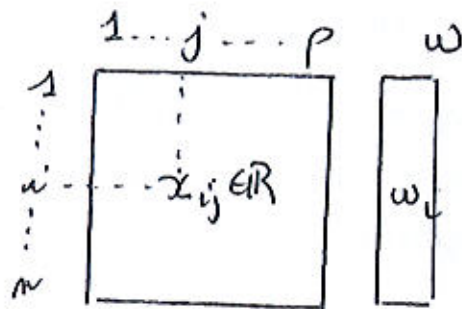


# Notions de base pour l'analyse de données quantitatives

On se place dans le cadre d'un tableau de données numériques où  $n$  individus sont décrits sur  $p$  variables.



On notera :

- \*  $X = (x_{ij})_{n \times p}$  la matrice des données "brutes" où  $x_{ij}$  = valeur de l'ième individu sur la j-ème variable
- \*  $w_i > 0$  : poids de l'individu  $i$
- \*  $x_i^t = (x_{i1}, \dots, x_{ij}, \dots, x_{ip}) \in \mathbb{R}^p$  le  $i$ -ème vecteur ligne de  $X$   
Il s'agit de la description de l' $i$ ème individu
- \*  $x^j = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{nj} \end{pmatrix} \in \mathbb{R}^n$  le  $j$ -ème vecteur colonne de  $X$ .  
Il s'agit de la description de la  $j$ -ème variable.

Exemple : On a mesuré la tension artérielle diastolique, systolique et le taux de cholestérol de 6 patients. Les résultats sont présentés dans le tableau ci-dessous :

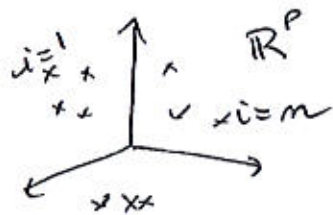
Diast	Syst	chol
90	140	6
60	85	5,9
75	135	6,1
70	145	5,8
85	130	5,4
70	145	5,0

$n =$   
 $p =$

$x_3^t =$   
 $x^2 =$

### 1) Nuage des $n$ points - individus

des  $n$  lignes de  $X$  définissent un nuage de  $n$  points de  $\mathbb{R}^p$



En général, on pondère chaque individu  $i$  par  $w_i > 0$ .

En pratique,

$$\left\{ \begin{array}{l} w_i = \frac{1}{n} \text{ (parfois } \frac{1}{n_i}) \text{ en cas de} \\ \text{tirage aléatoire} \\ w_i \neq \frac{1}{n} \text{ pour des échantillons} \\ \text{redressés, des données} \\ \text{regroupées...} \end{array} \right.$$

On notera :

$$N = \text{diag}(w_i) = \begin{pmatrix} w_1 & & 0 \\ & \ddots & \\ 0 & & w_n \end{pmatrix}$$

### 1.4) Centre de gravité du nuage des individus pondérés

On notera

$$\bar{x} = (\bar{x}_1, \dots, \bar{x}_j, \dots, \bar{x}_p) \text{ ce centre de gravité}$$

$$\bar{x}_j = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i x_{ij} \quad \text{d'où si } w_i = \frac{1}{n} \text{ ou } w_i = \frac{1}{n_i} \forall i,$$

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$$

### 1.2) Matrice centrée $Y$

Données brutes :

$$X = \begin{matrix} & \begin{matrix} 1 & \dots & j & \dots & p \end{matrix} \\ \begin{matrix} 1 \\ \vdots \\ i \\ \vdots \\ n \end{matrix} & \begin{bmatrix} \vdots & & & & \\ \vdots & & & & \\ \vdots & & x_{ij} & & \\ \vdots & & & & \\ \vdots & & & & \end{bmatrix} \end{matrix}$$

$\bar{x} = \bar{x}_1, \bar{x}_j, \dots, \bar{x}_p$

Données centrées :

$$\Rightarrow Y = \begin{matrix} & \begin{matrix} 1 & \dots & j & \dots & p \end{matrix} \\ \begin{matrix} 1 \\ \vdots \\ i \\ \vdots \\ n \end{matrix} & \begin{bmatrix} \vdots & & & & \\ \vdots & & & & \\ \vdots & & y_{ij} & & \\ \vdots & & & & \\ \vdots & & & & \end{bmatrix} \end{matrix} \quad \text{avec } y_{ij} = x_{ij} - \bar{x}_j$$

$\bar{y} = 0 \dots 0$

$\Rightarrow$  Translation du nuage de points

Exemple:  $\bar{x} =$

$y =$

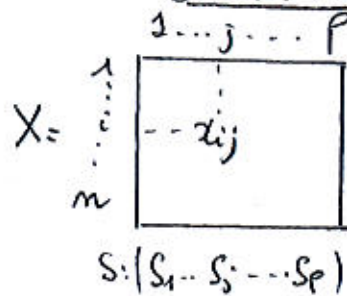


### 1.3) Variance empirique et matrice centrée-réduite Z

Variance empirique de la j-ème variable de suite par  $x^j$ :

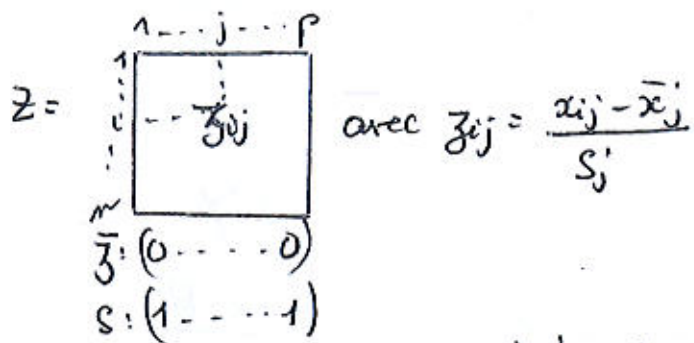
$$\begin{cases} S_j^2 = \sum_{i=1}^n w_i (x_{ij} - \bar{x}_j)^2 \\ S_j = \text{écart-type} = \sqrt{S_j^2} \end{cases} \quad \begin{array}{l} w_i = \frac{1}{n-1} : \text{estimateur sans biais} \\ w_i = \frac{1}{n} : \text{estimateur biaisé} \end{array}$$

$\Rightarrow S_j^2$  mesure la dispersion de la j-ème variable (colonne)  
Données brutes:



$\Rightarrow$

Données centrées-réduites



$\Rightarrow$  On divise chaque colonne de la matrice centrée par son écart-type

$\Rightarrow$  Variance = 1 dans toute les directions du nuage centré-réduit.

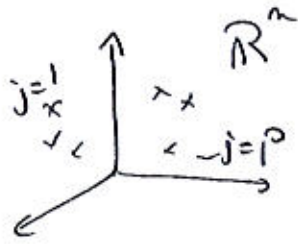
Exemple:  $S^2 =$

$S =$

$Z =$

2) Cloud de p points - nuages

Les  $p$  colonnes de  $X$  définissent un nuage de  $p$  points de  $\mathbb{R}^m$



On notera :

$$\underline{M} = \text{diag}(1/s_j^2) = \begin{pmatrix} 1/s_1^2 & & 0 \\ & \dots & \\ 0 & & 1/s_p^2 \end{pmatrix}$$

remarque:  $Z = Y M^{1/2}$

2.1) Matrice de covariance-covariance  $V$

\* Covariance empirique entre  $j$  et  $j'$ :

$$S_{jj'} = \sum_{i=1}^n w_i (x_{ij} - \bar{x}_j) (x_{ij'} - \bar{x}_{j'})$$

\* Matrice de var. cov :

$$V = (S_{jj'})_{p \times p} = Y^E N Y$$

2.2) Matrice des corrélations  $R$

\* Corrélation entre  $j$  et  $j'$ :

$$r_{jj'} = \frac{S_{jj'}}{S_j S_{j'}} = \sum_{i=1}^n w_i \overbrace{\left( \frac{x_{ij} - \bar{x}_j}{S_j} \right)}^{z_{ij}} \left( \frac{x_{ij'} - \bar{x}_{j'}}{S_{j'}} \right) = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ij'} - \bar{x}_{j'})}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} \sqrt{\sum_{i=1}^n (x_{ij'} - \bar{x}_{j'})^2}}$$

indépendant du choix de  $w_i$ !

\* Matrice de corrélations :

$$R = (r_{jj'})_{p \times p} = Z^t N Z$$

Exemple:

$V =$

$R =$

### ③ Métriques

#### 3.1) Métrique sur l'espace $\mathbb{R}^p$ des points-individus

\* Rappel: Soit  $\Gamma$  une matrice  $p \times p$ , symétrique définie positive.  
Alors  $\Gamma$  définit sur  $\mathbb{R}^p$ :

- un produit scalaire:  $\langle x, y \rangle_M = x^t M y$

- une norme:  $\|x\|_M = \langle x, x \rangle_M^{1/2}$

- une distance:  $d_M(x, y) = \|x - y\|_M$

- des angles  $\cos \theta_M(x, y) = \frac{\langle x, y \rangle_M}{\|x\|_M \|y\|_M}$

De plus, pour  $M$  donné, on peut définir:

- une matrice  $A$  est  $M$ -symétrique si  $(\Gamma A)^t = \Gamma A$

- deux vecteurs  $x$  et  $y$  sont  $\Gamma$ -orthogonaux si  $\langle x, y \rangle_M = 0$

- un vecteur  $x$  est  $\Gamma$ -normé à 1 si  $\|x\|_M = 1$ .

\* On munit l'espace  $\mathbb{R}^n$  de la métrique  $\Gamma$  pour mesurer la distance entre deux individus  $i$  et  $i'$ :

$$d_M^2(x_i, x_{i'}) = (x_i - x_{i'})^t M (x_i - x_{i'})$$

Donc,

$\rightarrow$  si  $\Gamma = I$ ,  $d_M^2(x_i, x_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \rightarrow$  distance Euclidienne simple

exemple:  $d_I^2(1, 2) =$

$\rightarrow$  si  $\Gamma = D_{1/5^2}$ ,  $d_M^2(x_i, x_{i'}) = \sum_{j=1}^p \frac{1}{5^2} (x_{ij} - x_{i'j})^2 \rightarrow$  distance Euclidienne normalisée par l'inverse de la variance:  
 $= \sum_{j=1}^p (z_{ij} - z_{i'j})^2 \rightarrow d_{D_{1/5^2}}(x_i, x_{i'}) = d_I(z_i, z_{i'})$

exemple:  $d_{D_{1/5^2}}^2(1, 2) =$

### 3.2) Métrique sur l'espace des points pondérés

En général, on mesure la proximité entre deux variables  $j$  et  $j'$  par leur covariance empirique  $S_{jj'}$ , ou par leur corrélation  $r_{jj'}$ :

On munit  $\mathbb{R}^m$  de la métrique  $N = \text{diag}(w_i)$  et on a

$$\begin{aligned} * S_{jj'} &= \sum_i w_i (\underbrace{x_{ij} - \bar{x}_j}_{y_{ij}}) (\underbrace{x_{ij'} - \bar{x}_{j'}}_{y_{ij'}}) \\ &= (y_j)^T N y_{j'} = \langle y_j, y_{j'} \rangle_N \end{aligned}$$

$\Rightarrow$  covariance entre  $x^j$  et  $x^{j'}$  est le produit scalaire (associé à  $N$ ) entre les variables centrées  $y^j$  et  $y^{j'}$ .

\*  $S_j^2 = \|y_j\|_N^2 \Rightarrow$  variance de  $x^j$  est égale à la  $N$ -norme de la variable centrée  $y^j$

$$\begin{aligned} * r_{jj'} &= \frac{S_{jj'}}{S_j S_{j'}} = \frac{\langle y_j, y_{j'} \rangle_N}{\|y_j\|_N \|y_{j'}\|_N} = \cos \theta_N(y_j, y_{j'}) \\ &= \langle z^j, z^{j'} \rangle_N \end{aligned}$$

$\Rightarrow$  corrélation entre  $x^j$  et  $x^{j'}$  est le cosinus de l'angle entre les variables centrées  $y^j$  et  $y^{j'}$ .

corrélation entre  $x^j$  et  $x^{j'}$  est le produit scalaire entre les variables centrées-oblitées  $z^j$  et  $z^{j'}$

### 4) Inertie du nuage des points-individus

L'inertie totale du nuage est la moyenne pondérée des carrés des distances des  $n$  points de  $\mathbb{R}^p$  au centre de gravité  $\bar{x}$ :

$$I(X) = \sum_{i=1}^n w_i d_M^2(x_i, \bar{x})$$

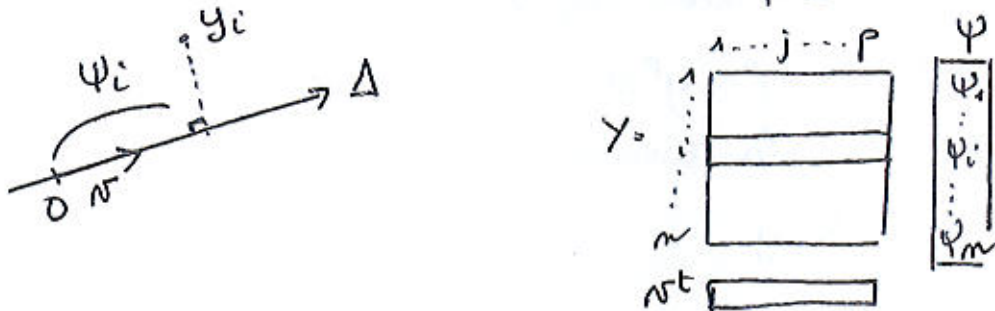
Montrer que si:

$M=I, w_i = \frac{1}{n} (\text{ou } \frac{1}{m}) \forall i$ , alors  $I(X) = S_1^2 + \dots + S_p^2$

$M=D_{1/2}$ , alors  $I(X) = P$

On considère le nuage des  $n$  points individuels centrés  $y_1, \dots, y_n \in \mathbb{R}^p$ , pondérés par  $\underline{N} = \text{diag}(\omega_i)$  et munis de la matrice  $\underline{M}$

- la projection orthogonale des  $n$  points individuels centrés  $y_1, \dots, y_n \in \mathbb{R}^p$ , sur un axe  $\Delta$  engendré par un vecteur unitaire  $\nu$  de  $M$  norme égale à 1 (i.e.  $\nu^t M \nu = 1$ ) est un vecteur  $\psi \in \mathbb{R}^n$



On a:  $\psi_i = \langle y_i, \nu \rangle = y_i^t M \nu$  ( $b = M \nu$ )

$$\Rightarrow \psi = \begin{pmatrix} \psi_1 \\ \vdots \\ \psi_i \\ \vdots \\ \psi_n \end{pmatrix} = \underbrace{Y}_{n \times p} \cdot \underbrace{M}_{p \times p} \cdot \underbrace{\nu}_{p \times 1} = \underbrace{Y}_{n \times p} \underbrace{b}_{p \times 1} = \sum_{j=1}^p b_j y^j$$

$\Rightarrow \psi$  est une combinaison linéaire des colonnes de  $Y$  avec les coefficients  $b_j$

$\Rightarrow \psi$  est un "résumé" des colonnes de  $Y$  appelée variable "synthétique"

Exemple:  $\underline{Y} = \begin{pmatrix} 0,6407 \\ 0,720 \\ -0,265 \end{pmatrix}$   $M = \underline{I}_3$   $\omega_i = \frac{1}{n}$

$\Rightarrow$  Vecteur  $\psi_1$  des projections orthogonales des 6 points individuels centrés-répondants sur l'axe  $\Delta_1$  de vecteur directeur  $\nu_1$ :

$$\psi_1 = Z \nu = 0,64 \begin{pmatrix} \cdot \\ \cdot \\ \cdot \end{pmatrix} + 0,720 \begin{pmatrix} \cdot \\ \cdot \\ \cdot \end{pmatrix} - 0,265 \begin{pmatrix} \cdot \\ \cdot \\ \cdot \end{pmatrix} =$$

$$\|\nu_1\|^2 = 1$$

$\vec{v}_2 = \begin{pmatrix} 0,445 \\ 1 \\ -0,065 \\ 0,894 \end{pmatrix}$  Les paires de  $\vec{v}_2$  est I-orthogonal à  $\vec{v}_1$  :  $\vec{v}_1^t \vec{v}_2 = 0$

Calculer  $\psi_2^T$  le vecteur des projections des 6 points-individus centrés-réduits sur l'axe  $\Delta_2$  de vecteur directeur  $\vec{v}_2$ .

$$\psi_2^T = \begin{pmatrix} \\ \\ \\ \\ \\ \end{pmatrix}$$

• Variance de la variable synthétique  $\psi$  :

→ Comme les colonnes  $y^1 \dots y^p$  de  $Z$  et les colonnes  $z^1 \dots z^p$  de  $Z$  sont centrés, la combinaison linéaire,

$$\psi = \sum_{j=1}^p b_j y^j \quad \text{Exemple: } \bar{\psi}_1 = \bar{\psi}_2 = 0$$

est centré :  $\bar{\psi} = 0$

TP : Reprendre l'exemple du cours dans R (en.vers aidant du code se trouvent dans le fichier exemple-notions-bases.R) et dans Excel !