

Analyse en Composantes Principales (ACP)

①

Il s'agit d'une méthode d'analyse et de description multidimensionnelle d'un tableau de données quantitatives.

L'ACP permet (entre autre) des représentations graphiques des similitudes entre les lignes (des distances euclidiennes entre les individus) et des liaisons entre les colonnes (les correlations entre les variables)

Il s'agit également d'une méthode de rédaction de dimension (construction d'un "petit" nombre de variables synthétiques "rencontrant" au mieux toutes les variables initiales).

1) Notations

On considère (n) individus $\{1, \dots, i, \dots, n\}$ décrits sur (p) variables $\{1, \dots, j, \dots, p\}$. On notera w_i , $i=1, \dots, n$ le pieds d'un individu i (généralement $\frac{1}{n}$).

1.1) Trois matrices de données

Données brutes: $X_{m \times p}$ Données centrées: $Y_{m \times p}$ Données centrées-rebâties:

$$X = \begin{pmatrix} 1 & \dots & j & \dots & p \\ \vdots & & \vdots & & \vdots \\ 1 & \dots & x_{ij} & \dots & \vdots \\ \vdots & & \vdots & & \vdots \\ m & \dots & \vdots & & \vdots \end{pmatrix}$$

$$Y = \begin{pmatrix} 1 & \dots & j & \dots & p \\ \vdots & & \vdots & & \vdots \\ 1 & \dots & y_{ij} & \dots & \vdots \\ \vdots & & \vdots & & \vdots \\ m & \dots & \vdots & & \vdots \end{pmatrix}$$

$$y_{ij} = x_{ij} - \bar{x}_j$$

$$Z = \begin{pmatrix} 1 & \dots & j & \dots & p \\ \vdots & & \vdots & & \vdots \\ 1 & \dots & z_{ij} & \dots & \vdots \\ \vdots & & \vdots & & \vdots \\ m & \dots & \vdots & & \vdots \end{pmatrix}$$

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j} = \frac{y_{ij}}{s_j}$$

\Rightarrow 3 matrices différentes pour décrire

les mêmes individus \Rightarrow 3 images de points différents

(2)

Dans chaque matrice :

- un individu i est décrit par un vecteur de \mathbb{R}^p noté x_i, y_i, z_i
- une variable j est décrite par un vecteur de \mathbb{R}^n noté x^j, y^j, z^j

On parle toujours de vecteurs colonnes (même si ils décrivent les lignes des tableaux X, Z ou Y).

1.2) Deux métriques

On muni l'espace $\begin{cases} \mathbb{R}^p \text{ d'une métrique } M, \text{ de dimension } p \times p \\ \mathbb{R}^n \text{ d'une métrique } N, \text{ de dimension } n \times n \end{cases}$

L'ACP va consister à analyser les n points-individus (les lignes) et les p points-relevés (les colonnes) de la matrice des données entrées-redressées, Z , avec les métriques :

$$M = \underbrace{\begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix}}_p = I_p \text{ sur } \mathbb{R}^p$$

$$N = \begin{pmatrix} w_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & w_n \end{pmatrix} \text{ sur } \mathbb{R}^n \text{ avec en général } w_i = \frac{1}{n} \quad \forall i=1 \dots n,$$

donc $N = \underbrace{\frac{1}{n} \begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix}}_n = \frac{1}{n} I_n$.

On dit qu'on fait l'analyse du triplet $(Z, I_p, N) \Leftrightarrow \begin{cases} \text{ACP normée} \\ \text{ACP sur matrice des corrélations} \end{cases}$

Remarque: Faire l'analyse du triplet (Y, M, N) avec $M = I_p$

$$\Leftrightarrow \begin{cases} \text{ACP non normée} \\ \text{ACP sur matrice des covariances} \end{cases}$$

Exemple

	Intensité bulles	Savon Salé	Appréciation globale
S ⁺ Yorre	3.9	6.4	2.9
Vichy	1.4	6.0	2.8
Quercy	5.1	6.7	3.5
Selvaticat	2.9	6.1	3.4
Perrier	8.2	6.9	2.8
moyenne de type	4.28 2.29	5.22 0.35	3.08 0.3

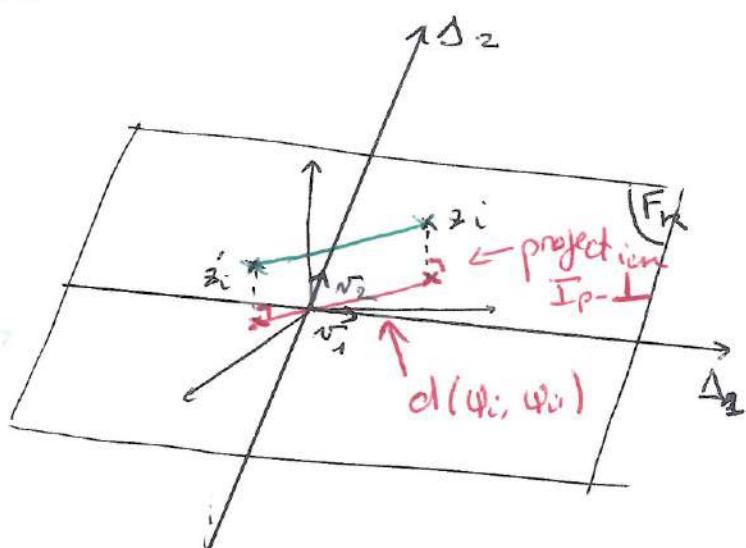
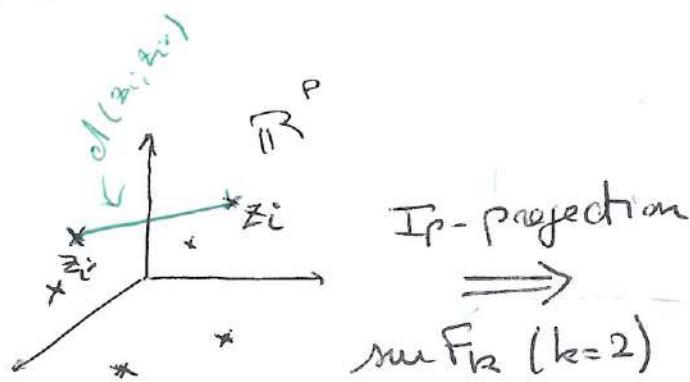
$$\Rightarrow Z = \begin{pmatrix} -0.17 & 1.38 & -0.58 \\ -1.26 & 0.91 & -0.91 \\ 0.35 & -0.61 & 1.37 \\ -0.6 & -1.31 & 1.04 \\ 1.69 & -0.37 & -0.91 \end{pmatrix} \quad N = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \frac{1}{5} I_5$$

$$M = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = I_3$$

2) Analyse du nuage des n points-individus de \mathbb{R}^P

2.1) Le principe

On va chercher un sous-espace vectoriel F_k de \mathbb{R}^P , de dimension $k \leq r$ (avec $r = \text{rang}(Z)$), c'est à dire
 k axes $\Delta_1, \dots, \Delta_k$ (I_P -orthogonaux) tels que
 le nuage des n points individus I_P -projétés
 sur F_k , "d'forme le moins possible" les distances
 euclidiennes entre les individus.



Nuage de n points-individus
centré-reduits

Matrice centrée-réduite

$$Z = \begin{pmatrix} 1 & \dots & j & \dots & P \\ \vdots & & \vdots & & \vdots \\ n & & z_i^t & \in \mathbb{R}^P \end{pmatrix} \quad \text{projection} \quad \Rightarrow \quad \text{sur } F_k$$

Matrice des coordonnées factorielles des individus ④

$$\Psi = \begin{pmatrix} 1 & \dots & d & \dots & k \\ \vdots & & \vdots & & \vdots \\ n & & \Psi_{ik} & \leftarrow \Psi_i^t \\ & & \vdots & & \Psi_i^t \end{pmatrix}$$

On appelle $v_1, \dots, v_d, \dots, v_k$ les vecteurs directeurs de $\delta_1, \dots, \delta_k$ les axes principaux

On appelle $\Psi_{n \times k}$ la matrice des coordonnées des n individus projetés sur ces axes :

$\left\{ \begin{array}{l} \text{matrice des coordonnées factorielles des individus} \\ \text{matrice des scores des individus sur les} \\ \text{1 premières composantes principales} \end{array} \right.$

On appelle $\Psi^k \in \mathbb{R}^n$ la "d-ième composante principale" c'est à dire le vecteur des scores des n individus sur cette composante

On a $\Psi_i^k \in \mathbb{R}^k$ le vecteur des coordonnées factorielles de l'individu i sur les k premiers axes = vecteur des scores de l'individu i sur les k premières composantes principales.

Que veut dire "déformer le moins" et quel critère va-t-on optimiser?

On a la relation suivante :

$$\sum_i \sum_{i'} w_i w_{i'} d^2(z_i, z_{i'}) = 2 \underbrace{\sum_i w_i d^2(z_i, \bar{z})}_{\text{Inertie} = I(z) = P \text{ (déjà démo)}}$$

Somme pondérée des carrés des distances entre paires d'individus

Somme des carrés des distances entre les individus et le centre de gravité = mesure de dispersion des nuages.

On a toujours que la distance entre i et i' dans \mathbb{R}^p (en vert) est plus grande que la distance "en projection" c'est à dire entre i et i' dans \mathbb{R}^k (en rouge) : (5)

$$d^2(z_i, z_{i'}) \geq d^2(\psi_i, \psi_{i'})$$

$\Rightarrow I(z) = p$ est \oplus grande que $I(\psi)$

\Rightarrow on va chercher les axes qui maximisent $I(\psi)$

Donc "déformer" le moins possible le nuage des n points ind. dans \mathbb{R}^p (contés-reduits) \Leftrightarrow maximiser $I(\psi)$

$$\text{Or } I(\psi) = \sum_{i=1}^n w_i \sum_{d=1}^k (\psi_{id} - \bar{\psi}^d)^2 = \sum_{d=1}^k \underbrace{\sum_{i=1}^n w_i (\psi_{id} - \bar{\psi}^d)^2}_{\text{var}(\psi^d)}$$

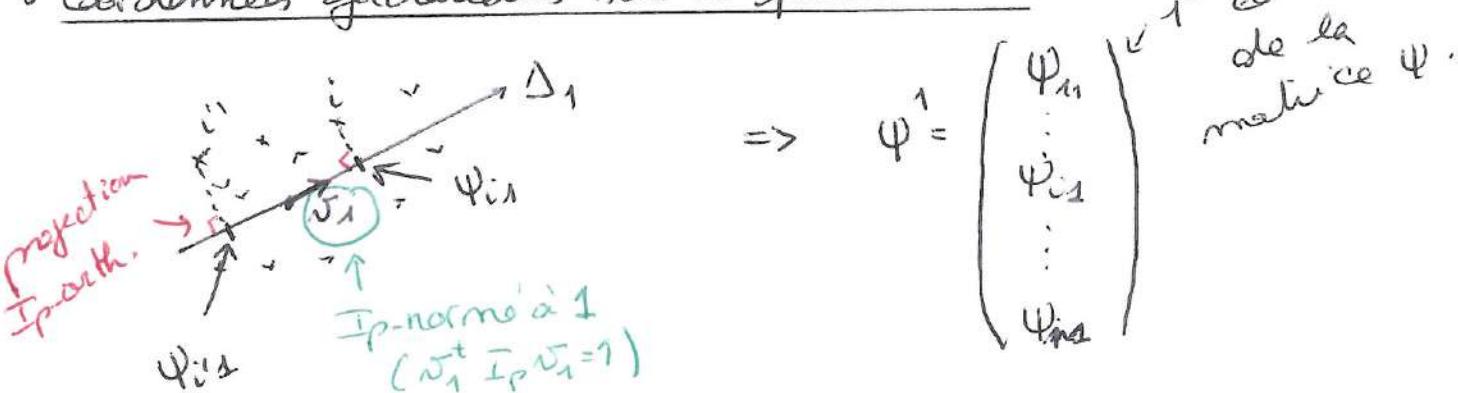
$$= \sum_{d=1}^k \text{var}(\psi^d)$$

\Rightarrow On va chercher à maximiser $\text{var}(\psi^d)$ pour $d=1 \dots k$.

Remarque: si $k=p$, $I(\psi) = I(z) = p$

2.2) Calcul des coordonnées factorielles des individus (scores)

• Coordonnées factorielles sur le premier axe



Par définition :

$$\psi^1 = z^T I_p w_1 = \underline{\underline{z^T w_1}}$$

Remarque: $\psi^1 = w_1 \left[\begin{array}{c} z^1 \\ z^2 \\ \vdots \\ z^p \end{array} \right] + \dots + w_p \left[\begin{array}{c} z^1 \\ z^2 \\ \vdots \\ z^p \end{array} \right]$

1^{re} colonne de z.

= combinaison linéaire des
colonnes de z

= variante "synthétique" regroupant
les colonnes de z

Le problème d'optimisation: Trouver le vecteur $\vec{v}_1 \in \mathbb{R}^p$,

\mathbb{I}_p -normé à 1, tel que $\text{var}(\Psi^1)$ soit maximum.

Remarque: $\bar{\Psi}^1 = 0 \Rightarrow \text{var}(\Psi^1) = \|\Psi^1\|_N^2$

Solution: \vec{v}_1 est le vecteur propre associé à la plus grande valeur propre λ_1 de la matrice des corrélations

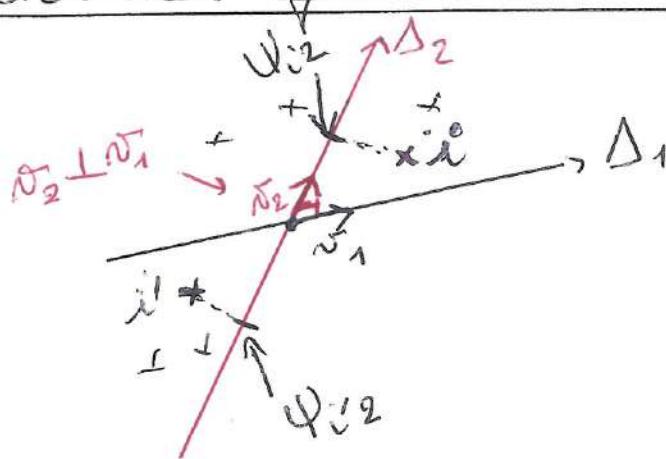
$$R = Z^t N Z$$

Exercice 1: Démontrer ce résultat

Finalement on a :

- $\Psi^1 = Z \vec{v}_1$, \vec{v}_1 est le vecteur propre associé à la plus grande valeur propre de $R = Z^t N Z$
- $\bar{\Psi}^1 = 0$
- $\text{var}(\Psi^1) = \lambda_1$

• Coordonnées factorielles sur le second axe



$$\Psi^1 = \begin{pmatrix} \Psi_{11} \\ \vdots \\ \Psi_{12} \\ \vdots \\ \Psi_{n1} \end{pmatrix} \quad \leftarrow \begin{array}{l} \text{2ème colonne} \\ \text{de } \Psi \end{array}$$

Par définition: $\Psi^2 = Z \vec{v}_2$

Le problème d'optimisation: Trouver le vecteur $\vec{v}_2 \in \mathbb{R}^p$,

\mathbb{I}_p -normé à 1 et \mathbb{I}_p -orthogonale à \vec{v}_1 , tel que $\text{var}(\Psi^2)$ soit maximum

Solution: \vec{v}_2 est le vecteur propre associé à la grande valeur propre λ_2 de la matrice des corrélations R

Final element on a :

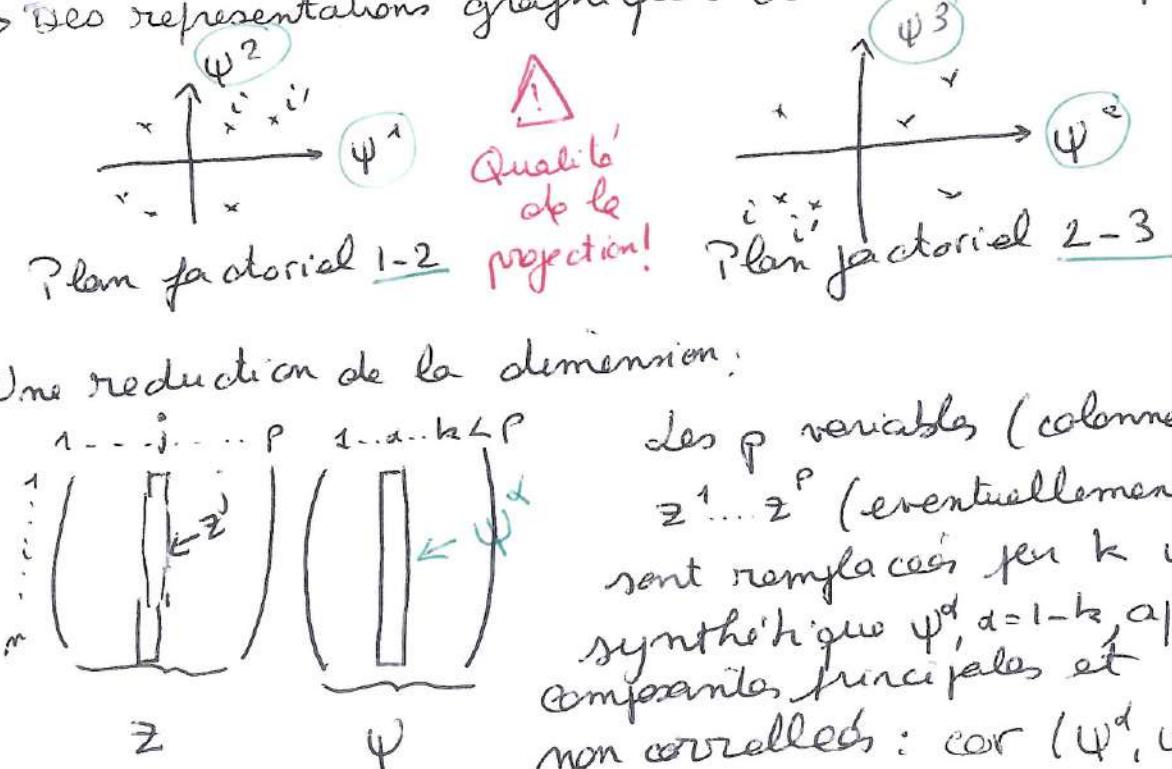
$$\left\{ \begin{array}{l} \cdot \psi^2 = 2 \nu_2 \\ \cdot \bar{\psi}^2 = 0 \\ \cdot \text{var}(\psi^2) = \sigma_2 \end{array} \right.$$

- Idem pour les axes suivants c'est à dire pour $\alpha = 3, \dots, k$ (et $k \leq r$).

• Conclusion: Pour calculer la matrice $\Psi_{m \times k}$ des correlances factorielles des individus (scores des individus), on effectue la décomposition spectrale de la matrice $R = Z^T N Z$ des corrélations: on cherche les I_p -normes à 1 et k vecteurs propres v_1, \dots, v_k de R associés aux plus grandes valeurs propres triées par ordre décroissant $\lambda_1, \dots, \lambda_k$ et on calcule: $\left\{ \begin{array}{l} \Psi^k = Z v_k \\ \text{var}(\Psi^k) = \lambda_k \end{array} \right.$

Ces catalyseurs permettent :

→ Des représentations graphiques des individus projetés



→ Une réduction de la dimension:

des p variables (colonnes)
 $z^1 \dots z^p$ (éventuellement corrélées)
 sont remplacées par k variables
 synthétiques ψ^d , $d=1-k$, appelées
 composantes principales et
 non corrélées : $\text{cor}(\psi^d, \psi^{d'}) = 0 \forall d, d'$

Exemple:

Matrice des corrélations R : Int. bulles
Savon relé app. globale

	Intensité Savon Appac. bulles Sels globale
Int. bulles	1 -0.29 -0.09
Savon relé	-0.29 1 -0.70
app. globale	-0.09 -0.70 0

 $N_1 \quad N_2 \quad N_3$

Int. bulles	0.20 0.92 -0.32
Sav. Sels	-0.71 -0.08 -0.68
Appac. globale	0.66 -0.37 -0.64

Norme: 1 1 1
et $N_1 + N_2 + N_3$

$$\Rightarrow \sqrt{t} V = I_3$$

Matrice Ψ des coordonnées factorielles (scores)

$$\begin{aligned} \Psi^1 &= Z N_1 \\ \Psi^2 &= Z N_2 \\ \Psi^3 &= Z N_3 \end{aligned} \quad \left\{ \Rightarrow \Psi = Z \sqrt{V} \quad n \times p \quad p \times 3 \right.$$

$$\begin{array}{l} \text{1re Composante principale} \\ \hline \begin{matrix} \Psi^1 & \Psi^2 & \Psi^3 \end{matrix} \\ \begin{matrix} \text{St Yorre} \\ \text{Vichy} \\ \text{Quézac} \\ \text{Salvetat} \\ \text{Pennier} \\ \text{Moyenne} \\ \text{Varianc} \end{matrix} \end{array} \Rightarrow \begin{pmatrix} -1.42 & -0.06 & -0.51 \\ -1.52 & -0.90 & 0.37 \\ 1.42 & -0.13 & -0.58 \\ 1.51 & -0.89 & 0.42 \\ 0.00 & 1.94 & 0.30 \\ 0 & 0 & 0 \\ 1.73 & 1.06 & 0.18 \\ (\lambda_1) & (\lambda_2) & (\lambda_3) \end{pmatrix} = \Psi$$

$$\text{et } \Psi^1 + \Psi^2 + \Psi^3$$

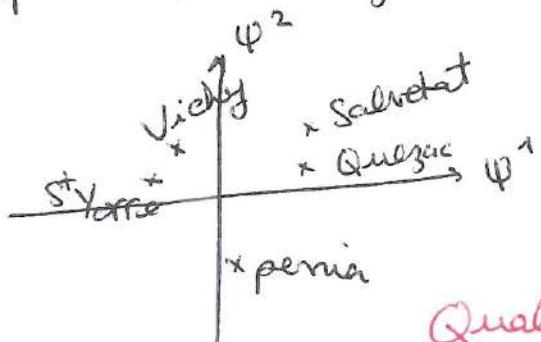
$$\text{et } \underline{\rho} = 3 = \underbrace{\lambda_1 + \lambda_2 + \lambda_3}_{I(z)} \quad \underbrace{I(\Psi)}$$

Variables synthétiques:

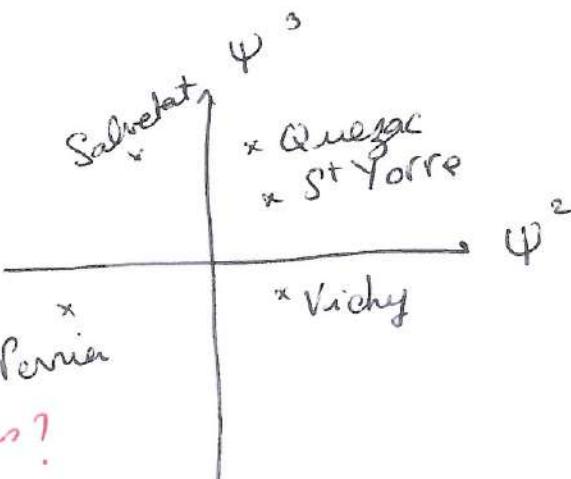
$$\Psi^1 = 0.2 Z^1 + 0.71 Z^2 + 0.66 Z^3$$

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ \text{Int.} & \text{Sav.} & \text{App.} \\ \text{bulles} & \text{Sels} & \text{globale} \end{matrix}$$

Représentations graphiques :



Qualité des projections?



2.3) Inerties associées aux axes.

2.3.1) Pourcentage de l'inertie des nuage de points initial expliquée par l'axe d:

- Ce pourcentage vaut :

$\frac{\lambda_d}{\lambda_1 + \dots + \lambda_p} \times 100 = \frac{\lambda_d}{p}$: mesure le "pouvoir explicatif" de l'axe d.

$$\begin{aligned} & \lambda_d \leftarrow \text{inertie (variance) du nuage des individus} \\ & \text{projeté sur } \Delta^d = I(\Psi^d) \\ & = p = I(\Psi) \\ & = \text{inertie du} \\ & \text{nuage initial} \end{aligned}$$

- Le pourcentage d'inertie expliquée moyen par axe vaut $\frac{1}{p} \times 100$ car la valeur propre moyenne vaut 1

⚠ Il faut faire attention lorsqu'on interprète le pourcentage d'inertie expliquée par un axe car cela dépend du nombre p de variables. Par exemple 10% peut être beaucoup si $p=100$ et peu si $p=10$. En effet $0,1 = \frac{10}{100}$ et $0,1 = \frac{1}{10}$

2.3.2) Pourcentage d'inertie expliquée par les k premiers axes.

- Ce pourcentage vaut

$$\frac{\lambda_1 + \dots + \lambda_k}{p} \times 100$$

- Pour "savoir" combien d'axes retenir pour visualiser le nuage de points (ou réduire la dimension) sans "manquer" d'informations, il existe plusieurs méthodes.

Par exemple:

→ Choisir k pour avoir 90% d'incertes expliquées
par exemple

→ Choisir k tel que $\lambda_k > 1$ et $\lambda_{k+1} < 1$: règle de Kaiser

→ Utiliser un critère statistique

Exemple des casse:

	λ_k	% vagues	% ages cumulés
Axe 1	1.73	57.7	57.7
Axe 2	1.06	35.4	93.3
Axe 3	0.20	6.7	100

2 axes retenus

3) Analyse du nuage des p points-variables de \mathbb{R}^n

On va "réaliser" l'analyse des individus pour faire cette analyse des variables.

3.1) Le principe

On va chercher un sous-espace vectoriel G_k de \mathbb{R}^n , $\Delta \Delta^t + \Delta \Delta^t$
de dimension $k \leq r$, c'est à dire la axes $\Delta^1, \Delta^2, \dots, \Delta^k$

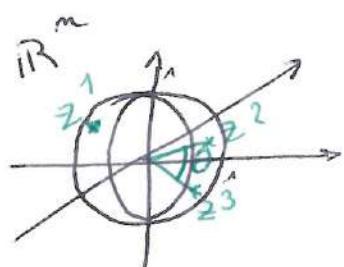
(N-orthogonaux) tels que le nuage des p points-variables N-projetés sur G_k , "déforme" le moins possible les distances entre les variables et donc leurs correlations.

En effet, dans une matrice centré-réduite:

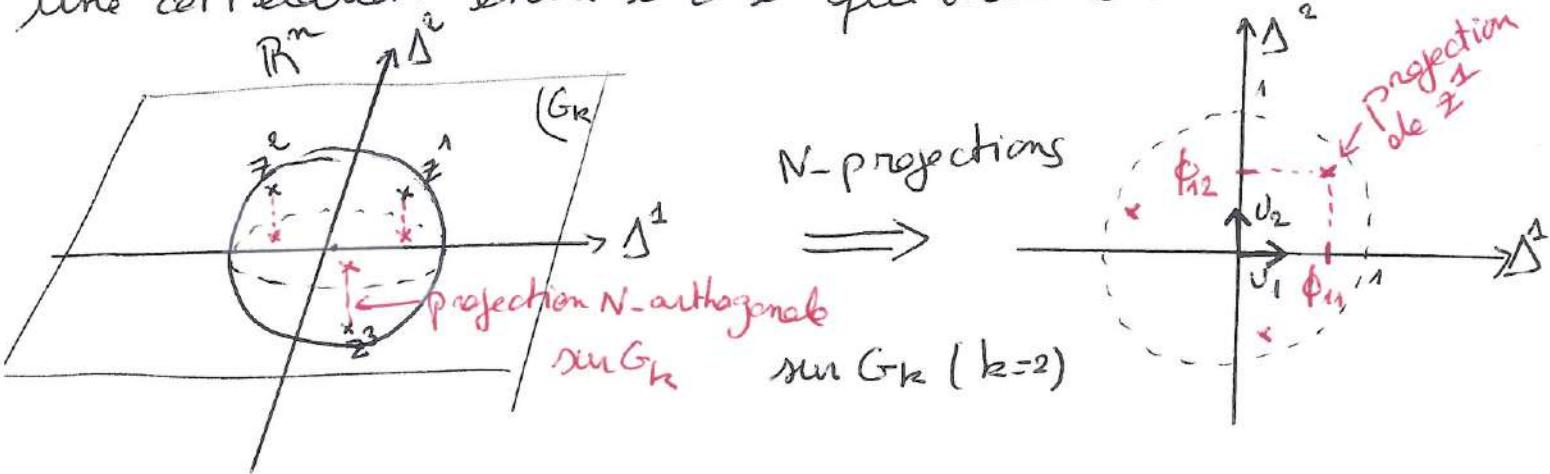
* les points-variables z^i sont sur l'hypersphère unité
 $\forall j, \|z^j\|_N = \text{var}(z^j) = 1$

* la correlation entre deux variables x^i et $x^{i'}$
est égale au cosinus de l'angle entre z^i et $z^{i'}$.
 $r_{jj'} = \langle z^j, z^{j'} \rangle_N = \text{cos} \theta$

\Rightarrow Dans \mathbb{R}^n , si z^i et $z^{i'}$ sont proches en terme de distance (avec la métrique N), alors x^i et $x^{i'}$ sont corrélés.



Exercice 2: Montrer que $d_N^2(z^i, z^{i'}) = 2(1 - r_{jj'})$ où $r_{jj'} = \text{corrélation}$ entre z^i et $z^{i'}$. En déduire que la distance entre deux points variables z^i et $z^{i'}$ varie entre 0 et 2 pour une corrélation entre z^i et $z^{i'}$ qui varie entre -1 et 1.



Matrice centrée-reduite:

$$z = \begin{pmatrix} 1 & \dots & j & \dots & p \\ \vdots & & \vdots & & \vdots \\ n & & z_{ij} & & \vdots \\ \vdots & & \vdots & & \vdots \\ 1 & & z^i \in \mathbb{R}^n \end{pmatrix} \xrightarrow{\text{N-projections sur } G_N}$$

Matrice des coordonnées factuelles des variables

$$\phi = \begin{pmatrix} 1 & \dots & d & \dots & k \\ \vdots & & \vdots & & \vdots \\ n & & \phi_{jd} & & \vdots \\ \vdots & & \vdots & & \vdots \\ 1 & & \phi & & \vdots \end{pmatrix}$$

- On note $v_1, \dots, v_d, \dots, v_k$ les vecteurs directeurs de axes $\Delta^1, \dots, \Delta^k$.

On verra que $v^d = \frac{\psi^d}{\sqrt{\lambda_d}} = d^{\text{ème}} \underbrace{\text{composante principale standardisée}}$

- On appelle ϕ_{pxk} la matrice des coordonnées des p points-variables projetés sur ces axes =
 } matrice des coordonnées factuelles des variables
 } matrice des loadings (saturations) des variables sur les k premières composantes principales.

- On notera ϕ^d la $d^{\text{ème}}$ colonne de ϕ et on verra que les loadings sont des corrélations : $\phi_{j,d} = \text{corr}(\psi^d, x^j)$

3.2) Calcul des coordonnées factuelles des variables (loadings)

- coordonnées factuelles sur le 1^{er} axe:

1^{re} colonne de Φ

$$\Rightarrow \Phi^1 = \begin{pmatrix} \phi_{11} \\ \phi_{j1} \\ \vdots \\ \phi_{p1} \end{pmatrix} \in \mathbb{R}^p$$

Par définition :

$$\Phi^1 = Z^T N U_1$$

Le problème d'optimisation : Trouver le vecteur $U_1 \in \mathbb{R}^n$, N-normé à 1, tel que $\|\Phi^1\|_{\mathbb{R}^p}^2$ soit maximale.

Solution : U_1 est le vecteur propre associé à la plus grande valeur propre λ_1 de la matrice $Z Z^T N$.

Remarque : On verra plus tard que λ_1 est aussi la plus grande valeur propre de $R = Z^T N Z$

Exercice 3 : Démontrer le résultat "relaxion".

- coordonnées factuelles sur le second axe:

Par définition :

$$\Phi^2 = Z^T N U_2$$

de problème : Trouver U_2 tel que $U_2^T N U_2 = 1$ et $U_2^T N U_1 = 0$ qui maximise $\|\Phi^2\|_{\mathbb{R}^p}^2$

Solution : U_2 est le vecteur propre, N-normé à 1, de $Z Z^T N$ associé à la seconde plus grande valeur propre de $Z Z^T N$ pour $d=3, \dots, k$.

- Conclusion: Pour calculer la matrice $\Phi_{p \times k}$ des coordonnées factuelles des variables (loadings des variables), on effectue la décomposition svd'itaire de la matrice $Z Z^t N$ (appelée matrice des produits scalaires des individus): on cherche les k vecteurs propres U_1, \dots, U_k de $Z Z^t N$, N -normés à 1, associés aux valeurs propres $\lambda_1, \dots, \lambda_k$ triées par ordre décroissant et on calcule $\Phi^d = Z^t N U_d$.

- Parallèle avec les coordonnées factuelles des individus,

$$\begin{aligned} v_d, & \left\{ \begin{array}{l} \text{vecteur propre de } Z^t N Z I_p \\ I_p \text{-normé à 1} \end{array} \right. \quad \text{et } \Psi^d = Z I_p v_d \\ u_d, & \left\{ \begin{array}{l} \text{vecteur propre de } Z I_p Z^t N \\ N \text{-normé à 1} \end{array} \right. \quad \text{et } \Phi^d = Z^t N u_d \end{aligned}$$

\Rightarrow on "échange" lignes et colonnes et les motifs

- Écriture matricielle des scores des individus et des loadings des variables

$$\text{Matrice } n \times k \text{ des scores: } \underbrace{\Psi}_{n \times k} = \underbrace{Z V}_{n \times p} \quad \text{où } V = \text{matrice } p \times k \text{ dont les colonnes sont les vect. propres de } R = Z^t N Z$$

$$\text{Matrice } p \times k \text{ des loadings: } \underbrace{\Phi}_{p \times k} = \underbrace{Z^t U}_{p \times n} \quad \text{où } U = \text{matrice dont les colonnes sont les vect. propres de } Z Z^t N.$$

Exemple:

Matrice des produits scalaires des individus $Z Z^t N$:

	Saint-Yorre	Vichy	Quercy	Salvetat	Perrier
Saint-Yorre	0.46	0.41	-0.34	-0.47	-0.06
Vichy	0.41	0.65	-0.45	-0.28	-0.33
Quercy	-0.34	-0.45	0.48	0.41	-0.09
Salvetat	-0.47	-0.28	0.41	0.64	-0.30
Perrier	-0.06	-0.33	0.09	-0.3	0.77

matrice symétrique.

Décomposition spectrale de $Z Z^t N$:

$$\begin{cases} \lambda_1 = 1.73 \\ \lambda_2 = 1.06 \\ \lambda_3 = 0.20 \\ \lambda_4 = 0 \\ \lambda_5 = 0 \end{cases}$$

$$Z Z^t N \xrightarrow{\text{Décomp. spectrale}} U \begin{pmatrix} U_1 & U_2 & U_3 \end{pmatrix} \begin{pmatrix} 1.73 & & \\ & 1.06 & \\ & & 0.20 \end{pmatrix}^{-1} \begin{pmatrix} U_1 & U_2 & U_3 \end{pmatrix}^t$$

$$I_5\text{-Norme: } 1 \quad 1 \quad 1 \quad \underline{N\text{-norme: }} 1 \quad 1 \quad 1$$

pour N -normer à 1,
on \otimes chaque colonne

$$\text{par } \sqrt{n} \left(a = \frac{b}{(b + N_b)^{1/2}} \right)$$

$$\begin{cases} U_1^t N U_2 = 0 \\ U_2^t N U_3 = 0 \\ U_1^t N U_3 = 0 \end{cases}$$

$$U^t N U = I_3$$

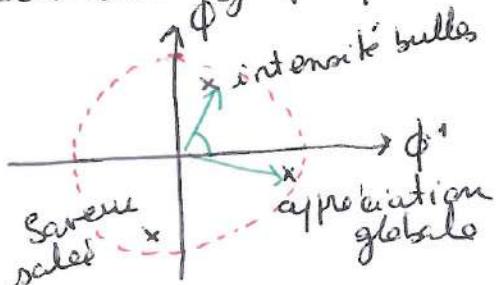
Matrice ϕ des coordonnées

factuelles des variables (loadings):

$$\begin{cases} \phi^1 = Z^t N U_1 \\ \phi^2 = Z^t N U_2 \\ \phi^3 = Z^t N U_3 \end{cases} \Rightarrow \phi = \underbrace{\begin{matrix} Z^t \\ \vdots \\ Z^t \end{matrix}}_{p \times m} \underbrace{\begin{matrix} U_1 \\ \vdots \\ U_3 \end{matrix}}_{m \times 3} = \phi \in \mathbb{R}^{p \times 3}$$

$$\begin{matrix} \phi^1 & \phi^2 & \phi^3 \\ \text{Intensité bulles} & 0.26 & 0.95 & -0.15 \\ \text{Saveur salé} & -0.95 & -0.09 & -0.29 \\ \text{Appréciation globale} & 0.88 & -0.38 & -0.29 \end{matrix} = \phi$$

Représentations graphiques:



* Asymétrie de flèche pour aider à l'interprétation: si points bien projetés, un angle faible (\Rightarrow) forte corrélation

* Asymétrie d'un cercle (clit des corrélations): si point bien projeté alors il est proche du cercle.

4) Lien entre l'analyse des individus et des variables

4.1) Décomposition en valeurs singulières de Z (DVS)

La DVS (ou SVD pour Singular Value Decomposition) de la matrice réelle Z , de rang r , est avec les matrices N sur \mathbb{R}^n et I_p sur \mathbb{R}^p :

$$Z = U \Lambda V^t$$

$m \times p$ $m \times r$ $r \times r$ $r \times p$

où valeurs singulières,

$$\Lambda = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r}) \text{ avec } \lambda_i \text{ valeur propre de } \begin{cases} Z^t N Z = I_p & \\ Z I_p Z^t N = Z Z^t N & \end{cases}$$

U est la matrice dont les colonnes sont les vecteurs propres de $Z Z^t N$ (tris par ordre décroissant des valeurs propres), N -normés et N -orthogonaux : $U^t N U = I_r$

V est la matrice dont les colonnes sont les vecteurs propres de $Z^t N Z$ (tris par ordre décroissant des valeurs propres) et $V^t V = I_r$

→ Cette décomposition en valeurs singulières donne directement les matrices Ψ et Φ des scores et loadings :

$$\begin{cases} \Psi = Z V \\ \Phi = Z^t N U \end{cases}$$

⚠ En pratique, les logiciels font les SVD avec les matrices I_m et I_p .
Donc on pose :

$$\begin{aligned} \tilde{Y} &= N^{1/2} V \Rightarrow \tilde{Y}^t \tilde{Y} = U N U = I_r \\ \Rightarrow Z &= N^{-1/2} \tilde{Y} \Lambda V^t \\ \Rightarrow \underbrace{N^{1/2} Z}_{\tilde{Z}} &= \tilde{Y} \Lambda V^t \Rightarrow \text{SVD "classique"} \\ \text{et } \tilde{Z} &\text{ et } U = N^{-1/2} \tilde{Y} \end{aligned}$$

16

4.2) Familles de passage (ou de transition)

$$\left\{ \begin{array}{l} N_d = \frac{1}{\zeta_d} \\ \mathcal{N}^t \cup \zeta_d \end{array} \right.$$

$$U_d = \frac{1}{\lambda_d} \underbrace{\sum n_d}_{\Psi^d}$$

Exercice 4 : Démontrer ce résultat.

On en déduit deux résultats importants :

$$* \quad v_d = \frac{\Psi^d}{\sqrt{d_d}} = \text{d}^{\text{ème}} \text{ composante principale standardisée}$$

car $d_d = \text{var}(\Psi^d) \Rightarrow \sqrt{d_d}$ est son ecart-type

$\Rightarrow U$ est la matrice des scores standardisés

$$*\Phi^d = \lceil \alpha_d \rceil \wedge_d \quad \left. \begin{array}{c} \\ \end{array} \right\} \Rightarrow \boxed{\begin{array}{l} \phi = V \wedge \\ \psi = U \wedge \end{array}}$$

La SVD de
 Z donne
 directement
 ϕ et ψ .

Exemple:

Exemple:

intensité bulles seuil ratei appréhension globale	$\begin{pmatrix} N_1 \\ 0.20 \\ -0.71 \\ 0.66 \end{pmatrix} \times \sqrt{1.73} =$	$\begin{pmatrix} \phi^1 \\ 0.26 \\ -0.95 \\ 0.88 \end{pmatrix}$
	U_1	Ψ^1

idem pour N_2, N_3 , ϕ^2, ϕ^3

$$\begin{array}{l} \text{St Yorre} \\ \text{Vichy} \\ \text{Quercac} \\ \text{Salvretat} \\ \text{Pennia} \end{array} \left(\begin{array}{c} -1.08 \\ -1.16 \\ 1.08 \\ 1.15 \\ 0.00 \end{array} \right) \times \sqrt{1.73} = \left(\begin{array}{c} -1.42 \\ -1.52 \\ 1.42 \\ 1.51 \\ 0.00 \end{array} \right) \text{ idem pour } U_2, U_3 \\ \Psi^2, \Psi^3$$

Remarque : Si on prend les colonnes de U , Λ et V ,

on a $Z = \underbrace{U^{(k)} \Lambda^{(k)} V^{(k)}} + E_k$

\hat{Z} : meilleure approximation des mounds connus de Z par une matrice de rang $k \leq r$, au sens de la norme de Hilbert-Schmidt.

$$\|A\|_{M,N} = \sqrt{\lambda_r(AM^T A^T N)} \quad \text{avec } \begin{cases} M = I_p \\ N = \text{diag}(c_i) \end{cases}$$

$\Rightarrow \|Z - \hat{Z}\|_{M,N}$ est minimum.

4.3) Correlations variables - composantes principales

Les coordonnées factorielles des variables (les loadings) vérifient la propriété suivante :

$$\Phi_{j1} = \text{cor}(\psi^1, x^j)$$

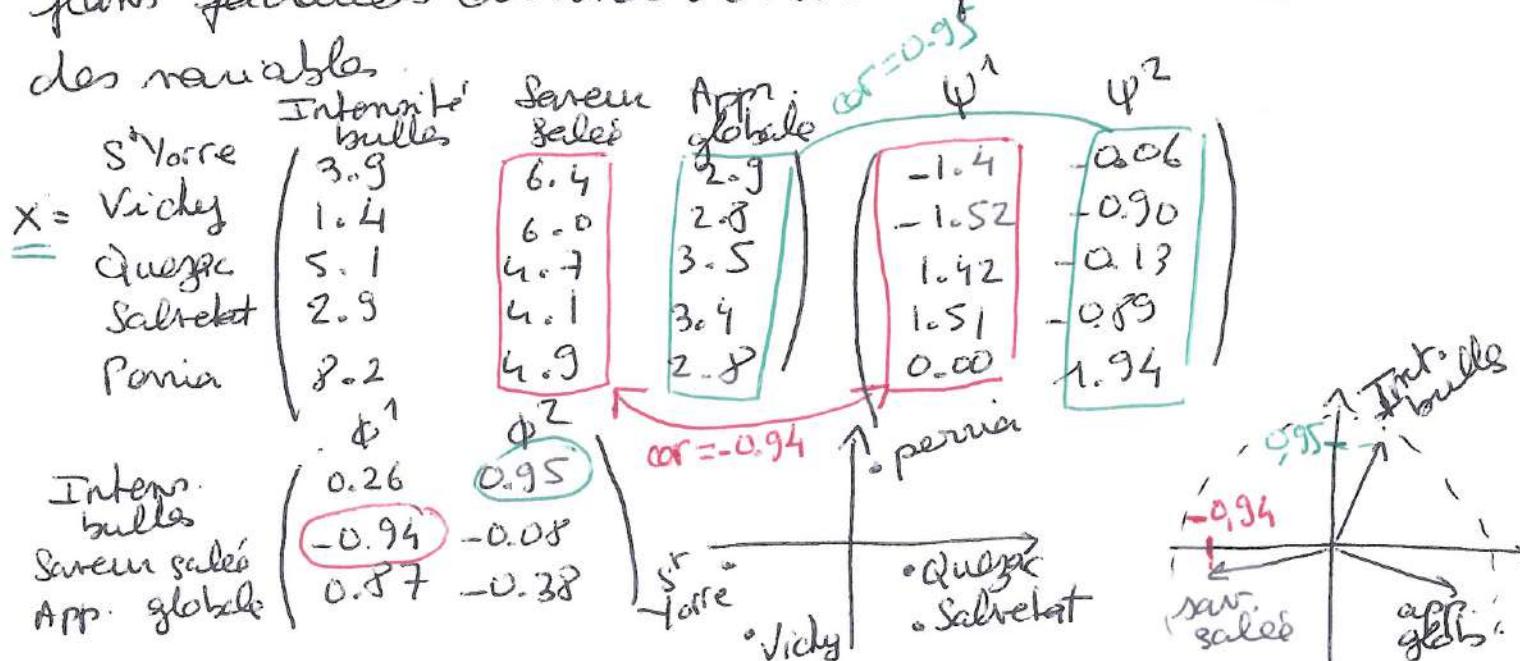
corrélation.

Remarque : En ACP non normée,

$$\Phi_{j1} = s_j \text{cor}(\psi^1, x^j)$$

Exercice 5 : Démontrer ce résultat.

Ce résultat est indispensable pour pouvoir interpréter les plans factoriels des individus en fonction de celle des variables.



Exercice 6 : A notre avis, pourquoi les eaux Quoyac et Salvetat sont-elles appréciées et moins salées tandis que St Yann est déplaisant ? Qu'est-ce qui caractérise l'eau de Pennet?

5) Interprétation des résultats

5.1) Plans factoriels des individus

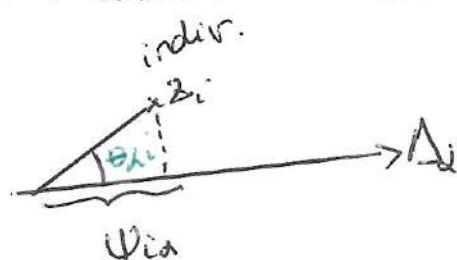
a) Deux individus bien projetés et proches en projection sont effectivement proches de \mathbb{R}^P

⇒ Comment savoir si ils sont bien projetés ? on calcule le cosinus cancé de l'angle entre le point-individu et l'axe (ou le plan)

* Qualité de la représentation de l'individu i sur l'axe A_i :

$$QLT_{A_i}(i) = \cos^2(\theta_{A_i}) = \frac{\Psi_{id}}{\|\mathbf{z}_i\|^2}$$

$$\Rightarrow \begin{cases} \cos^2 = 1 & \text{si } i \text{ parfaitement projeté} \\ = 0 & \text{si très mal projeté} \end{cases}$$



* Qualité de la représentation de l'individu i sur le plan (A_1, A_2):

$$QLT_{A_1, A_2}(i) = \frac{\Psi_{id}^2 + \Psi_{id'}^2}{\|\mathbf{z}_i\|^2}$$

cos²

distance au carré de "i projeté" à l'origine du plan

⚠ Dans SPAD, $QLT_{A_1, A_2}(i)$ déjunit le diamètre des points dans les graphiques des individus en choisissant en option que leur taille soit proportionnelle au \cos^2 .

b) Les individus qui contribuent de manière excessive aux axes factoriels sont source d'instabilité (si on les retire, cela change "fort"). Ils sont donc souvent retirés de l'analyse ou mis en illustratif

\Rightarrow Comment savoir si un individu contribue à un axe ? : on calcule la part (le pourcentage si on \otimes par 100) de l'inertie de l'axe expliquée par l'individu i.

* Contribution relative d'un individu i à la construction de l'axe Δ_2 :

$$\underbrace{CTR_{\Delta_2}(i)}_{\text{contribution}} = \frac{w_i \Psi_{id}^2}{\lambda_2}$$

$$\begin{aligned} \text{Rappel : } \lambda_2 &= \sum_{i=1}^n w_i (\Psi_{id} - \bar{\Psi})^2 \\ &= \sum_{i=1}^n w_i \Psi_{id}^2 \end{aligned}$$

$$\Rightarrow \sum_{i=1}^n CTR_{\Delta_2}(i) = 1$$

\Rightarrow s'exprime bien en pourcentage si \otimes par 100.

⚠ Si tous les poids w_i sont identiques, ce qui est le cas avec $w_i = \frac{1}{n}$, les individus les plus excentrés (pour lesquels $|\Psi_{id}|$ est max) sont ceux qui contribuent le plus.

Remarque : Contribution absolue = $w_i \Psi_{id}^2$

* Contribution relative d'un individu i à la construction du plan (Δ_2, Δ_1) : $\text{CTR}_{\Delta_2, \Delta_1}(i) = \frac{w_i (\Psi_{id}^2 + \Psi_{id1}^2)}{\lambda_2 + \lambda_1}$ Là encore, $\sum_i \text{CTR}_{\Delta_2, \Delta_1}(i) = 1$

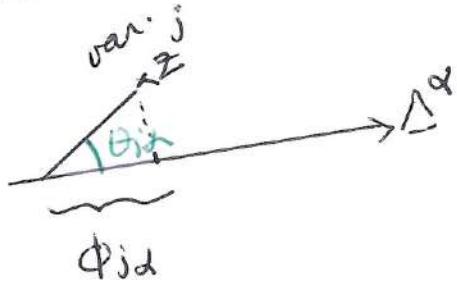
⚠ Dans SPAD, $CTR_{dd,i}(i)$ définit le diamètre des points dans les graphiques des individus en choisissant en option que leur taille soit proportionnelle à leur contribution

5.2) Plans factoriels des variables.

l'angle entre deux variables bien projetées donne une idée de leurs corrélations

* Qualité de la projection de la variable j sur l'axe Δ^d : \cos^2 de l'angle entre le point-variable et l'axe:

$$\underbrace{QLT_d(j)}_{\cos^2} = \cos(\theta_{jd}) = \frac{\phi_{jd}^2}{\underbrace{\|z^j\|^2}_{1 \text{ par construction}}} = \phi_{jd}$$



* Qualité de la projection de la variable j sur le plan $(\Delta^d, \Delta^{d'})$:

$$\underbrace{QLT_{dd'}(j)}_{\cos^2} = \underbrace{\phi_{jd}^2 + \phi_{jd'}^2}_{\text{carré de distance entre point-projecté et l'origine du plan.}}$$

Donc $QLT_{dd'}=1$ si le point-variable projeté est sur le cercle \Rightarrow Qualité de projection est d'autant meilleure que la flèche est proche du centre des corrélations

5.3) Interprétation des plans des individus à partir de certaines variables: On utilise $\phi_{jd} = \text{corr}(u^d, z^j)$