

CHAPITRE 2

LES OUTILS DE REPRESENTATION D'UN ECHANTILLON

Lorsque l'on fait des statistiques, on travaille en général sur des **observations** sur lesquelles sont mesurées ou observées des **variables** en plus ou moins grand nombre. Une observation est l'unité élémentaire ; nous verrons qu'elle peut avoir des formes multiples : une personne ou un groupe de personnes, un animal ou une portée d'animaux, une plante, un organe ou une seule cellule. Une variable est une caractéristique ou une propriété qu'il est possible de mesurer. Mesurer quelque chose signifie que nous faisons un modèle numérique de la chose à mesurer : nous assignons un nombre à un niveau de la caractéristique à mesurer. Le poids d'une personne est une variable ; nous assignons à deux personnes de même poids la même valeur numérique.

Dans ce chapitre nous allons définir ce qu'est la structure d'un ensemble **observations / variables** que nous appellerons tableau de données (§1) ; nous verrons ensuite comment résumer un tableau de données à l'aide d'un petit nombre de caractéristiques soit unidimensionnelles (§2) soit multidimensionnelles (§3) ; nous présenterons un outil commode pour représenter le tableau de données (§4), puis nous dresserons un bilan des aspects importants de ce chapitre (§5).

1. STRUCTURE D'UN TABLEAU DE DONNEES

1.1 Questions préables

Toute expérience ou toute enquête se fonde sur la mesure de quantités, qui sont les réalisations d'une ou de plusieurs variables. Mais avant toute analyse de ces données, il faut se poser un certain nombre de questions préalables :

- **Que sont les données ?** Sont-elles en nombre important, quelle est leur précision ?
- **Quelles sont les unités ?** Kilomètre ou mile marin, degré Celsius ou degré Fahrenheit, etc. ?
- **D'où proviennent-elles ?** Les valeurs sont-elles raisonnables ? Des valeurs comme 0.7 ou 44 pour le poids d'une dinde en kilogrammes indiquent-elles une simple erreur de transcription ou mettent-elles en doute l'ensemble des relevés ? Pour un relevé pluviométrique, l'emplacement du pluviomètre peut-il entraîner des variations importantes des mesures obtenues ?
- **Comment et quand ont-elles été mesurées ?** Des biâis sont-ils possibles ? L'observateur sait-il arrondir les nombres ? Que signifie la température d'un jour donné, est-ce la valeur maximale, la valeur à une heure précise ? Les enregistrements sont-ils différents en fin de semaine, ont-ils été réalisés par la même personne ?

- **Comment ont-elles été acquises ?** A-t-on mesuré toutes les unités de la population, ou bien a-t-on fait un échantillonnage ? Dans ce dernier cas, le choix de l'unité expérimentale peut-il avoir une influence sur le résultat ?
- **Existe-t-il une structure logique entre les observations ?** S'agit-il du poids d'un animal mesuré à plusieurs époques ou du poids d'animaux différents ? Les unités ont-elles un rapport entre elles, comme les animaux d'une même portée ou les fruits d'une même variété ? Les observations ont-elles en commun certains facteurs d'environnement ?

C'est seulement quand on a répondu à toutes ces questions, que l'on peut pousser plus loin l'analyse.

1.2 La forme du tableau de données

Il est généralement toujours possible de présenter des données sous la forme d'un tableau rectangulaire à I lignes et J colonnes ; quelquefois nous utiliserons les notations n au lieu de I et p au lieu de J . Formellement, nous pouvons écrire qu'un tableau de données \mathbf{X} est de la forme :

$$\mathbf{X} = [x_{ij}^I], i = 1, \dots, I \text{ (ou } n) ; j = 1, \dots, J \text{ (ou } p)$$

Ci-dessus, les deux indices i et j repèrent la valeur x_{ij}^I qui se trouve à la ligne i et à la colonne j . Mais la représentation d'une ligne ou d'une colonne est différente :

- une ligne i , notée $x_i = [x_i^1 \quad \dots \quad x_i^J]$, représente une unité de base : c'est l'**observation élémentaire** sur laquelle ont été mesurées, observées ou contrôlées J caractéristiques.
- une colonne j , notée $x^j = [x_1^j \quad \dots \quad x_I^j]^T$, représente la valeur soit d'une **variable**, soit d'un **indicateur**. La différence entre les deux provient de l'origine expérimentale de l'observation, c'est-à-dire de la façon dont on l'a obtenue : relever l'âge d'une personne fournit une variable, mais sélectionner des personnes d'un âge fixé, fournit un indicateur. C'est un problème que nous évoquerons plusieurs fois dans cet ouvrage, dans la mesure où la distinction peut ne pas être toujours évidente.

1.3 Notion de type

La valeur d'une variable, comme celle d'un indicateur, dépend de la structure mathématique qui permet de la définir. Est-elle définie sur l'ensemble des nombres réels \mathbb{R} ou simplement sur l'ensemble des réels positifs \mathbb{R}^+ ? Ne peut-elle prendre qu'un nombre limité de valeurs entières ? Quelquefois, la valeur ne correspond pas à une structure numérique : on peut toutefois s'y ramener par un simple **codage**. Par exemple, pour une couleur qui peut ne présenter que les trois caractéristiques vert, jaune ou marron, on peut soit établir la correspondance suivante :

$$\{\text{vert, jaune, marron}\} \Leftrightarrow \{1, 2, 3\}$$

¹ Dans la terminologie du modèle linéaire, on emploie le mot **facteur**.

soit utiliser pour chaque réalisation une variable à deux valeurs seulement, 0 ou 1. Chaque modalité de la variable définit une variable dichotomique, la correspondance est alors :

TABLEAU 1 - Exemple de passage d'un codage simple à un codage disjonctif complet.

	codage initial	vert	jaune	marron
vert :	1	1	0	0
jaune :	2	0	1	0
marron :	3	0	0	1

Pour un traitement statistique, on utilise ce dernier codage appelé **codage disjonctif complet**. Ainsi un objet vert sera analysé par l'ensemble des trois valeurs (1,0,0), un jaune par (0,1,0) et un marron par (0,0,1), comme indiqué dans le tableau 1.

Physiquement, une mesure numérique sur \mathbb{R} ne peut prendre qu'un nombre fini de valeurs. En pratique, ce qui distingue les valeurs sur \mathbb{R} de celles sur \mathbb{N} dépend essentiellement du nombre de valeurs possibles : si ce nombre est assez grand (> 10), on pourra considérer que la mesure est définie sur \mathbb{R} , sinon elle l'est sur \mathbb{N} ; le nombre d'enfants d'une famille est généralement défini sur \mathbb{N} alors que la taille de la portée d'un animal prolifique peut l'être sur \mathbb{R} . Néanmoins, rien ne s'oppose à dire que le nombre moyen d'enfants d'une famille européenne est de 1,8, bien que cette valeur n'ait aucun sens pour une famille particulière. La limite entre la définition sur \mathbb{R} ou sur \mathbb{N} dépend du contexte de l'application qui devra toujours être précisé. Cette remarque ne s'applique pas à un codage ; elle n'a de sens que pour une caractéristique sur laquelle on peut faire une opération arithmétique comme la sommation, donc un calcul de moyenne. Pour une utilisation statistique, il est utile de remplacer la notion de structure mathématique par celle de type de variable, schématisée dans le tableau 2.

TABLEAU 2 - Les différents types de variables ou d'indicateurs, nature et exemples.

variable ou indicateur	Type	Nature	Exemples
quantitative	continu		<i>poids, taille, revenu</i>
	entier	discrète	<i>nombre d'enfants</i>
	polytomie ordonnée	ordinale	<i>(mauvais,moyen,bon), (malade,atteint,sain)</i>
qualitative	polytomie non ordonnée	nominale	<i>profession, couleur, variété végétale</i>
	dichotomie	binaires	<i>(oui/non), (présence,absence)</i>

1.4 Représentation d'un tableau de données

Le tableau X est donc constitué par deux sous-tableaux :

$$X = [X^1 : X^2]$$

X^1 est un tableau à l lignes et p colonnes ; X^2 est un tableau à l lignes et q colonnes. Donc X est un tableau à l lignes et $j = p + q$ colonnes. L'unité expérimentale n^i_l , la ligne i de, X est donc formée d'un ensemble de p variables et de q indicateurs ; les p variables définissent le **vecteur observation** X^1_i lui-même éventuellement repéré par un **vecteur de contrôle** X^2_i . L'ensemble des méthodes statistiques se rattache :

- soit à des **méthodes unidimensionnelles** pour l'analyse séparée de chaque colonne de X^1 ;
- soit à des **méthodes multidimensionnelles**² pour l'analyse simultanée de l'ensemble des colonnes de X^1 .

Dans les deux cas, l'existence d'un vecteur de contrôle induit une structure dans le tableau de données : connue *a priori*, elle sera généralement utilisée dans l'analyse des données ; mais elle pourra aussi ne fournir que des informations non prises en compte dans l'analyse, néanmoins utiles pour l'interprétation.

Quelquefois, on peut compléter un premier tableau X par d'autres observations sur les mêmes j colonnes. Si ces observations ne sont pas prises en compte dans l'analyse initiale, on dira que l'on a affaire à des **observations passives** ou **observations supplémentaires**, par opposition aux premières qui sont des **observations actives**. Naturellement, on peut aussi ajouter d'autres colonnes à X sur les mêmes unités de base ; on parlera de **variables passives** ou **variables supplémentaires** (les secondes) et de **variables actives** (les premières). La même chose peut se produire pour des indicateurs. Schématiquement nous pouvons écrire que le tableau X est complété de la manière suivante :

$$\left[\begin{array}{c|c} X & X^{c\sup} \\ \hline X^{c\sup} & \end{array} \right]$$

Bien que nous ne l'ayons pas indiqué, le symbole « / » peut être remplacé par des variables supplémentaires associées à des observations supplémentaires.

2. RESUMES UNIDIMENSIONNELS

Nous débuterons ce paragraphe par l'étude des variables quantitatives ; puis nous verrons ensuite de manière plus succincte au §2.5 le cas des variables qualitatives. Avant d'étudier un tableau comme X^1 , il est toujours indispensable d'analyser séparément chacune de ses colonnes ; pour des variables quantitatives, on commence par étudier séparément chaque **distribution** des variables du tableau, soit p distributions. Il existe plusieurs manières de représenter ces distributions. A titre d'exemple regardons le tableau 3 qui représente 20 valeurs de 6 variables chimiques de composition d'eau minérale. On peut étudier ce tableau de deux manières : à l'aide de graphiques ou à l'aide de résumés numériques ; nous commencerons par les premiers.

Dans ce paragraphe nous allons considérer les colonnes de X^1 l'une après l'autre, même si ne nous interdisons pas des regards croisés ; aussi, pour éviter l'usage d'indices superflus, notons x^i ($i = 1, \dots, p$) ; les n valeurs de la colonne j (n est noté l au début du §1.4).

L'analyse des colonnes de X^2 est aussi possible, s'il existe une structure dans le tableau de données. Elle permet de voir si cette structure induit un équilibre entre les différents groupes définis par X^2 ; ainsi, lorsque X^2 n'a qu'une colonne (alors $q = 1$), si la variable a k modalités, l'écriture disjonctive de X^2 comportera k colonnes associées aux k

² On trouve aussi l'adjectif *multivarié*, provenant directement du *multivariate* anglais.

variables définies comme dans le tableau 1. Par exemple, si c'est une couleur : existe-t-il autant d'observations dans chaque niveau de couleur ? Nous verrons plus loin que l'existence d'un équilibre est souvent un élément important pour faire les « meilleures » statistiques. Dans la suite de ce chapitre, nous n'utiliserons pas cette possibilité de structuration, nous utiliserons la notation \mathbf{X} à la place de \mathbf{X}^1 . De plus, nous noterons n le nombre lignes (=observations, au lieu de J) et p celui des colonnes (=variables).

TABLEAU 3 - **Eaux1** : Echantillon de 20 eaux minérales non gazeuses pour 6 variables chimiques.

Eau	HCO3	SO4	Cl	Ca	Mg	Na
ALX	341	27	3	84	23	2
Bac	263	23	9	91	5	3
Cay	287	3	5	44	24	23
Cha	298	9	23	96	6	11
Ce1	200	15	8	70	2	4
CyE	250	5	20	71	6	11
Evl	357	10	2	78	24	5
FeE	311	14	18	73	18	13
HiP	256	6	23	86	3	18
Lau	186	10	16	64	4	9
Oge	183	16	44	48	11	31
Ord	398	218	15	157	35	8
Per	348	51	31	140	4	14
R1b	168	24	8	55	5	9
Spa	110	65	5	4	1	3
Tho	332	14	8	103	16	5
Ver	196	18	6	58	6	13
V11	59	7	6	16	2	9
V1c	402	306	15	202	36	3
Vo1	64	7	8	10	6	8

2.1 Graphiques

Nous verrons plus loin que les méthodes statistiques dépendent beaucoup des **suppositions** que l'on peut faire sur les variables. Il est donc indispensable de faire, avant tout calcul, des graphiques permettant de voir si ces suppositions sont plus ou moins bien respectées, c'est ce qu'on appelle une analyse exploratoire des données³. Ces graphiques permettent de répondre aux deux questions suivantes :

- La distribution des données est-elle sensiblement Normale, c'est-à-dire symétrique et relativement concentrée autour de la position centrale ?
- Existe-t-il des **valeurs suspectes**⁴ ?

On peut faire quatre types de graphiques : un **histogramme**, un **graphique de densité**, un **graphique de Normalité**⁵ et un **diagramme en boîte**⁶.

³ En anglais : *Exploratory Data Analysis* (EDA).

⁴ En anglais : *outliers*.

⁵ Qu'on appelle aussi *qq-plot*.

⁶ Qu'on appelle aussi *boîte à moustaches* ou *Box and Whisker plot* ou *box-plot*.

Un histogramme n'est parlant que si le nombre n d'observations est suffisamment élevé (>100), aussi est-il utile de le remplacer par le graphique de densité qui en est une version lissée dont nous parlerons plus loin.

Un graphique de Normalité est un graphique des valeurs ordonnées de la variable en fonction des quantiles correspondants de la loi Normale standard (de moyenne 0 et de variance 1) : la distribution est d'autant plus proche de la Normalité que le graphique est proche d'une droite.

Un diagramme en boîte est aussi très facile à interpréter : si nous appelons $Q_{0.25}$, $Q_{0.50}$ (qui est la **médiane**) et $Q_{0.75}$ les valeurs telles que le quart, la moitié et les trois-quarts des observations leur soient inférieures (ce que nous appellerons plus loin les quantiles) on trace une boîte comme celle de la figure 1 ; les données sont celles de la variable `Ca` du fichier **Eaux1** (Tab.3). Les extrémités indiquées par « | » sont à des positions : pour l'inférieure à $Q_{0.25} - 1.5(Q_{0.75} - Q_{0.25})$ et pour la supérieure à $Q_{0.25} + 1.5(Q_{0.75} - Q_{0.25})$ et les valeurs qui sont extérieures sont indiquées par des « * ». On repère ainsi très facilement les valeurs suspectes ou extrêmes, et on voit très rapidement si la distribution est symétrique. Ici $Q_{0.25} = 53.25$; $Q_{0.50} = 72.00$ et $Q_{0.75} = 92.25$. Apparaissent à droite deux observations suspectes correspondant aux valeurs 157 et 202.

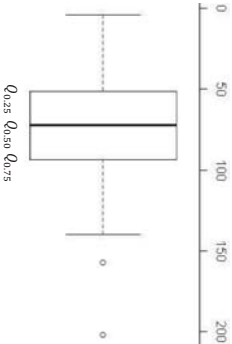


FIGURE 1 Représentation schématisée d'une boîte à moustaches pour la variable `Ca` du fichier **Eaux1**.

La comparaison de plusieurs distributions est facilitée. Ainsi si nous regardons les deux premières variables du fichier de données **Eaux1** la distribution de `HCO3` apparaît assez symétrique, alors que celle de `SO4` ne l'est pas du tout. Les quatre autres boîtes à moustaches donnent une vision plus claire de la dissymétrie (Tab.4A, Fig.2).

TABLEAU 4A - **Eaux1** : Boîtes à moustaches.

```
> library(MASS); library(car); library(corr); library(facto); library(ggplot2)
> Eaux1<-read.table("Eaux1.txt",h="T",row.names=7)
# Tracé des six boîtes à moustaches (Fig.2)
> boxplot(Eaux1)
```

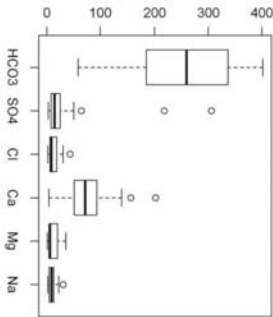


Figure 2. **Eaux1** : Boîtes à moustaches des six variables.

Un regard sur les graphiques complets de ces mêmes variables montre bien les différences entre les deux : HCO_3 a une distribution voisine de la Normale, alors que SO_4 est très dissymétrique avec un fort étalement vers les grandes valeurs (Tab.4B, Fig.3). Les graphes de Normalité vont dans le même sens : la forme de la distribution de HCO_3 est beaucoup proche de celle de la loi Normale que celle de SO_4 (Fig.3).

```
TABLEAU 4B - Eaux1 : programmes permettant de faire les graphiques des figures 2 et 4.

# Tracé sur un même graphe de l'histogramme et de la densité de la
# gaussienne correspondante
> Dist.forme<-function(x)
+ {
+   par(mfrow=c(2,2))
+   ## pb données manquantes
+   x.naomi<-na.omit(x)
+   hist(x.naomi, col="gray", prob=TRUE, xlim=c(min(x.naomi),
+   max(x.naomi)), main="")
+   curve(dnorm(x,mean=mean(x.naomi),sd=sd(x.naomi)),add=TRUE,lwd=2,
+   col="red")
+   boxplot(x.naomi)
+   logc<-summary(x.naomi)[5] - summary(x.naomi)[2]
+   points(mean(x.naomi), col = "orange", pch = 18)
+   plot(density(x.naomi,width=2*logc), xlab="x", ylab="", type="l", main="")
+   qqnorm(x.naomi)
+   qqline(x.naomi)
+ }
# Tracé des quatre graphiques réalisés sur les variables HCO3 et SO4
# (Fig. 3)
> Dist.forme(Eaux1[,1])
> Dist.forme(Eaux1[,2])
```

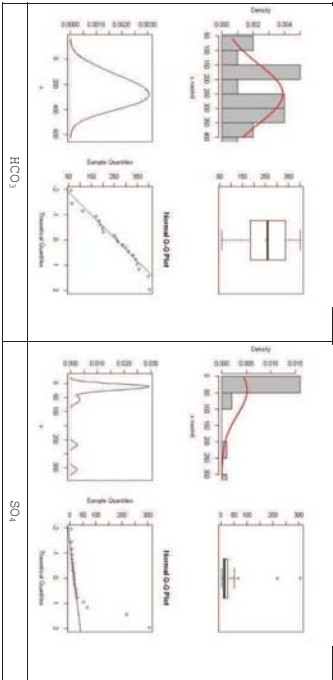


Figure 3. **Eaux1** : Graphiques EDA des deux variables HCO_3 et SO_4 .

2.2 Paramètres numériques classiques

Pour compléter ces visions graphiques, il faut calculer des paramètres numériques traduisant le mieux possible la distribution, c'est-à-dire les caractéristiques de ces n valeurs ; on peut regrouper ces caractéristiques selon trois types d'indics : **position**, **dispersion** et **forme**. Pour chaque type, plusieurs paramètres sont possibles avec des intérêts pratiques plus ou moins importants.

2.2.1 Paramètres de position : ils définissent la tendance générale de la distribution. Deux paramètres principaux sont utilisés, avec des propriétés différentes : la moyenne et la médiane d'une variable quantitative x . La **moyenne**, notée \bar{x} , est définie par :

$$\bar{x} = \sum_{i=1}^n x_i / n$$

Si on a n_i valeurs x_i identiques, et si le nombre total est $n = \sum n_i$, on associe à chaque observation une **pondération** $p_i = n_i/n$, et donc $\sum p_i = 1$; alors on écrit :

$$\bar{x} = \sum_{i=1}^n p_i x_i$$

Naturellement, quand toutes les observations sont individualisées, $p_i = 1/n$. On appelle **variable centrée** la variable x_c de coordonnées : $x_i - \bar{x}$; la moyenne de cette variable est donc nulle.

La **médiane**, notée Me , est la valeur qui divise les n observations en deux parties égales (50 % lui sont inférieures et 50 % supérieures), elle est déterminée lorsque l'on a rangé les n observations par ordre croissant :

$$x_{[1]} \leq x_{[2]} \leq \dots \leq x_{[n-1]} \leq x_{[n]}$$

alors si n est impair et égal à $2m + 1$, la médiane vaut :

$$Me = x_{(m+1)}$$

et si n est pair et égal à $2m$, la médiane vaut :

$$Me = (x_{(m)} + x_{(m+1)})/2$$

Cette formulation montre que le calcul de la médiane revient à pondérer les valeurs rangées de la façon suivante : toutes les valeurs ont une pondération nulle $p_{(j)} = 0$, sauf pour :

- n impair, alors $p_{(m+1)} = 1$ et
- n pair $p_{(m)} = p_{(m+1)} = 1/2$

En fait, la médiane est un paramètre beaucoup plus stable que la moyenne et beaucoup moins sensible qu'elle aux valeurs suspectes car les valeurs extrêmes n'ont pas de poids dans son calcul. On dit que c'est un paramètre **robuste** (Tab.5A).

Tableau 5A - **Baux1** : détermination de la médiane de HCO₃.

> Baux1[,1] # valeurs de HCO3 dans le fichier	
[1] 341 263 287 298 200 250 357 311 256 186 183 398 348 168 110 332 196 59 402 64	
> rank(Baux1[,1]) # rang des valeurs de HCO3 dans le fichier	
[1] 16 11 12 13 8 9 18 14 30 6 5 19 17 4 3 15 7 1 20 2	
> sort(Baux1[,1]) # valeurs rangées par ordre croissant de HCO3	
[1] 59 64 110 168 183 186 196 200 250 256 263 287 298 311 332 341 348 357 398 402	
> median(Baux1[,1]) # médiane de HCO3	
[1] 259.5	

La dernière remarque conduit à imaginer d'autres paramètres en jouant de façon moins drastique sur les pondérations des valeurs rangées. On peut ainsi calculer une **moyenne équilibrée**⁷, en supprimant un certain nombre de valeurs extrêmes, en général un pourcentage de n , par exemple 5 % ; si nous prenons un pourcentage de 50 % nous obtenons la médiane. Voici ce que donnent ces calculs sur les deux variables HCO₃ et SO₄ (Tab.5B) :

Tableau 5B - **Baux1** : Moyennes et moyennes équilibrées des deux variables HCO₃ et SO₄.

HCO ₃		SO ₄	
> mean(Baux1[,1])		> mean(Baux1[,2])	
[1] 250.45		[1] 42.4	
> mean(Baux1[,1],trim=0.5)		> mean(Baux1[,2],trim=0.5)	
[1] 259.5		[1] 14.5	
> mean(Baux1[,1],trim=0.05)		> mean(Baux1[,2],trim=0.05)	
[1] 252.6667		[1] 29.94444	

Une autre idée est celle de la **winsorization** qui ne modifie pas la pondération, mais qui donne à un certain nombre (ou à un certain pourcentage) de valeurs extrêmes la même valeur. Par exemple, on donne à $x_{(1)}$ la même valeur que celle de $x_{(n/2)}$ et à $x_{(n)}$ la même valeur que celle de $x_{(n-1)}$. Ces différents paramètres ont des propriétés intéressantes et, quelquefois, ils peuvent remplacer la moyenne.

Quand les observations sont discrètes ou regroupées en classes de même amplitude, on définit le **mode** comme la valeur ou la classe de fréquence maximale. Pour

⁷ En anglais : *trimmed mean*.

une variable quantitative, c'est la valeur la plus probable. Mais contrairement aux autres paramètres de position, le mode peut ne pas être unique : on parlera alors de **distribution multimodale**. Cet inconvénient peut devenir un avantage, puisque la présence de deux valeurs de mode peut être l'indice d'un mélange provenant de deux populations différentes. Pour une distribution Normale : Moyenne = Médiane = Mode.

Une grande différence entre ces paramètres (en particulier entre moyenne et médiane) est un indice certain de non-Normalité de la distribution.

2.2.2 Paramètres de dispersion : ils définissent la variabilité de la distribution. Le plus utilisé est la **variance** notée s^2 , définie par :

$$s^2 = var(x) = \sum_{i=1}^n (x_i - \bar{x})^2 / (n - 1)$$

La division par $(n - 1)$ sera justifiée ultérieurement pour des raisons statistiques ; quelquefois la division est faite par n , la variance est alors la moyenne des carrés des écarts à la moyenne⁸. Nous utiliserons toujours le facteur $(n - 1)$, sauf cas particulier justifié par le contexte ; de toute manière, dès que n est assez grand la différence est minime. On note si l'**écart-type** qui est la racine carrée de la variance, et s'exprime donc dans la même unité que la variable étudiée (Tab.5C).

Tableau 5C - **Baux1** : variances et écart-types des deux variables HCO₃ et SO₄.

HCO ₃		SO ₄	
> var(Baux1[,1],na.rm=T)		> var(Baux1[,2])	
[1] 10529.63		[1] 6079.516	
> sd(Baux1[,1])		> sd(Baux1[,2])	
[1] 102.6140		[1] 77.97125	

On appelle **variable centrée réduite** notée x_{cr} la variable de coordonnées : $(x_i - \bar{x})/s$; la moyenne de cette variable est donc nulle et son écart-type vaut 1.

Un autre paramètre de dispersion est l'**étendue** notée E , définie par :

$$E = x_{(n)} - x_{(1)}$$

L'étendue est donc la différence entre la valeur maximale et la valeur minimale des x_i : c'est un paramètre très sensible aux **valeurs extrêmes** ; on peut lui substituer des paramètres dérivés de la **distribution équilibrée**. En pratique, on considère surtout des paramètres dérivés des **α -percentiles** définis par la valeur de la variable telle que la proportion α ($\alpha < 0.5$) des observations lui soit inférieure. Si on prend la valeur entière (arrondie à l'entier le plus proche) de $v = \alpha n$, alors l' α - percentile $Q_\alpha = x_{(v)}$. L'importance des queues de distribution peut être étudiée par l'écart inter α - percentile :

$$D_\alpha = Q_{1-\alpha} - Q_\alpha$$

⁸ En anglais : *standard deviation (sd)*.

On peut démontrer, sous des conditions très larges (Kendall & Stuart, 1958), des propriétés des distributions de ces quantités ; en particulier, pour la médiane ($\alpha = 0.5$), on obtient la variance :

$$var(Me) = \frac{\pi}{2} var(\bar{X})$$

Ce résultat nous indique que la variance de la médiane est moins précise que celle de la moyenne.

Un paramètre de dispersion, souvent utilisé et beaucoup plus stable que l'étendue, est l'écart interquartile :

$$D_{0.25} = Q_{0.75} - Q_{0.25}$$

C'est celui-ci qui est représenté dans le graphique de la boîte à moustaches.

Nous mettons à part un autre paramètre de mesure de la variabilité, appelé **coefficient de variation**, qui est le rapport exprimé en % de l'écart-type à la moyenne :

$$CV(\%) = 100s/\bar{x}$$

Il peut être intéressant pour comparer deux échantillons d'une même variable ; de plus, étant sans dimension, il donne une mesure absolue de la variabilité. Néanmoins, il peut être d'un usage dangereux. En effet, la valeur de la moyenne est modifiée si on ajoute une constante à la mesure d'une variable ; ce n'est pas le cas de la valeur de l'écart-type. Donc CV sera différent selon que l'on fournit une température en degré Celsius ou en degré Fahrenheit, puisqu'une température de t degrés Fahrenheit correspond à $5(t - 32)/9$ degrés Celsius.

2.2.3 Paramètres de forme : Leur intérêt provient principalement du fait que leur valeur théorique pour la loi Normale est connue. On les utilise donc beaucoup pour vérifier le caractère gaussien d'une distribution ce sont :

- le coefficient d'**asymétrie**⁹ ou d'*obliquité* estimé par :
$$b_1 = \frac{m_3}{m_2^{3/2}} = \frac{1}{n} \left[\sum_{i=1}^n (x_i - \bar{x})^3 \right] / \left(\frac{1}{n} \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right] \right)^{3/2}$$
 - < 0 si mode > \bar{x} ; donc une queue de distribution étalée vers la gauche ;
 - = 0 ; pour une distribution symétrique ;
 - > 0 si mode < \bar{x} ; donc une queue de distribution étalée vers la droite.
- Notons qu'une forte asymétrie dans une distribution est un indicateur d'hétérogénéité au sein de la population étudiée.
- le coefficient d'**aplatissement**¹⁰ :
$$b_2 = \frac{m_4}{m_2^2} = \frac{1}{n} \left[\sum_{i=1}^n (x_i - \bar{x})^4 \right] / \left(\frac{1}{n} \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right] \right)^2$$

⁹ En anglais : *skewness*.
¹⁰ En anglais : *kurtosis*.

- < 3 si moins aplatie qu'une loi Normale ;
- = 3 ; si loi Normale ;
- > 3 si plus aplatie qu'une loi Normale.

Des valeurs respectives des coefficients d'asymétrie et d'aplatissement différentes de 0 et 3 indiquent que la distribution du caractère est plus ou moins éloignée de la distribution Normale ; évidemment la taille de l'échantillon joue un rôle si on souhaite faire un test de l'écart à la Normalité (Tab.6A).

TABLEAU 6A - **Eaux1** : résumé par variable et calcul des coefficients d'asymétrie et d'aplatissement.

```
# Résumé par variable
> summary(Eaux1)
# Ecart-type, asymétrie et aplatissement
> library(e1071)
# Bibliothèque nécessaire pour les fonctions skewness & kurtosis
> for (i in 1:6 )
+ {
+   print(sd(Eaux1[, i], na.rm=T))
+   print(skewness(Eaux1[, i], na.rm=T))
+   print(kurtosis(Eaux1[, i], na.rm=T))
+ }
```

6811. L'ensemble de ces résultats pour les six variables **Eaux1** est fourni au tableau

TABLEAU 6B - **Eaux1** : paramètres des 6 variables.

	HCO3	SO4	Cl	Ca	Mg	Na
X ₀₁	59.0	3.00	2.00	4.00	1.00	2.00
Q _{0.25}	185.2	8.50	6.00	53.25	4.00	4.75
Me	259.5	14.50	8.50	72.00	6.00	9.00
\bar{x}	250.4	42.40	13.65	77.50	11.85	10.10
Q _{0.75}	394.2	24.75	18.50	92.25	19.25	13.00
X ₀₄	402.0	306.00	44.00	202.00	36.00	31.00
s	102.61	77.97	10.58	48.09	11.09	7.32
b1	-0.33	2.45	1.24	0.76	0.93	1.24
b2-3	-1.02	4.84	1.02	0.39	-0.55	1.14

2.3 De l'intérêt des transformations des variables

Tous les paramètres présentés ci-dessus permettent de définir des indices synthétiques d'une distribution. Nous verrons ultérieurement que de nombreuses méthodes d'analyse statistique, notamment les méthodes paramétriques, se basent sur l'hypothèse de Normalité de la distribution de la variable étudiée. Même si on peut bien souvent ne pas en faire un usage exclusif en ayant recours à des méthodes statistiques qui

¹¹ Les valeurs des coefficients d'asymétrie et d'aplatissement n'ont ici qu'une signification limitée : il faudrait davantage d'observations pour se rendre compte si l'éloignement à la loi Normale a ou n'a pas un sens.

ne requièrent pas la Normalité des variables (par exemple les méthodes non paramétriques) ; il est souvent bien utile de s'en rapprocher en particulier en faisant une **transformation normalisatrice** des variables. Le type de transformation à adopter dépendra bien sûr de l'allure de la distribution de fréquences des données brutes. En particulier une distribution étalée vers les grandes valeurs, donc dissymétrique vers la droite, peut suggérer que si l'on utilisait une *transformation logarithmique*, on serait plus proche du référent Normal. Ceci conduit naturellement à choisir la moyenne des logarithmes des valeurs comme paramètre de position ; celle-ci est la **moyenne géométrique** notée g_x , des valeurs de départ :

$$\ln(g_x) = \frac{1}{n} \sum_{i=1}^n \ln(x_i) \Rightarrow g_x = \left(\prod_{i=1}^n x_i \right)^{1/n}$$

De la même façon, si c'est l'inverse qui suit une distribution Normale, un calcul du même type fournit la **moyenne harmonique**¹² notée h_x :

$$\frac{1}{h_x} = \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \Rightarrow h_x = n \left(\sum_{i=1}^n \frac{1}{x_i} \right)^{-1}$$

Si toutes les valeurs x_i sont positives, on peut montrer que :

$$h_x \leq g_x \leq \bar{x}$$

Il est à noter que les transformations normalisatrices courantes possèdent la propriété importante de réduire l'**hétéroscélasticité** des données, c'est-à-dire de stabiliser leur variance. Dans de nombreuses situations, la transformation normalisatrice n'opère pas une normalisation complète mais contribue à rendre la distribution de fréquences plus symétrique, ce dont on peut se satisfaire dans un certain nombre d'analyses.

Nous verrons au chapitre suivant qu'il est possible de faire un certain nombre de transformations simples avant de commencer le traitement véritable.

2.4 Estimation de la densité

Sur la figure 4, nous pouvons voir que le tracé d'un histogramme ne fournit pas de résultats très interprétables si le nombre n d'observations est trop faible. En effet pour le tracer on doit se fixer un nombre de classes, donc un intervalle de variation h dans lequel toutes les observations ont la même densité. La densité en un point est donc la valeur n_x/nh où n_x est le nombre de points de la classe. Son calcul vérifie donc la formule générale :

$$f(x) = \frac{1}{nh} \sum_{i=1}^n K \left(\frac{x - x_i}{h} \right)$$

où K est la **fonction indicatrice** de l'intervalle $[-1/2, 1/2[$: c'est-à-dire qu'elle vaut 0 hors de cet intervalle et 1 à l'intérieur. Si on applique cette formule, on peut calculer

¹² Naturellement pour ces deux moyennes $x_i > 0, \forall i$.

point par point la densité ; mais la courbe qui en résulte est peu régulière. Une idée pour la régulariser consiste à remplacer la fonction indicatrice par une autre fonction que l'on appelle **noyau**, un noyau fréquemment employé est le noyau gaussien :

$$K(u) = \frac{1}{\sqrt{2\pi}} \exp(-u^2/2)$$

Ensuite, le problème important consiste à choisir la valeur de h , que l'on appelle la **fenêtre** : trop faible elle fournira une courbe irrégulière, trop grande la courbe sera lisse. Les logiciels fournissent des options standard qui fonctionnent généralement bien.

Considérons les données archéologiques **ceram18** contenant en ligne la date de la monnaie, puis le tableau de contingence croisant les groupes techniques (soit 184 lignes) et les 132 contextes (ou ensembles stratigraphiques) étudiés en colonne (Tab.7). Les groupes techniques, aussi nommés productions, se définissent par les caractéristiques techniques des récipients que sont la nature de l'argile et le choix des couvertes (glacure, vernis, émail, engobe...) réalisés par le potier. Les contextes archéologiques correspondent aux dépôts anthropiques circonscrits dans l'espace et dans le temps, retenus pour la qualité de leur interprétation chrono-stratigraphique et fonctionnelle (niveau d'occupation d'un bâtiment, fosse dépotoir).

Tableau 7 - **ceram18** : données de base, tracés de 3 histogrammes et de la densité.

> ceram18 <- read_table("ceram18.txt", h="r", row.names=1, sep="\t")
> dim (ceram18) ; monnaie <- ceram18 [1,]
> summary (na.omit(monnaie))
monnaie
Min. : 353
1st Qu.: 1072
Median : 1316
Mean : 1188
3rd Qu.: 1468
Max. : 1643
> monnaie <- na.omit(monnaie)/100
> par (mfcow=c(2,2))
> hist (monnaie,nclass=4,xlab= "4 classes",prob=TRUE, col="lightblue",
border="pink")
> hist (monnaie,nclass=7,xlab= "7 classes",prob=TRUE, col="lightblue",
border="pink")
> hist (monnaie,nclass=10,xlab= "10 classes",prob=TRUE, col="lightblue",
border="pink")
> plot (x=c(3,16),y=c(0,0.7),type="n",bty="n",xlab="monnaie",ylab="Estimation de la densité")
> rug (monnaie)
lines(density(monnaie,width=width,\$f(monnaie,method="dpi"),n=200), lty=1)
> lines (density(monnaie,bw=0.15),col="blue")
> lines (density(monnaie,bw=0.25),col="green")
> lines (density(monnaie,bw=0.5),col="red")
> lines (density(monnaie,bw=1),col="pink")

L'examen des trois histogrammes (respectivement à 4, 7 et 10 classes) est peu révélateur de la forme de la distribution ; on soupçonne que la distribution de la datation des ensembles stratigraphiques (31 datés) est assez différente d'une gaussienne. L'examen des trois histogrammes est peu révélateur de la forme de la distribution ; par

contre celui de la courbe de densité (ici avec un noyau gaussien) permet de mieux détecter deux, voire trois, grandes périodes de déposition assez bien séparées (Fig.4).

Toutes les analyses précédentes, et en particulier les examens graphiques, peuvent conduire à effectuer des **transformations de variables** pour rendre plus symétriques et moins étalées certaines distributions. On pourra ainsi prendre les logarithmes, les racines carrées et ou toute autre transformation simple. Néanmoins, il faudra prendre garde à ce que l'analyse simultanée des variables peut alors s'avérer plus délicate ; non pas au niveau technique, mais au niveau de l'interprétation. L'utilisateur saura-t-il interpréter une somme, sans doute pondérée, de ces variables transformées ? Nous en reparlerons au chapitre suivant. On peut aussi envisager de travailler directement sur les rangs de toutes les variables : on remplace pour toutes les p variables, la valeur observée par son **rang** et ensuite on applique les méthodes prévues pour des variables continues. Cette pratique, qui n'est pas celle plus traditionnelle des statistiques non paramétriques sur les rangs, fournit quelquefois des résultats fort intéressants.

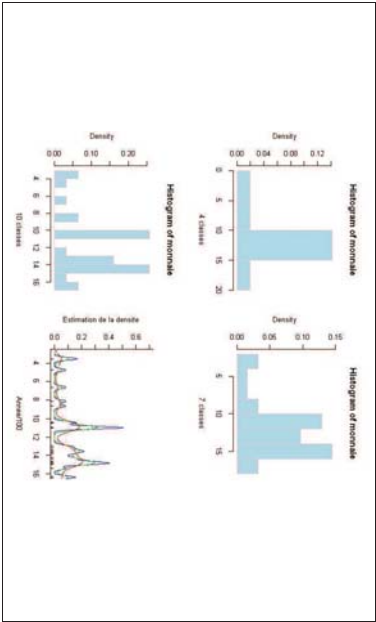


Figure 4. Cezan18 : Histogrammes et fonction de densité.

2.5 Les variables qualitatives

Nous pouvons enfin indiquer rapidement comment faire un résumé pour des variables qualitatives. Nous n'entrerons dans les détails, mais à partir d'un exemple nous allons voir comment faire une description simple. Considérons le fichier **Poumon** (72*7) provenant d'une étude sur la byssinose¹³, maladie propre aux personnes ayant travaillé sur le coton dans une atmosphère mal ventilée¹⁴. Le fichier comprend 72 lignes et 7 colonnes (2 variables et 5 facteurs) ; les 72 lignes du fichier représentent chacune une combinaison des niveaux des 5 facteurs (Everitt et Rabe-Hesketh, 2001). Les sept colonnes représentent :

- Oui : nombre de personnes malades ;
- Non : nombre de personnes non malades ;
- Pousiere : Haut, Moyen, Bas ;
- Race : Blanc, Autre ;
- Sexe : Homme, Femme ;
- Tabagie : Fumeur, Non-Fumeur ;
- Tempemploi : <10 ans, 10-20 ans, >20 ans.

On peut avoir des résumés numériques (dénombrements par niveaux de facteurs, par combinaisons de niveaux de facteurs) ou des graphiques comme des histogrammes ou des représentations en « camembert »¹⁴, cf. Tab.8.

Tableau 8 - Quelques analyses des données du fichier Poumon.

```
> Poumon<-read.table("Poumon.txt",h="?",sep="\t")
> dim(Poumon)
[1] 72 7
> attach(Poumon)
# Conversion de variables en facteurs
> Poumon$Pousiere<-
factor(Poumon$Pousiere,levels=1:3,labels=c("Haut","Moyen","Bas"))
> Poumon$Race<-factor(Poumon$Race,levels=1:2,labels=c("Blanc","Autre"))
> Poumon$Sexe<-factor(Poumon$Sexe,levels=1:2,labels=c("Homme","Femme"))
> Poumon$Tabagie<-factor(Poumon$Tabagie,levels=1:2,
labels=c("Fumeur","Non-Fumeur"))
> Poumon$Tempemploi<-factor(Poumon$Tempemploi,levels=1:3,
labels=c("<10 ans","10-20 ans","> 20 ans"))
> xtabs(Oui~Race+Pousiere)
      Race
Pousiere 1 2 3
         1 49 12 31
         2 56 6 11
         3 1 49 12 31
> summary(xtabs(Oui~Race+Pousiere))## Avec test d'indépendance ...
Call: xtabs(formula = Oui ~ Race + Pousiere)
Number of cases in table: 165
Number of factors: 2
Test for independence of all factors:
      ChiSq = 9.934, df = 2, p-value = 0.006963
> xtabs(Non~Race+Pousiere)
      Race
Pousiere 1 2 3
         1 218 843 2363
         2 346 439 1045
> library(Hmisc)
```

¹³ On l'appelle aussi *maladie brune des poumons* ou *fièvre du lundi*.
¹⁴ En anglais : *pie chart*.


```
> describe (Oui~Race+Pousiere)
3 Variables      72 Observations
-----
Oui
n missing unique Mean      .05      .10      .25      .50      .75      .90
72      0      12      2.292      0.0      0.0      1.0      3.0      4.9
.95
8.9      0      1      2      3      4      5      6      8      10      12      25      31
Frequency 3 12      10      1      1      2      1      1      1      1
%      47.17      10.14      1.11      3.11      1.11
Race
n missing unique Mean
72      0      2      1.5
1 (36, 50%), 2 (36, 50%)
Pousiere
n missing unique Mean
72      0      3      2
1 (24, 33%), 2 (24, 33%), 3 (24, 33%)

# Trouve les proportions de malade par Race et Pousiere
> prop<-propy(oui,list(race,Pousiere),sum)/apply(nen,list(race,Pousiere),sum)
> round(prop,5)
      1      2      3
1 0.223 0.014 0.013
2 0.162 0.014 0.011

# Trace de la figure 5
> barplot(prop,names=c("Haut","Moyen","Bas"),legend=c("Blanc","Autre"))
```

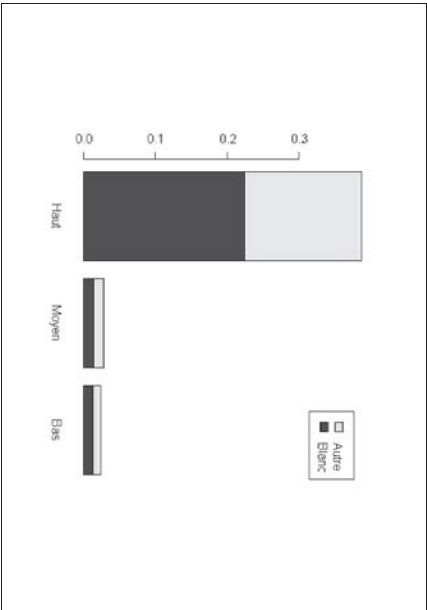


Figure 5. Graphique (barplot) de la proportion de malades suivant la race et le niveau de poussière.

3. RESUMES MULTIDIMENSIONNELS

3.1 Espace de représentation : point moyen \bar{x}

Maintenant, nous allons étudier le tableau dans son ensemble pour faire apparaître des relations entre les variables et entre les observations. Auparavant, il nous faut définir des espaces de représentation :

- celui des observations est de dimension p , noté \mathbb{R}^p . Chaque observation n^i est représentée par ses coordonnées sur les p axes de base de cet espace ; les valeurs de ces coordonnées sont celles des p variables de cette observation. C'est le vecteur ligne observation $x_i^T \in \mathbb{R}^p$.
- celui des variables est de dimension n , noté \mathbb{R}^n . Chaque variable n^j est représentée par ses coordonnées sur les n axes de base de cet espace ; les valeurs de ces coordonnées sont celles des n observations de cette variable. C'est le vecteur colonne variable $x^j \in \mathbb{R}^n$.

Le point moyen ou **barycentre** du nuage des n points, aussi appelé **centre de gravité**¹⁵, a comme coordonnées celles du vecteur moyenne \bar{x} , calculées par :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} X^T \mathbf{1}_{n \times 1} \in \mathbb{R}^p$$

où $\mathbf{1}_{n \times 1} \in \mathbb{R}^n$ est un vecteur de n lignes formé de 1.

En retranchant à chaque colonne sa moyenne, on obtient la *matrice des données centrées* :

$$X_C^T = [x_1 - \bar{x} \quad \dots \quad x_n - \bar{x}] = X^T - \bar{x} \mathbf{1}_{1 \times n} = X^T \left[\mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times 1} (\mathbf{1}_{n \times 1})^T \right] \\ = X^T \left[\mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n} \right]$$

où $\mathbf{1}_{n \times n}$ matrice formée de 1 et $\mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n}$ sont des matrices carrées de dimension n ; la première est de rang 1 et la seconde de rang $n - 1$. Ces matrices carrées sont **idempotentes** ($A \times A = A$) et symétriques $A^T = A$; ce sont des **matrices de projection**. En transposant la dernière équation, on a : $X_C = \left[\mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n} \right] X$.

L'instruction *R* permettant de centrer et/ou réduire les variables de **X** est :

```
> scale(x, center = TRUE, scale = TRUE)
```

Généralement, en option standard, toutes les valeurs d'une colonne sont centrées par rapport à leur moyenne et divisées par l'écart-type estimé ; les données manquantes éventuelles **NA** sont omises et la somme des carrés des écarts à la moyenne est divisée par $n-1$ pour calculer l'écart-type *s*.

Si toutes les observations ont la même importance, on dit qu'elles ont le même **poids**, c'est-à-dire $1/n$. Quelquefois, il peut être utile de leur donner des poids différents,

¹⁵ Aussi appelé *centroïde*.

p_i à chaque observation x_i . Ces poids sont des nombres positifs dont la somme est égale à 1 ; ils jouent un rôle analogue à celui de fréquence. On les regroupe dans une matrice diagonale **D** de dimension n :

$$\mathbf{D} = \begin{bmatrix} p_1 & 0 & \dots & 0 & 0 \\ 0 & p_2 & 0 & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \dots & 0 & p_n \end{bmatrix} \text{ avec } \sum_{i=1}^n p_i = 1$$

Avec un même poids pour chaque observation, on a bien sûr **D** = **I_n**/ n . Avec pondération, le vecteur moyenne est calculé par :

$$\bar{x} = \mathbf{X}^T \mathbf{D} \mathbf{1}_{n \times 1} \in \mathbb{R}^p$$

et la matrice centrée par :

$$\mathbf{X}_C = [\mathbf{I}_n - \mathbf{1}_{n \times 1} \mathbf{D}] \mathbf{X}$$

3.2 Première transformation : matrice de dispersion

Si maintenant, à partir de la matrice des données centrées, nous calculons :

$$\mathbf{S}_{p \times p} = \mathbf{X}_C^T \mathbf{D} \mathbf{X}_C = \left[y_{ij} = \sum_{i=1}^n p_i (x_i^j - \bar{x}^j)(x_i^k - \bar{x}^k) \right]$$

S est une matrice carrée symétrique, de dimension $p \times p$; elle contient sur la diagonale les p variances déjà vues plus haut et, hors diagonale les covariances. Elle s'appelle la **matrice de variances-covariances** ou **matrice de dispersion**. Elle va jouer un rôle fondamental dans toute l'analyse linéaire des données.

Souvent on travaille sur des données non seulement centrées, mais aussi réduites. En effet, les variables ne sont pas toujours mesurées avec les mêmes unités, il peut donc être utile de se ramener à des variables plus comparables. Pour ce faire, on les transforme pour qu'elles aient la même variabilité ; il suffit de diviser les colonnes de la matrice de données par leur écart-type. Formellement, pour y arriver nous prenons les p éléments de la diagonale de **S**, soit **diag(S)** la matrice diagonale formée par ces éléments. On a alors la matrice centrée réduite, celle que nous obtenons par la fonction **scale** de R, par :

$$\mathbf{X}_{CR} = \mathbf{X}_C (\mathbf{diag(S)})^{-1/2}$$

Cette matrice est telle que :

$$\sum_{i=1}^n |x_{CRi}| = 0 ; \sum_{i=1}^n (|x_{CRi}|)^2 = n - 1$$

Les variances sont alors toutes égales à 1. Quelquefois, il peut être utile d'avoir une matrice centrée réduite légèrement différente :

$$\mathbf{X}_{CRu} = \frac{1}{(n-1)^{1/2}} \mathbf{X}_C (\mathbf{diag(S)})^{-1/2}$$

de telle sorte que les sommes des colonnes soient nulles, et les sommes des carrés des éléments de la colonne soient égales à 1 :

$$\sum_{i=1}^n |x_{CRu}| = 0 ; \sum_{i=1}^n (|x_{CRu}|)^2 = 1$$

Enfin, on obtient la **matrice des coefficients de corrélation linéaire**, qui résume l'ensemble des dépendances linéaires entre les p variables :

$$\mathbf{R} = \mathbf{X}_{CRu}^T \mathbf{X}_{CRu} = \frac{1}{n} \mathbf{X}_{CR}^T \mathbf{X}_{CR} = \frac{1}{n} (\mathbf{diag(S)})^{-1/2} \mathbf{X}_C^T \mathbf{X}_C (\mathbf{diag(S)})^{-1/2}$$

Sous R, on peut utiliser la figure 6, appelée faute de mieux **scatterplot**¹⁶ pour résumer un certain nombre de caractéristiques unidimensionnelles et multidimensionnelles du corpus étudié. Les lignes de code suivantes permettent de l'obtenir sous R (cf. Tab. 9).

TABLEAU 9 - **Eaux1** : tracé des caractéristiques unidimensionnelles et multidimensionnelles.

```
#SCATTERPLOT
> Histogrammes sur la diagonale
> panel.hist <- function(x, ...)
+ {
+   usr <- par("usr"); on.exit(par(usr))
+   par(usr = c(usr[1:2], 0, 1.5) )
+   h <- hist(x, plot = FALSE, col="lightblue")
+   breaks <- h$breaks; nb <- length(breaks)
+   y <- h$counts; y <- y/max(y)
+   rect(breaks[-nb], 0, breaks[-1], y, col="cyan", ...)
+ }
# Coefficients de corrélation (en valeur absolue) sur la partie haute,
# Taille de police proportionnelle au coefficient de corrélation
> panel.cor <- function(x, y, digits=2, prefix="", cex.cor, ...)
+ {
+   usr <- par("usr"); on.exit(par(usr))
+   par(usr = c(0, 1, 0, 1))
+   x <- abs(cor(x, y))
+   txt <- format(c(x, 0.123456789), digits=digits) [1]
+   txt <- paste(prefix, txt, sep="")
+   if(missing(cex.cor)) cex.cor <- 0.8/strwidth(txt)
+   text(0.5, 0.5, txt, cex = cex.cor * x)
+ }
# Regression linéaire (lm) sur la partie basse ;
# d'où le signe du coefficient de corrélation.
> panel.cor <- function(x, y, digits=2, prefix="", cex.cor, ...)
+   points(x,y)
+   abline(lm(y~x))
+   #lines (lowess(y~x), col="red")
+ }
# Fig. 6
> pairs(Eaux1, diag, panel=panel.hist, cex.labels = 2,
+ font.labels=2, upper.panel=panel.cor, lower.panel=panel.lm )
```

¹⁶ On pourrait parler de *nuage de points* ou de *dispersion* des points.

La figure 6 produit une matrice de *scatterplots* contenant sur sa diagonale l'histogramme de chaque variable ; sur sa partie inférieure, les graphiques des couples avec droite des moindres carrés et sur sa partie supérieure, les coefficients de corrélation indiqués avec une taille de police plus ou moins importante selon l'intensité de leur valeur absolue.

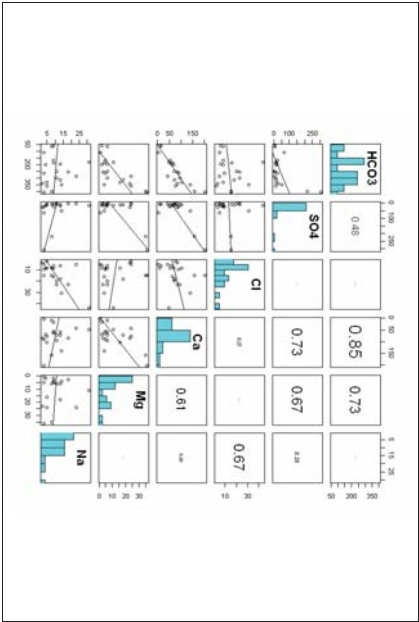


FIGURE 6. [Baux1](#) : *scatterplot*.

3.3 L'espace des observations \mathbb{R}^p : regards sur les lignes, inertie

Un des objectifs des analyses de tableau de données comme **X**, consiste bien souvent à regarder si deux lignes du tableau par exemple les deux premières, c'est-à-dire deux vecteurs ligne x_1 et x_2 de \mathbb{R}^p sont proches ou éloignées. Il est donc important de calculer une distance entre elles. Quand les unités de mesure des variables sont les mêmes, la façon la plus simple consiste à appliquer le théorème de Pythagore qui donne le carré de la distance entre deux observations :

$$d^2(x_1, x_2) = (x_1^1 - x_2^1)^2 + \dots + (x_1^p - x_2^p)^2 = (x_1 - x_2)^T (x_1 - x_2)$$

En statistique, le problème est un peu plus compliqué car, bien souvent les unités de mesure sont différentes ; on peut avoir un mélange de mensurations, de poids et d'âge. Il faut donc pondérer, d'une certaine manière, les différents carrés de l'expression précédente. De plus, la distance calculée n'a de sens que pour des axes perpendiculaires ; on aurait aussi bien pu prendre des axes obliques.

On utilise donc une formulation un peu plus générale en introduisant une matrice définie positive **Q**, symétrique de dimension $p \times p$. On dit que l'on a muni l'espace des

observations d'un **produit scalaire** : $\langle x_1, x_2 \rangle_Q = x_1^T Q x_2$, et le carré de la distance est définie par la forme quadratique :

$$d^2(x_1, x_2) = (x_1 - x_2)^T Q (x_1 - x_2)$$

Prendre **Q** = **I_p**, revient à conserver les unités initiales de chaque variable et, donc, à rendre dominantes les variables de plus grande variabilité.

Généralement, on prend **Q** = **(diag(S))⁻¹** qui est la métrique inverse des variances. Les distances sont alors calculées non plus sur **X**, mais sur la matrice centrée réduite **X_{cn}** : on impose une variance égale à 1 à chaque variable. Utiliser une métrique diagonale quelconque :

$$Q_a = \text{diag}(a_1, \dots, a_p)$$

revient à multiplier les variables par $\sqrt{a_j}$ et à utiliser ensuite la métrique **Q** = **I_p**.

On peut généraliser ce résultat à une métrique quelconque. On sait que toute matrice **Q** peut s'écrire **Q** = **T^TT**. Le produit scalaire entre deux observations s'en déduit :

$$\langle x_1, x_2 \rangle_Q = x_1^T Q x_2 = x_1^T T^T T x_2 = (T x_1)^T (T x_2)$$

La suite revient à utiliser la métrique **I_p** sur les données transformées **XT^T**.

Regardons maintenant ce que représente, pour une matrice de données **X**, munie de sa métrique **Q**, la somme I_g , appelée **inertie**, du carré des distances pondérées des vecteurs observations au point moyen du nuage :

$$I_g = \sum_{i=1}^n p_i d^2(x_i, \bar{x}) = \sum_{r=1}^n \sum_{s=1}^n p_r p_s (x_r - x_s)^T Q (x_r - x_s), r < s$$

La forme de l'inertie I_g nous dit que la somme pondérée du carré des distances des observations au point moyen est égale à la somme pondérée du carré des distances entre tous les couples d'observations. Regardons ce que cela donne pour le cas d'une seule variable ($p = 1$). Avec la définition de la variance (dans la formule dominant s^2 , le diviseur est $n - 1$ et non pas n) I_g est égale à $((n - 1)/n) \text{var}(x)$. La variance est donc une mesure de la cohésion interne des observations. Vérifions le sur un exemple ; avec les six observations : (19; 23; 24; 26; 28; 30) de moyenne $\bar{x} = 25$ et la variance vaut $76/5$; la somme de tous les carrés des différences entre les six valeurs vaut 456 ($=76 \times 6$).

L'inertie peut aussi s'exprimer en fonction de la trace de la matrice **QS** ou de celle de la matrice **sQ** :

$$I_g = \text{tr} \left(\sum_{i=1}^n p_i (x_i - \bar{x})^T Q (x_i - \bar{x}) \right) = \text{tr} \left(\sum_{i=1}^n Q p_i (x_i - \bar{x})^T (x_i - \bar{x}) \right) = \text{tr}(Q S)$$

Donc :

- si **Q** = **I_p** l'inertie est égale à la somme des p variances,
- si **Q** = **(diag(S))⁻¹**, l'inertie est égale à p , le nombre de variables ; elle est indépendante du nombre de variables.

3.4 L'espace des variables \mathbb{R}^n : regards sur les colonnes

Regardons maintenant le **vecteur variable** $x^j \in \mathbb{R}^n$. Que représente la proximité de deux vecteurs x^j et x^k ? Nous les supposons centrés et nous prenons comme métrique la matrice diagonale des poids **D** définie positive. Nous voyons que le produit scalaire entre ces deux vecteurs colonne $(x^j, x^k)_D = (x^j)^T D x^k = \sum_{i=1}^n p_i (x_i^j - \bar{x}^j)(x_i^k - \bar{x}^k)$ est proportionnel à la covariance $((n-1)/(n)s_j s_k)$ entre les deux variables. Le carré de la **D** norme (c'est-à-dire le carré de sa longueur) de chaque vecteur est proportionnelle à leur variance : $\|x^j\|_D^2 = ((n-1)/(n)s_j^2)$, leur longueur est donc égale à l'écart-type. L'angle θ_{jk} entre les deux vecteurs est donné par :

$$\cos(\theta_{jk}) = \frac{\langle x^j, x^k \rangle_D}{\|x^j\|_D \|x^k\|_D}$$

Le cosinus de l'angle entre deux vecteurs variable centrés est donc égal au coefficient de corrélation entre ces deux variables.

En conclusion, dans l'espace des observations, on s'intéresse aux distances entre points¹⁷ et dans celui des variables, à l'angle entre variables.

3.5 Transformation linéaire d'un tableau de données

Les analyses que nous allons étudier ultérieurement consistent, dans la majorité des cas, à trouver des combinaisons linéaires des p variables définissant les colonnes de **X**. Nous chercherons des combinaisons linéaires qui ont des propriétés particulières. Ces combinaisons linéaires seront des transformations des vecteurs x^j ($j = 1, \dots, p$) qui vont fournir de nouvelles variables. A chaque vecteur colonne x^j nous pouvons associer un axe de l'espace des observations \mathbb{R}^p ; une transformation linéaire consiste à projeter les observations sur un nouvel axe de vecteur unitaire ligne **a**. La coordonnée de l'observation i sur cet axe sera définie par $y_i = \mathbf{a}^T Q(x_i - \bar{x}) = (a, x_i - \bar{x})_Q$ et l'ensemble des n coordonnées par :

$$y = XQa = Xu = \sum_{j=1}^p x^j u_j \in \mathbb{R}^n \text{ avec } u = Qa$$

Nous avons donc créé une nouvelle variable avec :

- un axe de vecteur ligne unitaire $a \in \mathbb{R}^p$;
- un vecteur y de l'espace des variables \mathbb{R}^n ;
- une forme linéaire $u \in \mathbb{R}^p$ qui est appelée **facteur**.

Si, pour simplifier, nous supposons $Q = I_p$ et si nous faisons q transformations de l'ensemble du tableau $X_{n \times p}$ définies par la matrice $A = [a^1, \dots, a^q]$ telles que :

$$Y_{n \times q} = X_{n \times p} A_{p \times q}$$

on peut montrer que la matrice de dispersion S_Y de Y se déduit de S_X celle de **X** par :

¹⁷ Au chapitre consacré à la classification, nous donnerons d'autres types de distance. Nous donnerons aussi une extension de la notion de distance à la distance non plus entre observations mais à la distance entre tableaux en développant l'analyse procusteenne (chapitre 10).

$$S_{Y_{q \times q}} = A^T_{q \times p} S_{X_{p \times p}} A_{p \times q}$$

Il est souvent utile d'imposer des conditions à la matrice de transformation **A** :

- les q vecteurs unitaires : $\|a^k\|^2 = 1, k = 1, \dots, q$;
- les q vecteurs sont orthogonaux deux à deux : $(a^j)^T a^l = 0, k \neq l = 1, \dots, q$.

Quand $p = q$, la matrice **A** est ce que l'on appelle une **matrice orthogonale** : $AA^T = I_p$ et $A^T A = I_q$; la transformation est une simple **rotation** du repère de départ. Chaque vecteur observation x_i est alors dans ce cas transformé en autre vecteur observation $y_i = A^T x_i \in \mathbb{R}^p$. Chaque composante du nouveau vecteur est la projection de x_i sur les $p = q$ vecteurs ligne $a^k \in \mathbb{R}^p, k = 1, \dots, p$. Une propriété intéressante et importante est que la distance entre deux observations ligne est invariante dans une transformation linéaire :

$$(y_r - y_s)^T (y_r - y_s) = (x_r - x_s)^T A^T A (x_r - x_s) = (x_r - x_s)^T (x_r - x_s)^T$$

4. DECOMPOSITION D'UNE MATRICE DE DONNEES

Dans ce paragraphe, nous allons présenter quelques outils mathématiques d'un grand intérêt pour l'analyse d'un tableau de données. Dès maintenant, ils vont mettre en évidence la dualité entre les deux espaces de représentation. Ils permettent déjà de vérifier qu'il n'y a pas de redondance particulière entre les variables qui définissent le tableau et qu'il n'y a pas des relations trop fortes, voire fonctionnelles entre elles.

4.1 Décomposition en Valeurs Singulières (DVS)

A partir d'une matrice $X_{n \times n} = [x_i^j], i = 1, \dots, n, j = 1, \dots, p$; il est toujours possible de faire la décomposition suivante (théorème d'approximation de Schmidt (1907) et Eckart & Young (1936)) :

$$x_i^j = \theta_1 u_1^j v_1^k + \dots + \theta_r u_r^j v_r^k = \sum_{k=1}^r \theta_k u_i^k v_j^k ; \theta_1 \geq \theta_2 \geq \dots \geq \theta_r > 0$$

Cette décomposition s'appelle la **décomposition en valeurs singulières** de **X** ; les valeurs singulières de cette décomposition sont *uniques*, mais pas obligatoirement distinctes. Le nombre r de termes de **X** est le *rang* de la matrice ; il ne peut pas dépasser le plus petit des deux nombres n ou p : $r = \min(n, p)$. En général, on suppose que le nombre d'observations est supérieur au nombre de variables : $n > p$, donc $r \leq p$. Si $r = p$, on dit que la matrice est de **plein rang**. Les quantités u_i^k et v_j^k sont les composantes de vecteurs unitaires u^k et v^k de dimension respective n et p :

$$\sum_{i=1}^n (u_i^k)^2 = \sum_{j=1}^p (v_j^k)^2 = 1, k = 1, \dots, r$$

Les vecteurs u^k sont orthogonaux entre eux, de même que les vecteurs v^k , soit $U^T U = V^T V = I_r$.

On peut donc écrire :

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times r} \boldsymbol{\Theta}_{r \times r} (\mathbf{V}_{p \times r})^T$$

où :

- La matrice $\boldsymbol{\Theta}_{r \times r} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r})$ est une matrice diagonale contenant les valeurs singulières de la matrice $\mathbf{X}_{n \times p}$; c'est-à-dire les racines carrées des valeurs propres positives et non nulles¹⁸ λ_k ($k = 1, \dots, r$) de $\mathbf{X}\mathbf{X}^T$ ou de $\mathbf{X}^T\mathbf{X}$;
- $\mathbf{U}_{n \times r} = [\mathbf{u}^1 : \dots : \mathbf{u}^r]$ matrice unitaire $n \times r$, telle que \mathbf{u}^k est vecteur propre de $\mathbf{X}\mathbf{X}^T$ associé à la valeur propre non nulle λ_k ($k = 1, \dots, r$) ;
- $\mathbf{V}_{p \times r} = [\mathbf{v}^1 : \dots : \mathbf{v}^r]$ matrice unitaire $p \times r$, telle que \mathbf{v}^k est vecteur propre de $\mathbf{X}^T\mathbf{X}$ associé à la valeur propre non nulle λ_k ($k = 1, \dots, r$).

La relation fournissant \mathbf{X} signifie que ses lignes sont des combinaisons linéaires des lignes \mathbf{U} et ses colonnes des combinaisons linéaires des colonnes de \mathbf{V} . Les colonnes de \mathbf{U} forment donc une **base orthogonale** pour l'espace des colonnes \mathbb{R}^n de \mathbf{X} ; c'est-à-dire le sous-espace de dimension r de \mathbb{R}^n engendré par les colonnes de \mathbf{X} . Les colonnes de \mathbf{V} forment de même une base orthogonale pour l'espace \mathbb{R}^p des lignes de \mathbf{X} .

L'avantage de cette décomposition est exploité en statistique car elle est liée au système des valeurs propres de la matrice $\mathbf{X}^T\mathbf{X}$ (qui est directement liée, après centrage de \mathbf{X} , à la matrice de dispersion \mathbf{S}). Les θ_k sont les racines carrées des valeurs propres non nulles de $\mathbf{X}^T\mathbf{X}$ aussi bien que de $\mathbf{X}\mathbf{X}^T$. Les colonnes de \mathbf{U} sont les vecteurs propres de $\mathbf{X}\mathbf{X}^T$ et les colonnes de \mathbf{V}^T ceux de $\mathbf{X}^T\mathbf{X}$. Nous verrons que nombre de calculs statistiques font intervenir soit l'inversion soit la diagonalisation de $\mathbf{X}^T\mathbf{X}$; ces calculs s'expriment simplement en fonction des éléments de la décomposition singulière. Par exemple, voici deux expressions que nous utiliserons plus loin :

$$\begin{aligned} (\mathbf{X}^T\mathbf{X})^{-1} &= \mathbf{V}\boldsymbol{\Theta}^{-2}\mathbf{V}^T \\ \mathbf{H} &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T = \mathbf{U}\mathbf{U}^T \end{aligned}$$

La première est une autre forme de la matrice qui intervient dans la résolution des équations du modèle linéaire. On voit apparaître le rang r de la matrice des **équations normales** $\mathbf{X}^T\mathbf{X}$; c'est le nombre de valeurs propres non nulles, les r premiers éléments, de la matrice diagonale $\boldsymbol{\Theta}$. On notera que sous R , la fonction **svd** permet d'effectuer une DVS. Dans la formulation ci-dessus $\mathbf{X}^T\mathbf{X}$ et $\boldsymbol{\Theta}$ doivent être inversibles et donc $r = p$.

4.2 Inverse généralisée d'une matrice

Il est quelquefois utile de faire apparaître les relations linéaires entre les colonnes de \mathbf{X} ; c'est ce que l'on appelle la **colinéarité**. Ceci peut arriver quand, dans un fichier de données, une variable est la somme de deux autres. Quand le rang de $\mathbf{X}_{n \times p}$ vaut r , en supposant que le nombre d'observations n est plus grand que le nombre de variables p , il existe $p - r$ relations linéaires pas toujours facilement décelables. En pratique, il peut exister une **pseudo-colinéarité** qui traduit des relations étroites entre certaines colonnes

¹⁸ En pratique, seules les valeurs propres non nulles sont utilisées ; néanmoins, les valeurs propres nulles peuvent quelquefois avoir un intérêt en analyse des données.

de \mathbf{X} . Dans le premier cas, on ne peut pas inverser une matrice, dans le second le résultat est numériquement imprécis.

Pour traiter cette question, il est utile d'introduire la notion d'**inverse généralisée**, dont le calcul est simple à partir de la décomposition singulière. L'inverse généralisée, dite inverse de Moore-Penrose, d'une matrice quelconque $\mathbf{X}_{n \times p}$, notée \mathbf{X}^- , a les propriétés suivantes :

$$\left\{ \begin{array}{l} 1 : \mathbf{X}\mathbf{X}^- \mathbf{X} = \mathbf{X} \\ 2 : \mathbf{X}^- \mathbf{X}\mathbf{X}^- = \mathbf{X}^- \\ 3 : \mathbf{X}^- \mathbf{X} = (\mathbf{X}^- \mathbf{X})^T \\ 4 : \mathbf{X}\mathbf{X}^- = (\mathbf{X}\mathbf{X}^-)^T \end{array} \right.$$

Cette inverse est *unique*, et elle constitue une généralisation de l'inverse classique. Si \mathbf{X} est une matrice carrée ($n = p$) inversible, on a :

$$\mathbf{X}^- = \mathbf{X}^{-1}$$

Si le rang de \mathbf{X} est égal au nombre de colonnes ($r = p$) :

$$\mathbf{X}^- = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

Si dans la décomposition singulière de \mathbf{X} , on fait apparaître les valeurs singulières nulles :

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times n} \boldsymbol{\Theta}_{n \times p} (\mathbf{V}_{p \times p})^T ; \text{ où } \theta_{i,r+1} = \dots = \theta_p = 0$$

On peut démontrer que :

$$\mathbf{X}^- = \mathbf{V}\boldsymbol{\Theta}^{-1}\mathbf{U}^T$$

où $\boldsymbol{\Theta}^{-1}$ est l'inverse généralisée de $\boldsymbol{\Theta}$, obtenue en prenant dans cette dernière matrice l'inverse des éléments diagonaux non nuls. Donc \mathbf{X}^- joue donc un rôle analogue à \mathbf{X}^{-1} dans la résolution de systèmes rectangulaires d'équations linéaires. Si nous avons le système $\mathbf{y} = \mathbf{X}\mathbf{z}$, où $\mathbf{y} \in \mathbb{R}^n$ et $\mathbf{z} \in \mathbb{R}^p$, on peut montrer que $\mathbf{z} = \mathbf{X}^- \mathbf{y}$ est une solution particulière de ce système.

4.3 Exemples numériques

Nous allons illustrer ces résultats par deux exemples numériques que nous réutiliserons ultérieurement.

a) \mathbf{X} est de plein rang : soit la matrice \mathbf{X}_a ($n = 4$; $p = 2$) et un vecteur \mathbf{y}_a :

\mathbf{X}_a	\mathbf{y}_a
1	1
1	2
1	3
1	10

Les différents calculs des éléments de la décomposition singulière de sont présentés ci-dessous (Tab.10), ainsi que la résolution du système $\mathbf{y}_a = \mathbf{X}_a \mathbf{z}_a$.


```
> round(t(yb)**B$u[,1:2]**diag(1/B$d[1:2]))**%t(B$%v[,1:2]),3)
      [,1] [,2] [,3]
[1,] 13.333 1.667 11.667
> round(ginvreeoe(yb),3)
      [,1] [,2] [,3] [,4] [,5] [,6]
[1,] 0.111 0.111 0.111 0.111 0.111 0.111
[2,] 0.222 0.222 0.222 -0.111 -0.111 -0.111
[3,] -0.111 -0.111 -0.111 0.222 0.222 0.222
attr(,"rank")
[1] 2
```

4.4 Décomposition en valeurs singulières du triplet (X, Q, D)

La plupart des méthodes d'Analyse Factorielle des Données peuvent être présentées dans un cadre commun : celui de l'extension du théorème de la DVS au cadre d'espaces euclidiens généraux. Elle correspond à la DVS du triplet (X, Q, D) où :

- X matrice n x p des données ;
- Q est la matrice d'un produit scalaire de \mathbb{R}^p (l'espace des individus), soit une matrice carrée symétrique définie positive qui définit la fonction : $(x,y) = ([x_1 \dots x_p]^T, [y_1 \dots y_p]^T) \in \mathbb{R}^p \times \mathbb{R}^p \mapsto x^T Q y = \langle x,y \rangle_Q \in \mathbb{R}$
- D est la matrice d'un produit scalaire de \mathbb{R}^n (l'espace des variables), soit une matrice carrée symétrique définie positive qui définit la fonction :

$$(x,y) = \left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \right) \in \mathbb{R}^n \times \mathbb{R}^n \mapsto x^T D y = \langle x,y \rangle_D \in \mathbb{R}$$

DVS de (X, Q, D) : soit les matrices réelles $X_{n \times p}$ de rang r et les métriques $Q_{p \times p}$ et $D_{n \times n}$ de \mathbb{R}^p et de \mathbb{R}^n , il existe :

- Une matrice $U_{n \times r} = [u^1 \dots u^r]$ D-orthonormée, dont les colonnes sont les vecteurs propres de $X Q X^T D = W D$ associés aux valeurs propres $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$;
- Une matrice $V_{p \times r} = [v^1 \dots v^r]$ Q-orthonormée, dont les colonnes sont les vecteurs propres de $X^T D X Q = S Q$ associés aux valeurs propres $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$;
- Une matrice diagonale $\Theta_{r \times r} = \Lambda^{1/2} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r})$ contenant les r valeurs singulières du triplet (X, Q, D) telles que X se décompose en :

$$X_{n \times p} = U_{n \times r} \Theta_{r \times r} (V_{p \times r})^T = \sum_{k=1}^r \sqrt{\lambda_k} u^k (v^k)^T$$

5. BILAN

Dans ce chapitre, nous avons présenté les premiers outils de description d'un tableau de données défini par des observations et des variables. Des outils typiquement statistiques que l'on va retrouver dans tout l'ouvrage, mais aussi des outils de nature purement mathématique dont la signification en statistique apparaîtra plus loin. Pour l'instant, il est important de retenir les points les plus importants suivants :

- Avant toute analyse d'un tableau il faut bien connaître ce que représentent les observations et les variables.
- Toute étude, même celles portant sur des données multidimensionnelles, doit commencer par un examen détaillé de chaque variable, en faisant des analyses graphiques. Il ne faut pas oublier que les graphiques sont bien souvent le meilleur moyen de faire comprendre à des utilisateurs des résultats complexes ; il faut se souvenir qu'au XIX^e siècle c'est ainsi que l'illustre infirmière Florence Nightingale, plus connue sous le nom de *Dame à la lampe*, réussit à convaincre le corps médical d'améliorer ses pratiques. Première femme élue en 1859 membre de la *Royal Statistical Society*, elle est considérée comme *a true pioneer in the graphical representation of statistics* (Triola & Triola, 2012). On pourra consulter les très nombreux ouvrages sur la réalisation de graphiques (Chambers et al., 1983), en particulier avec R (Murrell, 2005).
- En particulier, les graphiques de densité et de Normalité doivent permettre de juger l'allure des distributions pour, éventuellement, transformer les variables.
- L'analyse globale du tableau, c'est-à-dire d'une matrice $X_{n \times p}$ doit être associée à la définition d'une pondération éventuelle des observations à l'aide d'une matrice diagonale $D_{n \times n}$ et d'une matrice définie positive $Q_{p \times p}$ qui permet de définir une métrique, c'est-à-dire comment calculer la mesure de ressemblance entre observations. Ces deux dernières matrices sont importantes. Ne pas s'en soucier revient à choisir implicitement les valeurs $D_{n \times n} = \frac{1}{n} I_n$ et $Q_{p \times p} = I_p$: ces choix ne sont pas toujours ceux qui permettent de faire l'analyse de la matrice de données la plus adaptée aux buts de l'étude.

Ce qu'il faut retenir :

- > Les questions préalables à toute analyse : quelles sont les unités de mesure ? Comment les données ont été acquises ? Existe-t-il une structure *a priori* dans tableau de données. Quel est le type des variables ?
- > Résumé de chaque variable : graphiques, paramètres numériques : position, dispersion, forme. Existe-t-il des valeurs suspectes et peut-on en justifier l'origine ?
- > Résumé de l'ensemble du tableau de variables : centre de gravité, matrice de dispersion.
- > Décomposition en valeurs singulières (DVS) d'une matrice de données.

PROBLEMES ET EXERCICES

1. Décrire le jeu de données **Eaux2010** dont **Eaux1** est extrait. En particulier, comparer les eaux selon leur nature (plate ou gazeuse).
2. Faire les graphiques des facteurs non analysés dans l'exemple des données **Poumon** du §2.5.
3. *Diagramme en boîte : pourquoi 1.5 ?*
Dans un *Diagramme en boîte*, les valeurs aberrantes sont définies comme les observations en dehors de l'intervalle $I = [Q_{0.25} - 1.5(Q_{0.75} - Q_{0.25}); Q_{0.25} + 1.5(Q_{0.75} - Q_{0.25})]$. Tukey, son inventeur, justifie l'introduction de la valeur 1.5 en disant que la valeur 1 est trop petite tandis que 2 est trop grande !
Supposer que les observations sont indépendantes et identiquement distribuées (*iid*) suivant une loi Normale $N(0,1)$
 - Calculer les quartiles théoriques $Q_{0.25}$ et $Q_{0.75}$ associés à $N(0,1)$ en utilisant la fonction **qnorm**.
 - En déduire l'écart interquartile $D_{0.25} = Q_{0.75} - Q_{0.25}$ et l'intervalle I .
 - Calculer $P\{X \in I\}$ quand $X \sim N(0,1)$.
 - Refaire les calculs en remplaçant 1.5 par 1 puis par 2.
 - Commenter les résultats.
4. La robustesse de la médiane et de la moyenne peuvent être étudiées en ajoutant des observations aberrantes. A l'aide du code suivant :

```
z=5 ;
x=norm(500) ;

moyenne=NULL ; mediane=NULL ;
m.max=25
for( m in 1:m.max)
{ xx=c(x,rep(z,m))
  moyenne=c(moyenne,mean(xx))
  mediane=c(mediane,median(xx))
}
```

maPlot(cbind(moyenne,mediane) , type="l", ylim=c(-0.3,.5))
abline(b=c(mean(x),median(x)) , col=1:2, lty=1:2)

 - Simuler un échantillon de taille $n = 1000$ suivant une loi $N(0,1)$ et ajouter m fois l'observation $z = 5$.
 - Représenter l'évolution de la médiane et de la moyenne en fonction de m .
 - Commenter les résultats obtenus.
5. *Nombre de classes dans un histogramme*
L'option **nclass** (ou **breaks**) de la fonction **hist** permet d'ajuster le nombre de classes n_c de l'histogramme ; par défaut fixé par la formule de Sturges à partie entière de $1 + \log_2 n$.
 - Simuler un échantillon de de taille $n = 100$ suivant une loi $N(0,1)$
 - Sur une même feuille graphique (**par(mfrow=c(3,3))**), tracer en faisant varier le nombre de classes l'histogramme correspond à l'échantillon simulé. On pourra prendre par exemple, **nclass** ∈ {3,5,8,10,15,20,25,30,50}
 - Reprendre la question b. avec un échantillon simulé suivant une loi exponentielle de paramètre 1, puis une loi de Cauchy.

- d. Commenter les résultats.
6. Matrices partitionnées : vérifiez que l'inverse de la matrice symétrique non singulière $A_{(p+1) \times (p+1)}$ avec $A_{11,p \times p}$, $a_{12,(p \times 1)}$ et a_{22} un scalaire :

$$A = \begin{pmatrix} A_{11} & a_{12} \\ a_{12}^T & a_{22} \end{pmatrix}$$

est donnée par :

$$A^{-1} = \frac{1}{b} \begin{pmatrix} bA_{11}^{-1} + A_{11}^{-1}a_{12}a_{12}^T A_{11}^{-1} & -A_{11}^{-1}a_{12} \\ -a_{12}^T A_{11}^{-1} & 1 \end{pmatrix} \text{ où } b = a_{22} - a_{12}^T A_{11}^{-1} a_{12}.$$

Une matrice non singulière de la forme $B + cc^T$, où B est non singulière a une inverse de la forme :

$$(B + cc^T)^{-1} = B^{-1} - \frac{B^{-1}cc^TB^{-1}}{1 + c^TB^{-1}c}$$



Informations sur l'eau minérale (Chambre Syndicale des Eaux Minérales (CSEM), 2008).

Fichiers : **Eaux1** & **Eaux2010**