rightarrow$  **Student**  $t_{n-1}$  ;
  - $n$  grand  $\Rightarrow$  **Normale** (approx.).
6. **Variance à partir de données normales**  $\Rightarrow \chi^2$ .

- Densité en  $x_0$  : `dnorm(x0, mean=mu, sd=sigma)`
- Probabilité  $P(X \leq x_0)$  : `pnorm(x0, mean=mu, sd=sigma)`
- Queue  $P(X > x_0)$  : `pnorm(x0, mean=mu, sd=sigma, lower.tail=FALSE)`
- Quantile  $q_p$  tel que  $P(X \leq q_p) = p$  : `qnorm(p, mean=mu, sd=sigma)`
- Simulation de  $n$  valeurs : `rnorm(n, mean=mu, sd=sigma)`

## Erreurs fréquentes (à éviter)

- **Confondre variance et écart-type** dans R : `dnorm(..., sd = écart-type)`.

- Oublier **la correction de continuité** pour comparer une loi continue (normale) à une loi discrète (binomiale).
- Utiliser la **normale** quand  $np$  ou  $n(1-p)$  sont trop petits ( $< 5$ ) : préférer **Poisson** ou exact binomial.
- Croire qu'un **mélange de normales** est normal : *faux* (souvent bimodal).

## 15 Loi normale $\mathcal{N}(\mu, \sigma^2)$

### 15.1 Définition et paramètres

On dit que  $X \sim \mathcal{N}(\mu, \sigma^2)$  si sa densité est

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}.$$

Paramètres :  $\mu$  (moyenne),  $\sigma > 0$  (écart-type), variance =  $\sigma^2$ .

**Faits clés.**

- $Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1)$  (standardisation).
- Environ 68% des valeurs dans  $[\mu - \sigma, \mu + \sigma]$ , 95% dans  $[\mu - 2\sigma, \mu + 2\sigma]$ , 99,7% dans  $[\mu - 3\sigma, \mu + 3\sigma]$ .
- **R et la normale** : `dnorm` (densité), `pnorm` (répartition), `qnorm` (quantiles), `rnorm` (simulation).

### 15.2 Quand l'utiliser ?

- Phénomènes continus, symétriques, bruit de mesure, somme de petits effets (TCL).
- Approximation de nombreuses statistiques (moyenne d'échantillon, grandes tailles).

### 15.3 Exemples :

**Exemple 1 — tracer une densité.**  $X \sim \mathcal{N}(19, 3)$  (ici 3 est la *variance*, donc  $sd = \sqrt{3}$ ).

```
1 mu <- 19; sd <- sqrt(3)
2 x <- seq(mu-4*sd, mu+4*sd, length.out = 400)
3 plot(x, dnorm(x, mean=mu, sd=sd), type="l", lwd=2,
4      main="Densité N(19, 3)", xlab="x", ylab="densité")
```

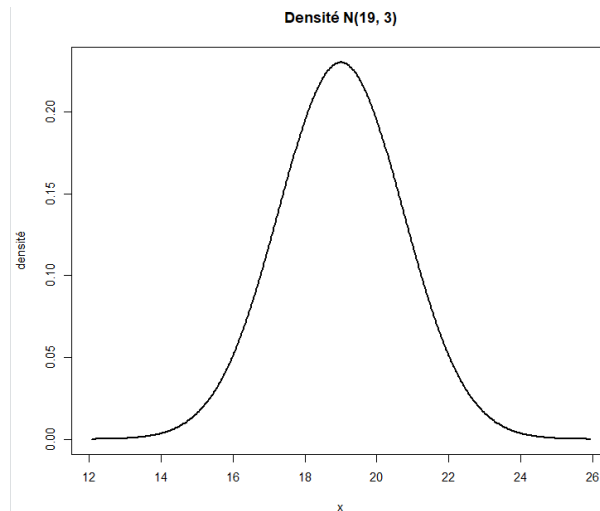


FIGURE 15 – Enter Caption

**Exemple 2 — probas avec standardisation.**  $X \sim \mathcal{N}(15, 9)$ , calculer  $P(16 \leq X \leq 20)$ .

```
1 mu <- 15; sd <- 3
2 pnorm(20, mu, sd) - pnorm(16, mu, sd)      # direct
3 # ou en normalisant :
4 pnorm((20-mu)/sd) - pnorm((16-mu)/sd)      # Z ~ N(0,1)
```

**Exemple 3 — quantiles (valeurs seuils).** Trouver  $q_{0.025}$  et  $q_{0.975}$  de  $N(0, 1)$ .

```
1 qnorm(0.025, 0, 1)      # ~ -1.96
2 qnorm(0.975, 0, 1)      # ~ +1.96
```

Pour  $Y \sim \mathcal{N}(19, 3)$  :  $q_{0.975}(Y) = 19 + \sqrt{3} q_{0.975}(Z)$ .

```
1 mu <- 19; sd <- sqrt(3)
2 mu + sd * qnorm(0.975)      # vérifie la relation de changement d'échelle
```

**Exemple 4 — simulation et histogramme avec densité.**

```
1 set.seed(123)
2 mu <- 15; sd <- sqrt(3)
3 x <- rnorm(1000, mean=mu, sd=sd)
4 hist(x, freq=FALSE, col="gray90", border="white",
5       main="Histogramme + densité théorique", xlab="x")
6 curve(dnorm(x, mean=mu, sd=sd), add=TRUE, lwd=2, col="red")
```

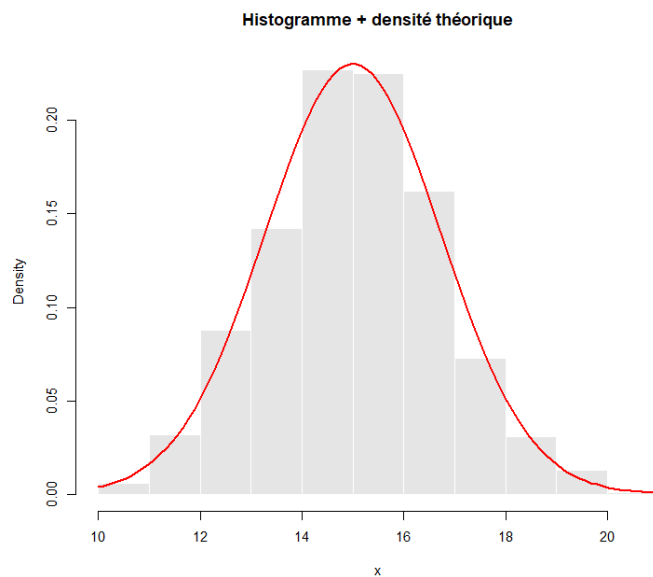


FIGURE 16 – Histogramme simulé (1000 tirages) et densité théorique de la loi  $\mathcal{N}(15, 3)$ .

**Exemple 5 — tracer la fonction de répartition  $F(x)$ .**

```
1 mu <- 0; sd <- 1
2 curve(pnorm(x, mu, sd), from=-4, to=4, lwd=2,
3       main="Fonction de répartition N(0,1)", xlab="x", ylab="F(x)")
4 abline(h=c(0,1), v=0, lty=3)
```

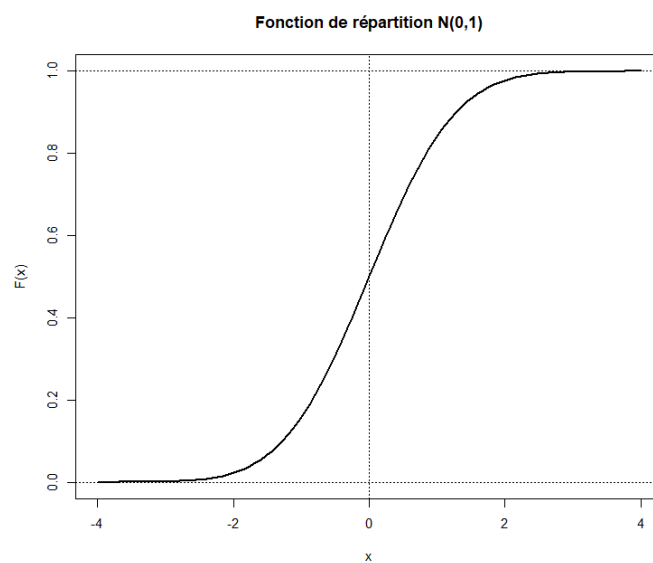


FIGURE 17 – Fonction de répartition de la loi normale standard  $N(0, 1)$ .

# Exemple d'analyse avec le jeu de données airquality

Dans cette section, nous allons mettre en application l'ensemble des notions vues dans ce document (structures de données, manipulation, statistiques, visualisation, modélisation, etc.) en utilisant un jeu de données réel déjà inclus dans R : `airquality`.

## 1. Accéder aux jeux de données intégrés dans R

R met à disposition de nombreux jeux de données intégrés, très utiles pour l'apprentissage et l'expérimentation.

- Pour voir la liste de tous les jeux de données disponibles dans votre session (pratique si vous cherchez un dataset pour vos projets personnels) :

```
1 data()
```

```
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
CO2                Carbon Dioxide Uptake in Grass Plants
ChickWeight        Weight versus age of chicks on different diets
DNase              Elisa assay of DNase
EuStockMarkets     Daily Closing Prices of Major European Stock Indices, 1991-1998
Formaldehyde       Determination of Formaldehyde
HairEyeColor       Hair and Eye Color of Statistics Students
Harman23.cor       Harman Example 2.3
Harman74.cor       Harman Example 7.4
Indomethacin       Pharmacokinetics of Indomethacin
InsectSprays       Effectiveness of Insect Sprays
JohnsonJohnson    Quarterly Earnings per Johnson & Johnson Share
LakeHuron          Level of Lake Huron 1875-1972
LifeCycleSavings   Intercountry Life-Cycle Savings Data
Loblolly           Growth of Loblolly Pine Trees
Nile               Flow of the River Nile
Orange             Growth of Orange Trees
OrchardSprays      Potency of Orchard Sprays
PlantGrowth        Results from an Experiment on Plant Growth
Puromycin          Reaction Velocity of an Enzymatic Reaction
Seatbelts          Road Casualties in Great Britain 1969-84
Theoph             Pharmacokinetics of Theophylline
Titanic            Survival of passengers on the Titanic
ToothGrowth        The Effect of Vitamin C on Tooth Growth in Guinea Pigs
UCBAdmissions      Student Admissions at UC Berkeley
UKDriverDeaths     Road Casualties in Great Britain 1969-84
UKgas              UK Quarterly Gas Consumption
USAccDeaths        Accidental Deaths in the US 1973-1978
USArrests          Violent Crime Rates by US State
USJudgeRatings     Lawyers' Ratings of State Judges in the US Superior Court
USPersonalExpenditure Personal Expenditure Data
UScitiesD          Distances Between European Cities and Between US Cities
VADeaths           Death Rates in Virginia (1940)
WWWusage           Internet Usage per Minute
WorldPhones        The World's Telephones
ability.cov        Ability and Intelligence Tests
airmiles           Passenger Miles on Commercial US Airlines, 1937-1960
airquality         New York Air Quality Measurements
anscombe           Anscombe's Quartet of 'Identical' Simple Linear Regressions
attenu             The Joyner-Boore Attenuation Data
attitude           The Chatterjee-Price Attitude Data
austres            Quarterly Time Series of the Number of Australian Residents
```

FIGURE 18 – Extrait du résultat de la commande `data()`.

Chaque jeu de données intégré possède une page d'aide détaillée que l'on peut consulter. Exemple avec `airquality` :

1 ?airquality

R: New York Air Quality Measurements Find in Topic

**Format**

A data frame with 153 observations on 6 variables.

[,1]	Ozone	numeric Ozone (ppb)
[,2]	Solar.R	numeric Solar R (lang)
[,3]	Wind	numeric Wind (mph)
[,4]	Temp	numeric Temperature (degrees F)
[,5]	Month	numeric Month (1--12)
[,6]	Day	numeric Day of month (1--31)

**Details**

Daily readings of the following air quality values for May 1, 1973 (a Tuesday) to September 30, 1973.

- Ozone: Mean ozone in parts per billion from 1300 to 1500 hours at Roosevelt Island
- Solar.R: Solar radiation in Langleys in the frequency band 4000–7700 Angstroms from 0800 to 1200 hours at Central Park
- Wind: Average wind speed in miles per hour at 0700 and 1000 hours at LaGuardia Airport
- Temp: Maximum daily temperature in degrees Fahrenheit at LaGuardia Airport.

**Source**

The data were obtained from the New York State Department of Conservation (ozone data) and the National Weather Service (meteorological data).

**References**

Chambers, J. M., Cleveland, W. S., Kleiner, B. and Tukey, P. A. (1983) *Graphical Methods for Data Analysis*. Belmont, CA: Wadsworth.

**Examples**

[Run examples](#)

```
require(graphics)
pairs(airquality, panel = panel.smooth, main = "airquality data")
```

FIGURE 19 – Documentation du jeu de données `airquality`.

## Téléchargement de jeux de données

Pour explorer d'autres jeux de données libres à télécharger :

- Rdatasets : collection de jeux de données R
- CRAN : site officiel de R
- data.gouv.fr : données publiques françaises

## 2. Chargement et description du jeu de données

Charger `airquality`, vérifier sa structure et comprendre les variables.

```
1 # 1) Charger le dataset
2 data("airquality")
3
4 # 2) Copie
5 data <- airquality
6
7 # 3) Dimensions, structure
```



```

8 dim(data)           # nb lignes x nb colonnes
9 nrow(data)          # nb d'observations (lignes)
10 ncol(data)          # nb de variables (colonnes)
11 str(data)           # types des variables
12 head(data, 10)      # les 10 premières lignes
13 tail(data, 10)      # les 10 dernières lignes
14 summary(data)       # stats descriptives (montre aussi les NA)
15 names(data)         # noms des colonnes

```

```

1 colnames(data)       # nom des colonnes = names(data)
2 rownames(data)       # noms des lignes
3
4 # 5) Types de données
5 sapply(data, class)   # classe de chaque variable
6 sapply(data, typeof) # type de chaque variable
7
8 # 6) Valeurs manquantes
9 colSums(is.na(data))  # nb de NA par variable
10 anyNA(data)          # TRUE si des NA sont présents
11
12 # 7) Doublons
13 anyDuplicated(data)   # nb de lignes dupliquées
14 sum(duplicated(data)) # comptage exact
15
16 # 9) Statistiques élémentaires par variable numérique
17 sapply(data, min, na.rm=TRUE) # minimum
18 sapply(data, max, na.rm=TRUE) # maximum
19 sapply(data, mean, na.rm=TRUE) # moyenne
20 sapply(data, median, na.rm=TRUE) # médiane
21 sapply(data, sd, na.rm=TRUE) # écart-type
22 sapply(data, var, na.rm=TRUE) # variance
23 sapply(data, IQR, na.rm=TRUE) # intervalle interquartile

```

### 3 Exemple de visualisations sur le dataset

Explorer visuellement la distribution des variables et leurs relations afin de mieux comprendre la structure globale du jeu de données.

**3.1 Histogrammes.** Les histogrammes permettent d'observer la distribution d'une variable numérique.

```

1 # Histogrammes (deux graphiques côte à côte)
2 par(mfrow = c(1,2))
3 hist(data$lOzone, main="Distribution de lOzone",

```

```

4     xlab="Ozone (ppb)", col="lightblue")
5 hist(data$Temp, main="Distribution de la Température",
6       xlab="Température (F)", col="lightcoral")
7 par(mfrow = c(1,1)) # revenir au mode 1 graphique

```

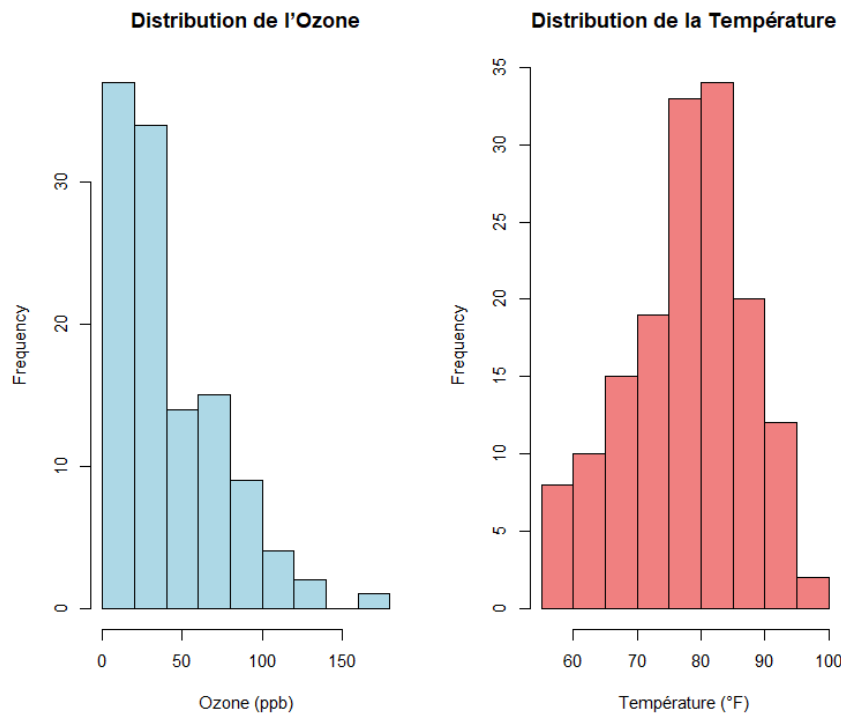


FIGURE 20 – Histogrammes de la distribution de l’ozone et de la température.

*Remarque : La distribution de l’ozone est très asymétrique avec beaucoup de valeurs faibles et quelques pics élevés. La température, en revanche, est concentrée entre 70 et 90 °F, ce qui traduit une répartition plus régulière et saisonnière.*

**3.2 Boxplots.** Les boxplots mettent en évidence la médiane, la dispersion et les valeurs extrêmes.

```

1 par(mfrow = c(1,2))
2 boxplot(data$Ozone, main="Boîte à moustaches Ozone", col="lightgreen")
3 boxplot(data$Temp, main="Boîte à moustaches Température", col="orange")
4 par(mfrow = c(1,1))

```

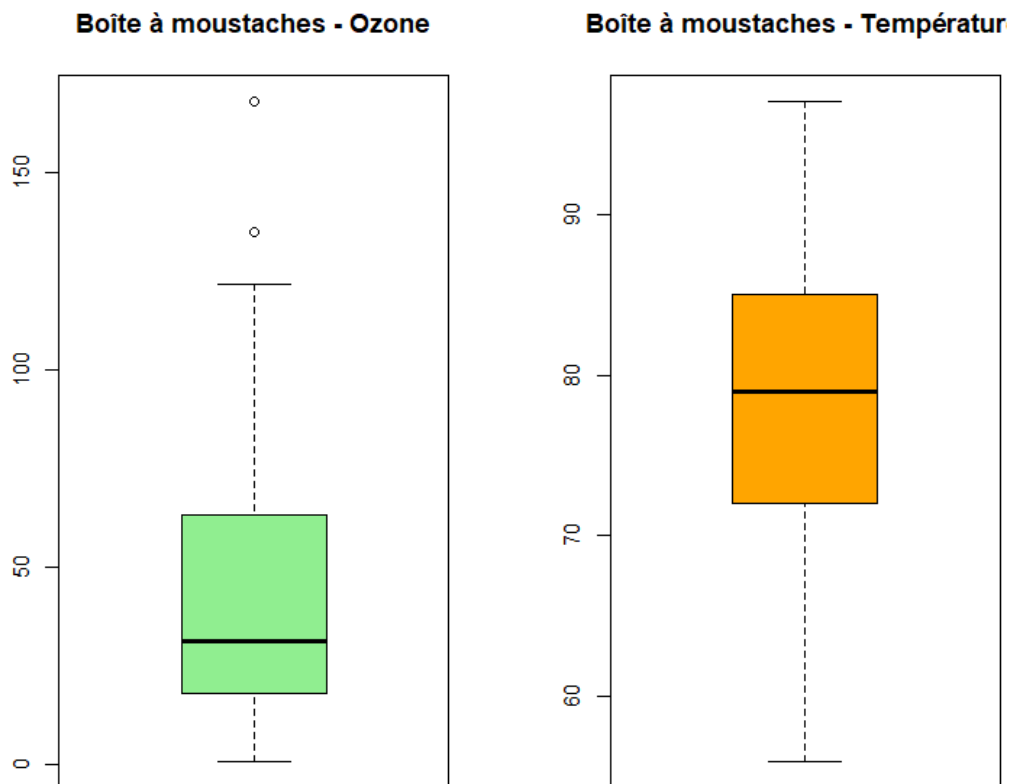


FIGURE 21 – Boxplots globaux de l’ozone et de la température.

*L’ozone présente une forte variabilité avec des valeurs extrêmes au-dessus de 100 ppb. La température est plus homogène, centrée autour de 80 °F, avec une dispersion modérée et peu de valeurs aberrantes.*

**3.3 Boxplots par mois.** Utile pour comparer la répartition de l’ozone selon les mois (mai à septembre).

```

1 # Boxplot Ozone par mois
2 boxplot(Ozone ~ Month, data = data,
3         main = "Ozone par mois",
4         xlab = "Mois", ylab = "Ozone (ppb)",
5         col = "lightblue")

```

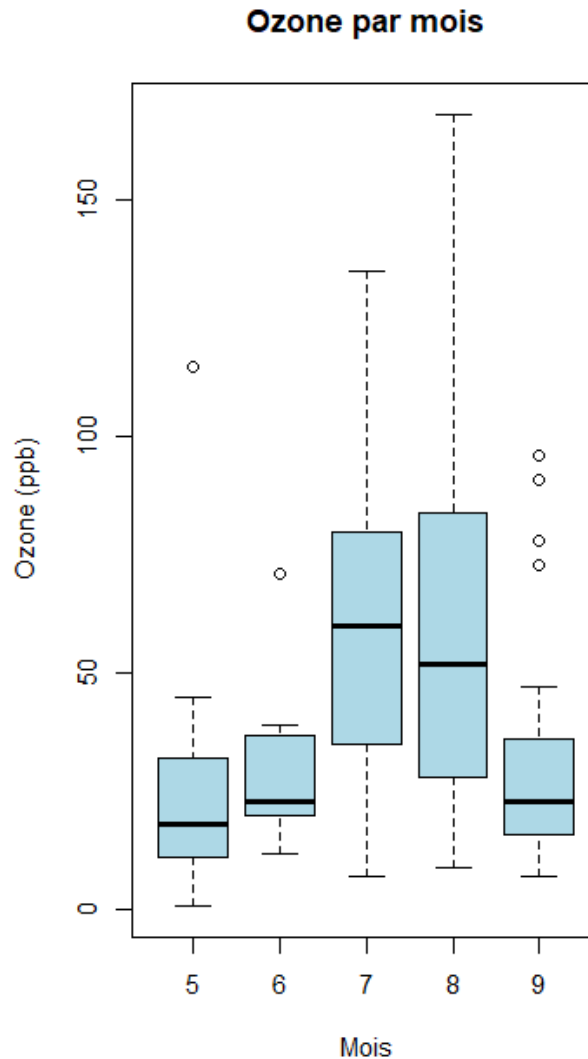


FIGURE 22 – Distribution de l’ozone par mois (mai–septembre).

*Les niveaux d’ozone sont faibles au printemps (mai–juin) puis augmentent fortement en juillet et août, avant de redescendre en septembre. La variabilité est particulièrement marquée en été, avec de nombreux outliers.*

**3.4 Nuage de points.** Permet d’examiner la relation entre la température et l’ozone.

```

1 # Scatterplot Ozone ~ Température
2 plot(data$Temp, data$Ozone,
3       pch=19, col="darkgreen",
4       main="Lien entre Ozone et Température",
5       xlab="Température (F)", ylab="Ozone (ppb)")

```

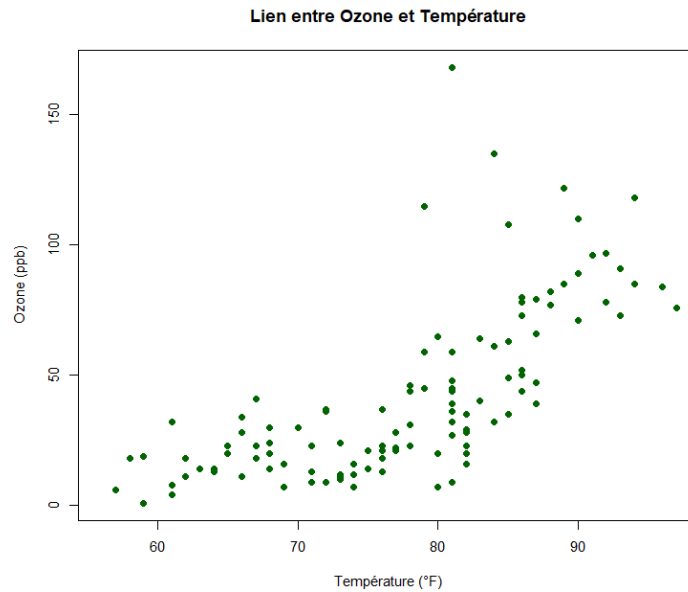


FIGURE 23 – Nuage de points entre la température et l’ozone.

*On observe une corrélation positive entre la température et l’ozone : plus la température augmente, plus les niveaux d’ozone tendent à s’élever. Cependant, la dispersion reste importante pour les fortes chaleurs.*

**3.5 Matrice de corrélation visuelle.** Visualise les relations entre les variables numériques principales.

```
1 # Matrice de nuages de points
2 pairs(data[, c("Ozone", "Solar.R", "Wind", "Temp")],
3       main = "Matrice de nuages de points Variables principales")
```

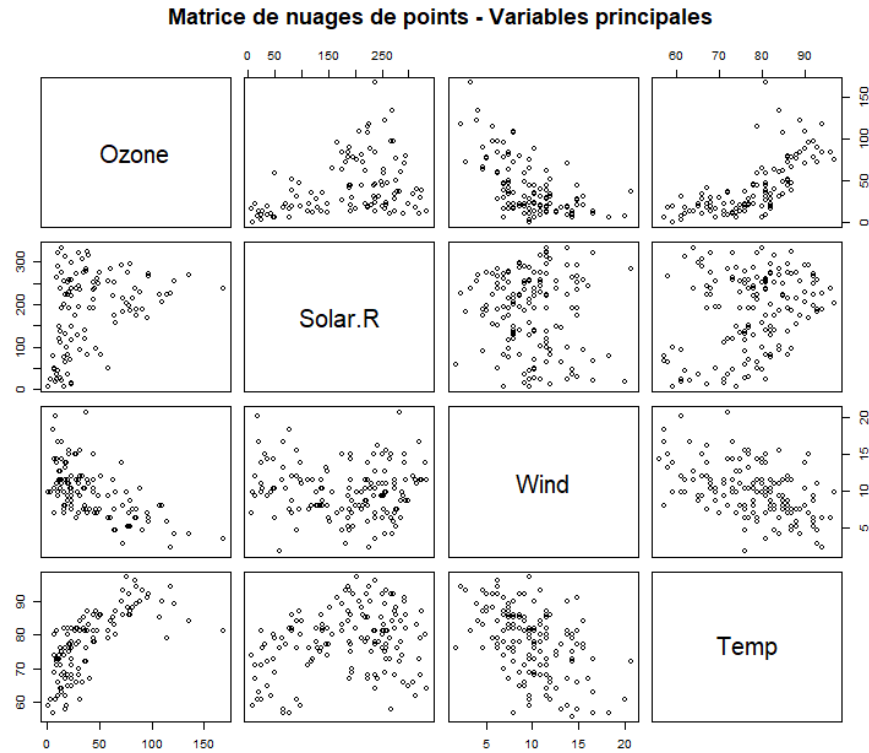


FIGURE 24 – Matrice de nuages de points pour Ozone, Solar.R, Wind et Temp.

*La matrice confirme une corrélation positive entre l’ozone et la température, tandis que le vent est plutôt associé négativement avec l’ozone. La variable Solar.R montre une relation plus diffuse avec les autres indicateurs.*

## 4. Exemple de nettoyage applicable au jeu de données

### 4.1 Valeurs manquantes

```
1 # Créer un dataset avec uniquement les lignes complètes
2 data_cc <- na.omit(airquality)
3
4 # Vérifier les dimensions
5 dim(data_cc)
```

### 4.2 Imputation simple des NA

```
1 # Imputation par la médiane de chaque variable
```

```

2 data_imp <- airquality
3 for (col in c("Ozone", "Solar.R")) {
4   data_imp[[col]][is.na(data_imp[[col]])] <- median(data_imp[[col]], na.rm = TRUE)
5 }
6
7 # Vérification des NA restants
8 colSums(is.na(data_imp))

```

### 4.3 Recodage des mois pour plus de lisibilité

```

1 # Recoder la variable Mois en facteur avec labels
2 data_imp$Month <- factor(data_imp$Month,
3                           levels = 5:9,
4                           labels = c("Mai", "Juin", "Juillet", "Août", "Septembre"))
5
6 # Vérification du nouveau codage
7 table(data_imp$Month)

```

### 4.4 Renommage des colonnes

```

1 names(data_imp) <- c("Ozone_ppb", "Rayonnement", "Vent",
2                     "Temperature_F", "Mois", "Jour")
3 names(data_imp)

```

### 4.5 Conversion d'unités

```

1 # Ajouter une colonne en C
2 data_imp$Temperature_C <- (data_imp$Temperature_F - 32) * 5/9
3 head(data_imp[, c("Temperature_F", "Temperature_C")], 5)

```

### 4.6 Détection et traitement des valeurs aberrantes

```

1 # Identifier les valeurs extrêmes d'ozone (>150)
2 which(data_imp$Ozone_ppb > 150)

```

```

1 # Supprimer les outliers d'ozone
2 data_no_outliers <- subset(data_imp, Ozone_ppb <= 150)
3 dim(data_no_outliers)

```

## 4.7 Normalisation d'une variable

```
1 # Normalisation Min-Max de lozone (entre 0 et 1)
2 ozone_norm <- (data_imp$Ozone_ppb - min(data_imp$Ozone_ppb)) /
3               (max(data_imp$Ozone_ppb) - min(data_imp$Ozone_ppb))
4 summary(ozone_norm)
```

```
1 # Standardisation (moyenne 0, écart-type 1)
2 ozone_z <- scale(data_imp$Ozone_ppb)
3 summary(as.numeric(ozone_z))
```

## 5 Visualisation sur une variable

```
1 # Comparer la distribution d'Ozone avant/après filtrage des outliers
2 op <- par(mfrow = c(1,2)) # sauvegarde de l'état graphique
3 boxplot(airquality$Ozone, main="Ozone brut", col="tomato")
4 boxplot(data_no_outliers$Ozone_ppb, main="Ozone nettoyé", col="lightgreen")
5 par(op) # on rétablit l'état graphique initial
```

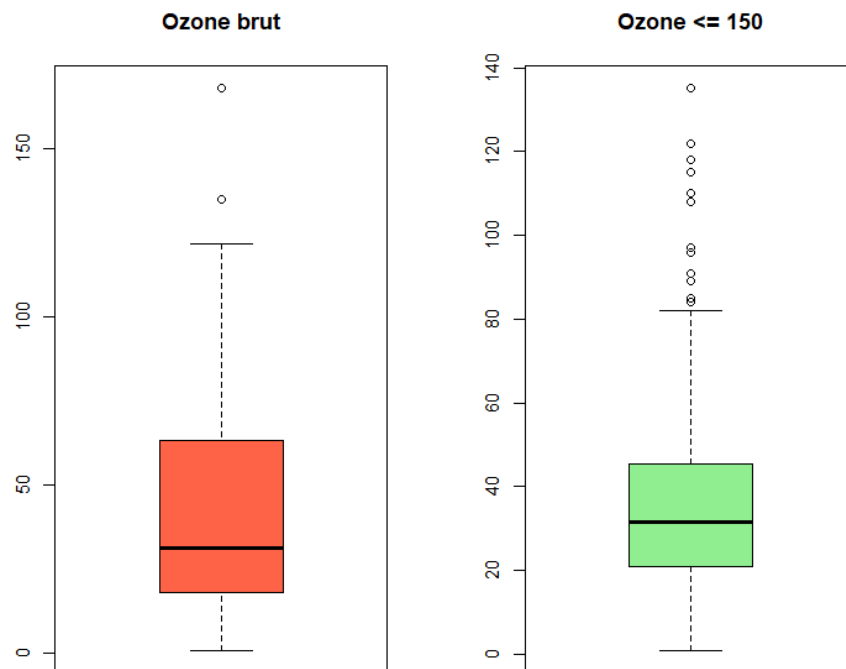


FIGURE 25 – Comparaison de la distribution de l'ozone avant et après filtrage des valeurs extrêmes.



## 6 Visualisation multiple

Le code suivant crée une figure en grille 2x3 avec six graphiques :

- Histogramme + densité pour l'Ozone, la Température, le Vent et le Rayonnement.
- Une comparaison de densités (Ozone vs Température).
- Un boxplot de l'Ozone par mois.

```
1 # Fonction utilitaire : histogramme + densité
2 histo_dens <- function(x, main, xlab, col_hist, col_dens) {
3   x <- x[!is.na(x)]
4   if (length(x) < 2) { plot.new(); title(main); return(invisible()) }
5   d <- density(x)
6   h0 <- hist(x, plot = FALSE, breaks = "FD")
7   ylim_max <- max(max(h0$density, na.rm = TRUE), max(d$y, na.rm = TRUE)) * 1.08
8   xlim_use <- range(c(range(h0$breaks, na.rm=TRUE), range(d$x, na.rm=TRUE)))
9   hist(x, probability = TRUE, col = col_hist, border = "white",
10        main = main, xlab = xlab, ylim = c(0, ylim_max), xlim = xlim_use)
11   lines(d, col = col_dens, lwd = 2)
12   rug(x)
13 }
14
15 op <- par(mfrow = c(2,3), mar = c(4,4,2,1), oma = c(0,0,1,0), xaxs = "i", yaxs = "i")
16
17 # (1) Ozone
18 histo_dens(airquality$Ozone, main = "Ozone", xlab = "ppb",
19            col_hist = "lightblue", col_dens = "red")
20
21 # (2) Température
22 histo_dens(airquality$Temp, main = "Température", xlab = "F",
23            col_hist = "lightgreen", col_dens = "blue")
24
25 # (3) Vent
26 histo_dens(airquality$Wind, main = "Vent", xlab = "mph",
27            col_hist = "lightyellow", col_dens = "darkgreen")
28
29 # (4) Rayonnement
30 histo_dens(airquality$Solar.R, main = "Rayonnement", xlab = "lang",
31            col_hist = "pink", col_dens = "purple")
32
33 # (5) Comparaison densités Ozone/Température
34 oz <- na.omit(airquality$Ozone)
35 tm <- na.omit(airquality$Temp)
36 d1 <- density(oz); d2 <- density(tm)
37 xlim_use <- range(c(d1$x, d2$x))
38 ylim_use <- c(0, max(d1$y, d2$y) * 1.08)
```

```

39 plot(d1, col = "red", lwd = 2, main = "Comparaison des densités",
40      xlab = "Valeur", xlim = xlim_use, ylim = ylim_use)
41 lines(d2, col = "blue", lwd = 2)
42 legend("topright", legend = c("Ozone", "Température"),
43      col = c("red", "blue"), lwd = 2, bty = "n")
44
45 # (6) Boxplot Ozone par mois
46 boxplot(Ozone ~ Month, data = airquality, col = "lightcyan",
47      main = "Ozone par mois", xlab = "Mois (5=Mai ... 9=Sept.)", ylab = "ppb")
48 par(op)

```

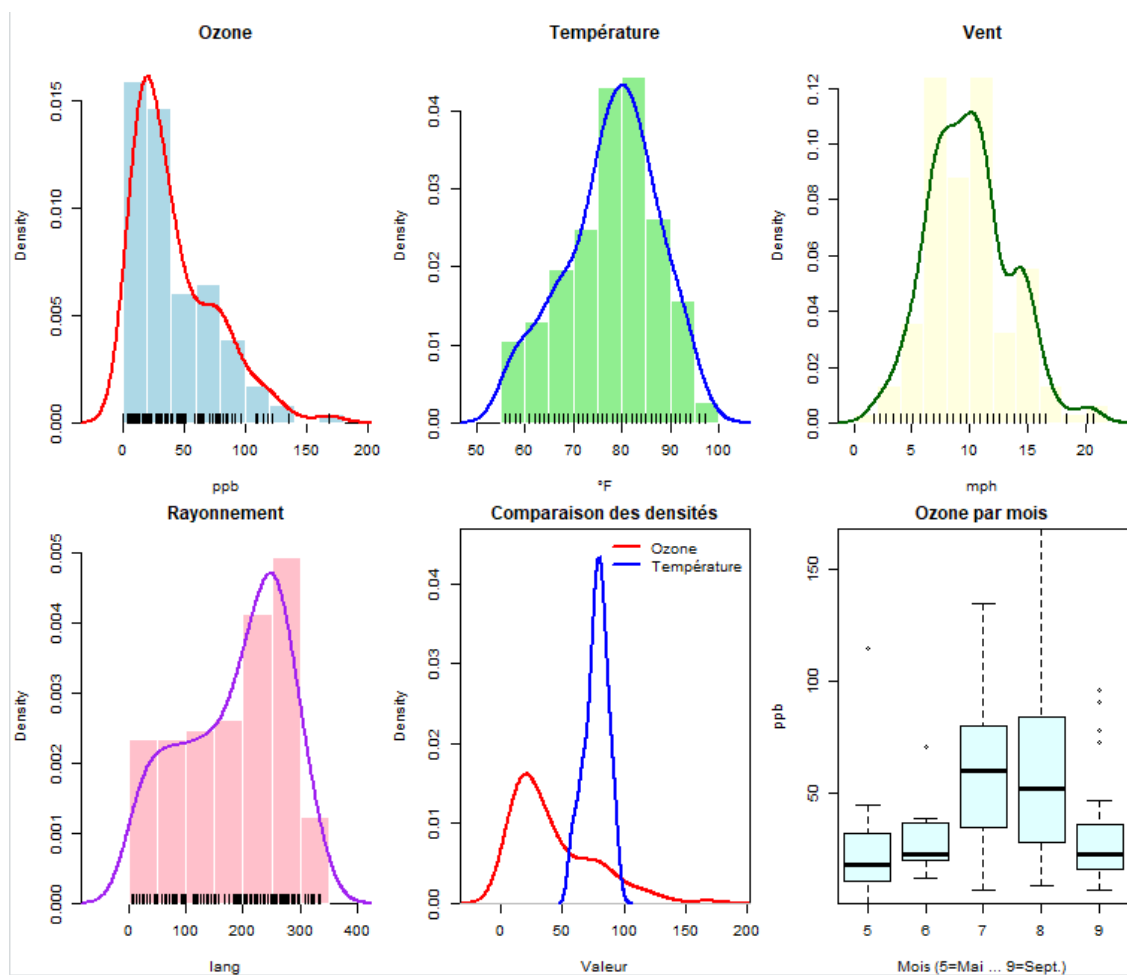


FIGURE 26 – Visualisation en grille (2x3) : distributions, densités et comparaison par mois.

*Interprétation :* Ces graphiques résument les distributions des variables principales du jeu de données `airquality`. Les histogrammes montrent que l’ozone est fortement concentré à de faibles valeurs, tandis que la température suit une distribution plus régulière centrée autour de 80°F. Le vent est globalement distribué entre 5 et 15 mph, et le rayonnement solaire présente une forte variabilité. La comparaison des densités met en évidence des profils très

différents pour l’ozone et la température. Enfin, le boxplot par mois révèle que l’ozone atteint ses valeurs les plus élevées en juillet et août, avec une forte variabilité et plusieurs valeurs extrêmes.

Après cette partie consacrée à l’exploration et à la visualisation des données, les sections suivantes aborderont des thématiques avancées : les *tests statistiques* pour valider des hypothèses, la *modélisation* afin de prédire et expliquer les phénomènes observés, et enfin l’outil **R Markdown**, qui permet de combiner code, résultats.