

Imputation des données manquantes

1. Imputation simple : remplacer la valeur manquante par une seule valeur "plausible"
- imputation par moyenne (quantitative)
ou mode (qualitative)
 - imputation par régression = remplacer la valeur manquante par une valeur prédite par un modèle de régression. (simple, en général)
 - imputation à l'aide des analyses multivariées (ACP, ACM)
idée = utiliser la formule de substitution de données à l'aide de l'ACP.

Dans la suite on va s'intéresser à cette dernière méthode qui sera appelée NIPALS (Algorithme NIPALS pour données manquantes).

NIPALS = Nonlinear Iterative Partial Least Squares.

L'algorithmme NIPALS

13

Observation importante : dans le cas des données sans valeur manquantes
NIPALS = formule de reconstitution des données en ACP.

NIPALS avec données complètes :

$$\text{Soit } X = \begin{pmatrix} x_{11} & x_{1j} & x_{1p} \\ x_{21} & x_{2j} & x_{2p} \\ \vdots & \vdots & \vdots \\ x_{i1} & x_{ij} & x_{ip} \\ \vdots & \vdots & \vdots \\ x_{n1} & x_{nj} & x_{np} \end{pmatrix}$$

x_1, x_2, \dots, x_p = variables quantitatives

- il est important de bien savoir
 - la régression linéaire
 - l'analyse en composantes principales

On rappelle ces deux aspects dans les pages suivantes.

Rappel 1 : Régression linéaire simple

(14)

On dispose de deux variables X et Y observées sur n individus :

	X	Y
1	x_1	y_1
2	x_2	y_2
\vdots	\vdots	\vdots
i	x_i	y_i
\vdots	\vdots	\vdots
n	x_n	y_n

L'équation de régression linéaire est

$$Y = \beta_0 + \beta_1 X + \varepsilon \text{ et sous les}$$

hypothèses classiques de la régression on obtient les estimateurs :

$$\left\{ \begin{array}{l} \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \\ \hat{\beta}_1 = \frac{\frac{1}{n} \sum x_i y_i - \bar{X} \bar{Y}}{S_x^2} \end{array} \right.$$

$$\text{avec } \bar{X} = \frac{1}{n} \sum x_i, \quad \bar{Y} = \frac{1}{n} \sum y_i$$

$$S_x^2 = \frac{1}{n} \sum (x_i - \bar{X})^2$$

Observons que ni les variables sont centrées et réduites alors :
($\bar{x}=0, \bar{y}=0$) ($s_x^2=1, s_y^2=1$)

$$\hat{\beta}_0 = 0 \quad \text{et} \quad \hat{\beta}_1 = \frac{1}{n} \sum x_i y_i$$

et le modèle s'écrit donc :

$$Y = \left(\frac{1}{n} \sum x_i y_i \right) \cdot X + \Sigma, \text{ ou, en forme matricielle}$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} = \left(\frac{1}{n} \sum x_i y_i \right) \cdot \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} + \begin{pmatrix} \Sigma_1 \\ \Sigma_2 \\ \vdots \\ \Sigma_m \end{pmatrix}$$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix} = \hat{b}_1 \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{pmatrix} + \begin{pmatrix} \Sigma_1 \\ \vdots \\ \Sigma_i \\ \vdots \\ \Sigma_n \end{pmatrix}$$

Notation : $\sum x_i y_i = \langle X, Y \rangle$ ← produit scalaire

Remarque importante : régression linéaire (16)
multipli avec variables indépendantes
 (ou non-corrélées).

On a :

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix} \quad \text{et} \quad X = \begin{matrix} & \begin{matrix} X_1 & X_2 & \dots & X_p \end{matrix} \\ \begin{matrix} \downarrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \end{matrix} & \begin{bmatrix} X_{11} & \dots & X_{1j} & \dots & X_{1p} \\ X_{21} & & X_{2j} & & X_{2p} \\ \vdots & & \vdots & & \vdots \\ X_{i1} & \dots & X_{ij} & \dots & X_{ip} \\ \vdots & & \vdots & & \vdots \\ X_{n1} & \dots & X_{nj} & \dots & X_{np} \end{bmatrix} \end{matrix}$$

On suppose que :

- X_1, X_2, \dots, X_p sont indépendantes et
 de plus centrées et réduites.

Alors, dans ce cas, la matrice de var-cov. de X

$$V = \text{Cov}(X) = R(X) = \frac{1}{n} X^T X = I_p = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}_{p \times p}$$

Dans le modèle linéaire

$$y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \Sigma$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} = (X^T X)^{-1} X^T y$$

ou encore :

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_j \\ \vdots \\ \hat{\beta}_p \end{pmatrix} = \begin{pmatrix} \frac{1}{n} & & & & \\ & \frac{1}{n} & & & \\ & & \ddots & & \\ & & & \frac{1}{n} & \\ & & & & 0 \end{pmatrix} \cdot \begin{pmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \quad (17)$$

$I_n : (n \times n)$ $X^T : (p \times n)$

On obtient que

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_j \\ \vdots \\ \hat{\beta}_p \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \sum x_{i1} \cdot y_i = \frac{1}{n} \langle x_1, Y \rangle \\ \frac{1}{n} \sum x_{i2} \cdot y_i = \frac{1}{n} \langle x_2, Y \rangle \\ \vdots \\ \frac{1}{n} \sum x_{ip} \cdot y_i = \frac{1}{n} \langle x_p, Y \rangle \end{pmatrix}$$

La leçon à tirer est que les coefficients des variables x_j , les $\hat{\beta}_j$, sont en fait les mêmes comme si on avait fait juste de régression linéaire simple entre Y et x_j .

Pour résumer en "deux" mots : la régression linéaire multiple avec p variables indépendantes se réduit à p régressions simples ($Y \sim x_j ; j = 1 \dots p$)

Fin Rappel 1.

Rappel 2 : Analyse en composantes principales (ACP) (18)

Construire des composantes qui resument "au mieux" l'information dans les données.

composante (linéaire) :

$$c = u_1 X_1 + u_2 X_2 + \dots + u_p X_p \quad \text{avec}$$

$u = \begin{pmatrix} u_1 \\ \vdots \\ u_p \end{pmatrix}$ des poids définissant la composante c .

les poids sont normalisés

$$\sum_{i=1}^p u_i^2 = 1 \quad \text{car le critère}$$

pour déterminer c est celui de variance (information) maximale.

Donc, ACP consiste à chercher c t.q.

$\text{Var}(c)$ soit maximale

parmi toutes les combinaisons linéaires possibles de X_1, \dots, X_p .

Solution : $\nabla u = \lambda u$, u associé à la valeur propre la plus grande de X .

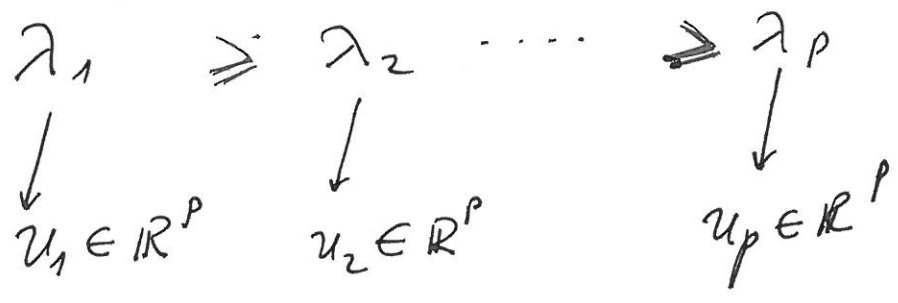
↓
matrice de covariance de X .

La matrice de variance-covariance V
 (qui coïncide avec la matrice de corrélation R quand
 les données sont centrées et réduites) a les
 propriétés suivantes :

- elle est symétrique et positive définie
- elle admet p valeurs propres positives ou nulles

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

- à chaque valeur propre on associe un vecteur propre u_i :



tels que :

$$\|u_i\|^2 = \sum_{j=1}^p u_{ij}^2 = 1, \quad \forall i = 1 \dots p$$

(a) $u_i = \begin{pmatrix} u_{i1} \\ u_{i2} \\ \vdots \\ u_{ip} \end{pmatrix}$ = le i ème vecteur propre de V .

et

(b) u_i et u_j sont orthogonaux entre eux pour $i \neq j$
 $\langle u_i, u_j \rangle = \sum_{k=1}^p u_{ik} \cdot u_{jk} = \underline{\underline{0}}$

Le point b) nous indique que



$\{u_1, u_2, \dots, u_p\}$ forment une base

dans \mathbb{R}^p .

Du coup, tout élément x de \mathbb{R}^p peut s'écrire dans la base $\{u_1, u_2, \dots, u_p\}$.

$$x = \alpha_1 \cdot u_1 + \alpha_2 u_2 + \dots + \alpha_p u_p$$

$$\begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$$

avec $\alpha_i, i=1 \dots p$, les coefficients de développement.

$$\alpha_i = \langle u_i, x \rangle = \sum_{k=1}^p u_{ik} \cdot x_k$$

→ De manière vectorielle :

$$\begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} = \alpha_1 \begin{pmatrix} u_{11} \\ \vdots \\ u_{1p} \end{pmatrix} + \dots + \alpha_p \begin{pmatrix} u_{p1} \\ \vdots \\ u_{pp} \end{pmatrix}$$

écriture en colonne.

ou

$$(x_1, \dots, x_p) = \alpha_1 (u_{11}, \dots, u_{1p}) + \dots + \alpha_p (u_{p1}, \dots, u_{pp})$$

écriture en ligne

On définit donc les composantes principales 21
de X , comme étant les variables définies
par les vecteurs $\{u_1, u_2, \dots, u_p\}$.

On a donc p composantes principales :

$$u_1 \rightarrow C_1 = u_{11}X_1 + u_{12}X_2 + \dots + u_{1p}X_p$$

$$u_2 \rightarrow C_2 = u_{21}X_1 + u_{22}X_2 + \dots + u_{2p}X_p$$

$$\vdots$$

$$u_p \rightarrow C_p = u_{p1}X_1 + u_{p2}X_2 + \dots + u_{pp}X_p$$

ou matriciellement :

$$C_1 = \begin{pmatrix} c_{1,1} \\ c_{2,1} \\ \vdots \\ c_{n,1} \end{pmatrix} = u_{11} \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{pmatrix} + u_{12} \begin{pmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{n2} \end{pmatrix} + \dots + u_{1p} \begin{pmatrix} x_{1p} \\ x_{2p} \\ \vdots \\ x_{np} \end{pmatrix}$$

\uparrow \uparrow \uparrow
 X_1 X_2 X_p

Vous remplacez C_1 par C_i et vous obtenez
l'expression de toute composante C_i .

$$\text{Notez que } C_1 = \begin{pmatrix} c_{11} \\ c_{21} \\ \vdots \\ c_{n1} \end{pmatrix} = \begin{pmatrix} \langle u_1, (x_{11}, x_{12}, \dots, x_{1p}) \rangle \\ \vdots \\ \langle u_1, (x_{n1}, x_{n2}, \dots, x_{np}) \rangle \end{pmatrix} = \begin{pmatrix} \langle u_1, \text{ind}_1 \rangle \\ \vdots \\ \langle u_1, \text{ind}_n \rangle \end{pmatrix}$$

Ces composantes principales sont donc des nouvelles variables aléatoires construites à partir des variables X_1, \dots, X_p et qui ont les propriétés suivantes :

	X_1	X_2	...	X_p	C_1	C_2	...	C_p
1	X_{11}	X_{12}	...	X_{1p}	c_{11}	c_{12}	...	c_{1p}
2	X_{21}	X_{22}	...	X_{2p}	c_{21}	c_{22}		c_{2p}
...	...							
n	X_{n1}	X_{n2}		X_{np}	c_{n1}	c_{n2}		c_{np}

Données d'origine Composantes principales

Propriétés :

- [données centrées et réduites] : $\bar{c}_i = 0$ et $V(c_i) = \lambda_i$
(moyenne nulle et variance = val propre).

⚠ Attention : discussion sur les val. propres nulles $\lambda_i = 0$.

- $\text{cor}(c_i, c_j) = 0 \quad \forall i, j \quad i \neq j$
variables non corrélées

Comment les composantes principales resument-elles l'information contenue dans X ?

Information = variance

Alors, l'information totale contenue dans le tableau de données X est donnée par l'inertie totale, qui après petit calcul, donne

$$I_{total} = \underline{V(x_1) + \dots + V(x_p)} \quad (= p \text{ si } x_i \text{ reduites})$$

A leur tour, l'information contenue dans les composantes C_1 est $Var(C_1) = \lambda_1$.

On dit que la composante C_1 explique $\frac{\lambda_1}{I_{tot}}$ de l'information.

Pareil pour une composante quelconque, c_i , à elle seule, elle explique $\frac{\lambda_i}{I_{total}}$

de l'information,

Maintenant, si on met ensemble plusieurs composantes, on a le pouvoir explicatif suivant :

• C_1 toute seule : $\frac{\lambda_1}{I_{tot}} = \frac{\lambda_1}{P}$
si ACP normée

Donc C_1 toute seule resume $\frac{\lambda_1}{I_{tot}}$ de l'information de X .

s_{ij}	X_1	...	X_p
1	x_{11}	...	x_{1p}
2	x_{21}	...	x_{2p}
...	x_{i1}	...	x_{ip}
h	x_{h1}	...	x_{hp}

 \approx

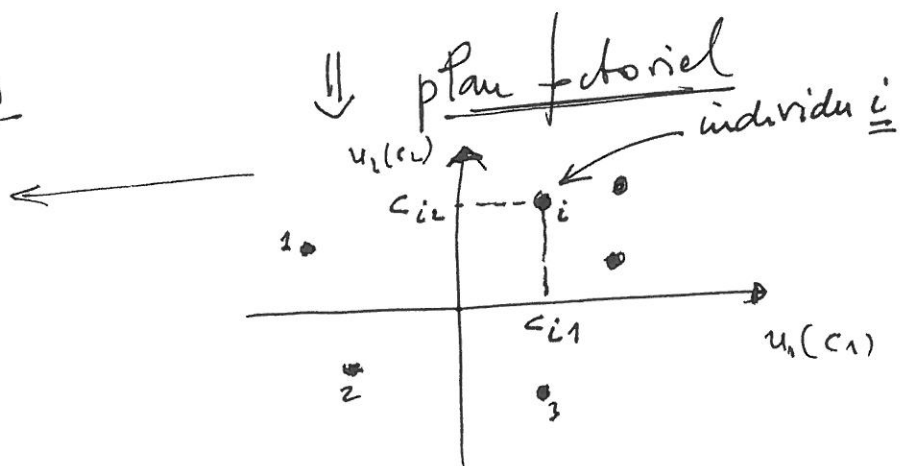
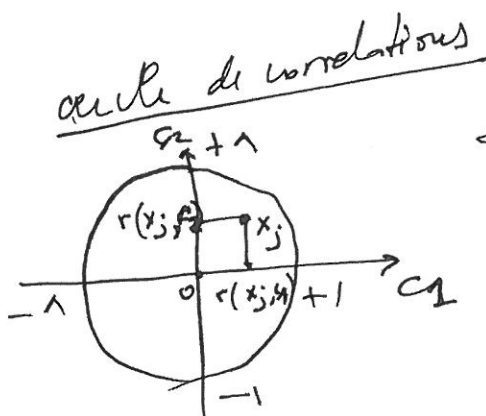
s_{ij}	C_1
1	c_{11}
2	c_{21}
...	c_{i1}
h	c_{h1}

• C_1 et C_2 : $\frac{\lambda_1 + \lambda_2}{I_{total}} = \frac{\lambda_1 + \lambda_2}{P}$
si ACP normalisée

	X_1	...	X_p
1	x_{11}	...	x_{1p}
2	x_{21}	...	x_{2p}
...	x_{i1}	...	x_{ip}
h	x_{h1}	...	x_{hp}

 \approx

	C_1	C_2
1	c_{11}	c_{12}
2	c_{21}	c_{22}
...	c_{i1}	c_{i2}
h	c_{h1}	c_{h2}



Et ainsi de suite.

On garde en général un nombre de composantes qui assurent un bon taux d'information exigée (90%, par exemple, mais il n'y a pas de règle précise).

- Supposons que $q \leq p$ les valeurs propres de V sont non-nulles : $\lambda_i > 0, \forall i=1 \dots q$.
 et que $\lambda_{q+j} = 0 \forall j \in [1, p-q]$.
 $\lambda_1 \geq \lambda_2 \dots \geq \lambda_q > 0 = \lambda_{q+1} = \dots = \lambda_p$.

Ce qui est un peu moins connu de l'ACP est la formule de reconstitution des données.

Le fait que $\{u_1, u_2, \dots, u_p\}$ forment une base de \mathbb{R}^p cela veut dire que l'individu i

$$ind_i = (x_{i1}, x_{i2}, \dots, x_{ip})$$

peut être exprimé dans cette base, et donc

$$ind_i = \langle ind_i, u_1 \rangle \cdot u_1 + \langle ind_i, u_2 \rangle \cdot u_2 + \dots + \langle ind_i, u_p \rangle \cdot u_p$$

$$ind_i = \underbrace{c_{i1}}_{\in \mathbb{R}} \cdot \underbrace{u_1}_{\in \mathbb{R}^p} + c_{i2} \cdot u_2 + \dots + c_{ip} \cdot u_p$$

ou encore :

$$(x_{i1}, \dots, x_{ip}) = c_{i1}(u_{11}, \dots, u_{1p}) + c_{i2}(u_{21}, \dots, u_{2p}) + \dots + c_{ip}(u_{p1}, \dots, u_{pp})$$

Remarque : évidemment, si $\lambda_j = 0$ alors $c_{ij} = 0$
 car $V(c_j) = 0$ et $c_j \equiv 0 \triangleq$
 pour tous les individus..

Donc, chaque ligne de la matrice X peut s'écrire comme une combinaison linéaire des vecteurs propres de V .

(On savait que chaque composante c_i est une combinaison des variables $x_j, j=1 \dots p$.)

En mettant la décomposition de (x_{i1}, \dots, x_{ip}) pour tous les i sous forme matricielle, on obtient :

$$\begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1p} \\ \vdots & & & & \\ x_{i1} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & & & & \\ x_{n1} & \dots & x_{nj} & \dots & x_{np} \end{pmatrix} = \begin{pmatrix} c_{11} \\ c_{21} \\ \vdots \\ c_{i1} \\ \vdots \\ c_{n1} \end{pmatrix} (u_{11}, \dots, u_{1p}) + \begin{pmatrix} c_{12} \\ c_{22} \\ \vdots \\ c_{i2} \\ \vdots \\ c_{n2} \end{pmatrix} (u_{21}, u_{22}, \dots, u_{2p}) + \dots$$

$$X = c_1 \cdot u_1^T + c_2 \cdot u_2^T + \dots + c_q \cdot u_q^T$$

⇒ Formule de reconstitution de l'AEP.

On obtient donc la formule de reconstitution (27)
 des données à l'aide de l'ACP:

Elle nous dit que les données d'origine peuvent être reconstituées à partir des composantes principales $(c_i, i=1 \dots q)$ et des facteurs principaux $(u_i, i=1 \dots p)$.

$$\begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1p} \\ \vdots & & \vdots & & \vdots \\ x_{i1} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & & \vdots & & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{np} \end{pmatrix} = \begin{pmatrix} c_{11} \\ \vdots \\ c_{i1} \\ \vdots \\ c_{n1} \end{pmatrix} \cdot (u_{11}, \dots, u_{1p}) + \dots + \begin{pmatrix} c_{1q} \\ \vdots \\ c_{iq} \\ \vdots \\ c_{nq} \end{pmatrix} \cdot (u_{q1}, \dots, u_{qp})$$

⚠ relation exacte

Si l'on pense que les dernières composantes principales (peu explicatives) sont en fait liées à un "bruit" présent dans les données d'origine, alors on peut retenir dans cette décomposition uniquement $r < q$ composantes pour "débriiter" les données. On aurait alors :

$$\underbrace{X}_{\text{données observées}} = \underbrace{c_1 \cdot u_1^T + \dots + c_r \cdot u_r^T}_{\text{données débriitées}} + \underbrace{c_{r+1} \cdot u_{r+1}^T + \dots + c_q \cdot u_q^T}_{\text{bruit}}$$

Fin rappel 2.

Maintenant on est en mesure de présenter l'algorithme NIPALS.

NIPALS sur un tableau de données sans valeurs manquantes.

NIPALS \Leftrightarrow ACP.

On s'intéresse dans un premier temps au calcul de la 1^{ère} composante principale, x_1 .

x_1 est donnée par ses valeurs prises sur chaque individu:

$$x_1 = \begin{pmatrix} c_{11} \\ \vdots \\ c_{i1} \\ \vdots \\ c_{n1} \end{pmatrix}$$

On rappelle la formule de régression (page 26, top)

$$(x_{i1}, \dots, x_{ip}) = \underbrace{c_{i1}}_{u_1} (u_{11}, \dots, u_{1p}) + c_{i2} (u_{21}, \dots, u_{2p}) + \dots + c_{ip} (u_{p1}, \dots, u_{pp})$$

Remarque importante:

On constate ici que c_{i1} , $i = 1, \dots, n$ est le coefficient de la régression linéaire simple du vecteur (x_{i1}, \dots, x_{ip}) sur le vecteur (u_{11}, \dots, u_{1p})

Donc, si on connaît le vecteur $u_1 = (u_{11}, \dots, u_{1p})$ on connaîtrait aussi le vecteur $c_1 = (c_{11}, \dots, c_{n1})$ en réalisant n régressions simples.

Comment trouver le vecteur $u_1 = (u_{11} \dots u_{1p})$

Toujours grâce à la formule de reconstruction, on constate que la variable X_j (comme colonne) s'écrit :

$$\begin{pmatrix} X_j \\ X_{1j} \\ \vdots \\ X_{ij} \\ \vdots \\ X_{nj} \end{pmatrix} = \begin{pmatrix} c_{11} \\ \vdots \\ c_{i1} \\ \vdots \\ c_{n1} \end{pmatrix} \cdot u_{1j} + \begin{pmatrix} c_{12} \\ \vdots \\ c_{i2} \\ \vdots \\ c_{n2} \end{pmatrix} \cdot u_{2j} + \dots + \begin{pmatrix} c_{1q} \\ \vdots \\ c_{iq} \\ \vdots \\ c_{nq} \end{pmatrix} \cdot u_{qj}$$

De nouveau, grâce à l'orthogonalité des composantes c_i , on constate que

Remarque

u_{1j} est le coefficient de la régression linéaire simple de la variable X_j sur la variable X_1 .

Donc pour ~~estimer~~ trouver $u_2 = (u_{21} \dots u_{2p})$ on doit faire p régressions linéaires simples entre $(X_j \text{ et } X_1)$, $j=1 \dots p$.

On obtient donc l'algorithme affine suivant :

1. Initialisation :

$c_1 = X_1$ (la premiere composante est egale à la variable X_1)
↳ on pourrait prendre une autre si l'on veut.

2. Calcul de u_1 : effectuer donc p regressions linéaires multiples.

for (j in 1:p)

{ $u_{1j} =$ coeff de la regression de X_j sur c_1

(*) $u_{1j} = \langle X_j, c_1 \rangle = \sum_{k=1}^n X_{kj} \cdot c_{k1}$

{ $u_{1j} = \text{coef}(\text{lm}(X_j \sim c_1))[2]$

on normalise u_1 :

$$u_{1j} = \frac{u_{1j}}{\sqrt{\sum u_{1j}^2}}$$

3. Calcul de c_1 : on effectue n regressions linéaires entre les individus et u_1

for (i in 1:n)

(**) { $c_{1i} = \sum_{k=1}^p X_{ik} \cdot u_{1k}$ # regression de l'individu i sur u_1

{ $c_{1i} = \text{coef}(\text{lm}(X[i,] \sim u_1))[2]$

④ Itérer les pas 2 et 3 jusqu'à la convergence.

Soit on se donne un nombre fixé d'itérations
Soit on regarde d'un pas à l'autre l'évolution de

$$\|u_1^{(s)} - u_1^{(s+1)}\|$$

et on s'arrête lorsque cette quantité est inférieure à un seuil fixé ($\epsilon = 10^{-5}$ par ex.).

Fin algorithme NIPALS pour le calcul de la composante 1


Comment on calcule la composante \underline{x}_2 après avoir calculé la composante \underline{x}_1 ?

On sait que

$$X = c_1 \cdot u_1^T + c_2 \cdot u_2^T + \dots + c_9 \cdot u_9^T$$

$$X - c_1 \cdot u_1^T = c_2 \cdot u_2^T + \dots + c_9 \cdot u_9^T$$

$$X_{\Delta} = c_2 \cdot u_2^T + \dots + c_9 \cdot u_9^T$$

et on applique NIPALS pour la 1^{ère} composante à la matrice X_{Δ} . 

Après avoir obtenu

$$(c_1, u_1), (c_2, u_2), \dots, (c_h, u_h)$$

pour obtenir c_{h+1} et u_{h+1} on applique l'algorithme NIPALS pour chercher la 1^{ère} composante au tableau

$$X_h = X - c_1 \cdot u_1^T - c_2 \cdot u_2^T \dots - c_h \cdot u_h^T$$

FIN NIPALS données complètes

Et en présence de données manquantes ?

La force de NIPALS est qu'il est basé uniquement sur des régressions linéaires simples. Et ces régressions sont faisables en présence de données manquantes (approximation de coefficients sur les données (couples complètes !)).

C'est à dire que si x et y sont deux variables avec des données manquantes (33)

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \text{ alors, en travaillant}$$

sur les couples (x_i, y_i) où les données sont présentes,

$\hat{\beta}_1$ peut être estimé par

$$\hat{\beta}_1 = \sum_{i=1}^n x_i \cdot y_i \quad \text{à une constante près}$$

$\left(\frac{1}{\text{nb couples non manquants}} \right)$

$\{i | (x_i, y_i) \text{ non manquants}\}$

Du coup, dans l'algorithme NIPALS (page 30) les pas (*) et (**) peuvent être modifiés par

$$(*) \quad u_{1j} = \sum_{k=1}^n x_{kj} \cdot c_{k1} \implies u_{1j} = \sum_{k=1}^n x_{kj} \cdot c_{k1}$$

$\{k: c_{k1} \text{ existe}\}$

$$(**) \quad c_{1i} = \sum_{k=1}^p x_{ik} \cdot u_{1k} \implies c_{1i} = \sum_{k=1}^p x_{ik} \cdot u_{1k}$$

$\{k: x_{ik} \text{ existe}\}$

A la fin de l'algorithme NIPALS appliqué pour le calcul des r composantes principales et r facteurs principaux on aura :

$$\hat{c}_1 = \begin{pmatrix} \hat{c}_{11} \\ \vdots \\ \hat{c}_{nr} \end{pmatrix}, \hat{u}_1 = (\hat{u}_{11}, \dots, \hat{u}_{1p})$$

$$\hat{c}_2 = \begin{pmatrix} \hat{c}_{21} \\ \vdots \\ \hat{c}_{nr} \end{pmatrix}, \hat{u}_2 = (\hat{u}_{21}, \dots, \hat{u}_{2p})$$

$$\hat{c}_r = \begin{pmatrix} \vdots \\ \hat{c}_{r1} \\ \vdots \\ \hat{c}_{nr} \end{pmatrix}, \hat{u}_r = (u_{r1}, \dots, u_{rp})$$

avec données complètes! (sans "trous").

Avec la formule de reconstruction des données on obtient

$$X \approx \hat{c}_1 \cdot \hat{u}_1^T + \hat{c}_2 \cdot \hat{u}_2^T + \dots + \hat{c}_r \cdot \hat{u}_r^T$$

↑
avec d.m.

pas des d.m.

On obtient donc une estimation des données manquantes :

Si x_{ij} est l'observation de la variable j sur l'individu i alors son estimation à l'aide des r composantes principales est donnée par (voir page 29)

$$\hat{x}_{ij} = \hat{c}_{i1} \cdot \hat{u}_{1j} + \hat{c}_{i2} \cdot \hat{u}_{2j} + \dots + \hat{c}_{ir} \cdot \hat{u}_{rj}$$

Si x_{ij} est manquante, alors il s'agit d'une imputation.

Si x_{ij} est présente alors \hat{x}_{ij} est une approximation de x_{ij} si $r < q$. (valeur exacte si $r = q$).

Remarque importante :

Nous avons considéré au départ que les données du tableau X sont centrées et réduites. Avant d'appliquer NIPALS il faut donc centrer et réduire les variables (avec les valeurs manquantes!).

$$X \xrightarrow[\substack{\text{centrer} \\ + \\ \text{réduire}}]{\downarrow \begin{matrix} \{\bar{x}_1, \dots, \bar{x}_p\} \\ \{s_1^2, \dots, s_p^2\} \end{matrix}} Z \xrightarrow{\text{NIPALS}} \hat{Z}$$

$$\hat{x}_{ij} = \bar{x}_j + s_j \cdot \hat{z}_{ij}$$

Fiche TP (R)

1 simulation :

1.1 générer un tableau X de $\frac{n=100}{\text{individus}} \times \frac{p=4}{\text{variables}}$

provenant d'une distribution normale multivariée

$N(\mu, \Sigma)$ avec

$$\mu = (1, 2, 4, 3) \text{ et}$$

$$\Sigma = \begin{pmatrix} 0.7 & 0 & 1.3 & 0.5 \\ 0 & 0.2 & -0.3 & -0.1 \\ 1.3 & -0.3 & 3.1 & 1.3 \\ 0.5 & -0.1 & 1.3 & 0.6 \end{pmatrix}$$

Nb : utiliser le package "mvtnorm"

1.2 Réaliser l'ACP normée de X

- utiliser le package "FactoMineR"

- regarder les composantes principales et les facteurs propres.

1.3 programmer NIPALS. Comparer les résultats fournis par NIPALS avec ceux donnés par FactoMineR.

1.4. générer aléatoirement des valeurs manquantes dans le tableau X . On va supposer que $P(X_{ij} \text{ soit manquante}) = 0.05$ (Données MCAR donc!).

①.5 On considère maintenant

\tilde{X} le tableau avec données manquantes obtenu en 1.4.

Appliquez NIPALS pour estimer les données manquantes. Comparez avec les vraies valeurs

②. Application réelle.

Imputez par NIPALS les données manquantes présentes dans la base "airquality" (base de données existante en R)

> help(airquality)

— enjoy! —

Fin imputation simple. Passons à l'imputation multiple.