

Reconnaissance des formes

Analyse en Composantes Principales

- Chapitre 2 -

ANALYSE EN COMPOSANTES PRINCIPALES

Objectifs

Contexte : Chaque individu x_i du tableau \mathbf{X} est considéré comme un point d'un espace vectoriel \mathcal{E} de dimension p .

L'ensemble des individus constitue un nuage de points dans \mathcal{E} .

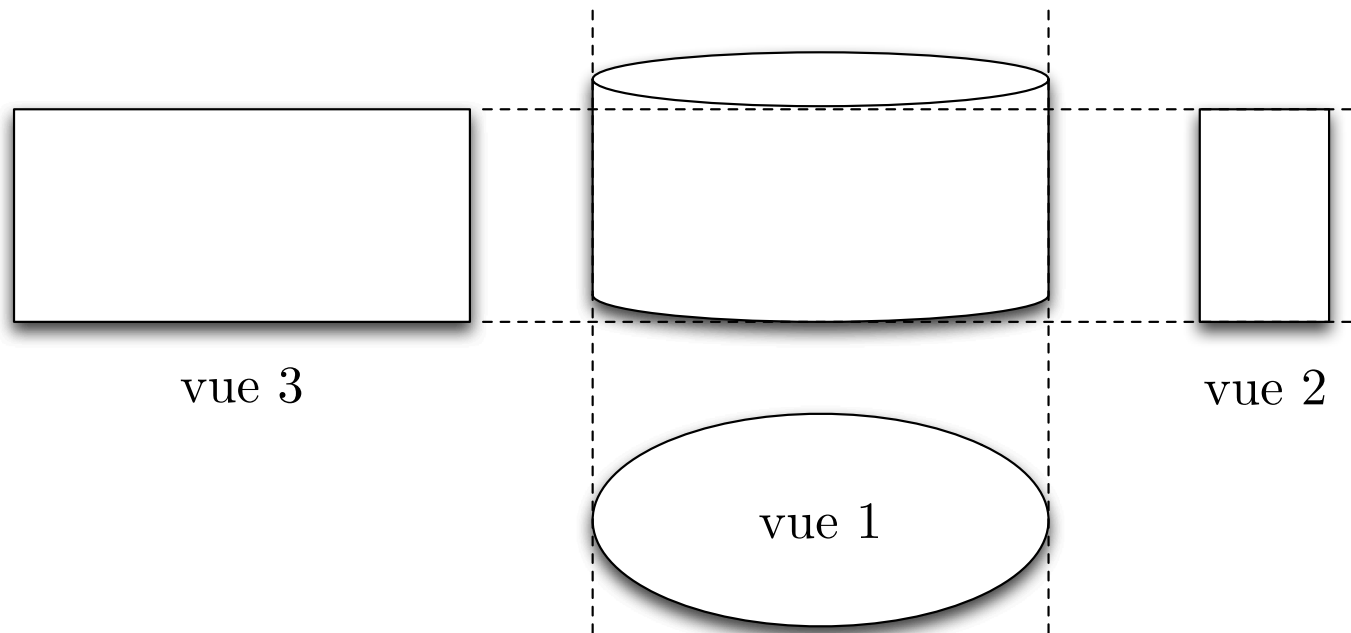
On note m son point moyen, appelé aussi *centre de gravité*.

Objectif : On cherche à réduire le nombre p de variables en préservant au maximum la structure du nuage, afin de le représenter le plus fidèlement possible.

Principe : L'ACP vise à projeter les données dans un sous-espace approprié, de dimension plus faible, préservant la topologie du nuage.

ANALYSE EN COMPOSANTES PRINCIPALES

Exemple



ANALYSE EN COMPOSANTES PRINCIPALES

Exemple plus réaliste

cidre	odeur	sucré	acide	amer	astringence	suffocante	piquante	alcool	parfum	fruité
1	2.14	1.86	3.29	2.29	2	0.14	2.29	1.86	1.29	1.29
2	2.43	0.79	2.71	2.57	2	0.43	2.57	2.86	0.43	0.14
3	2.71	3.14	2.57	2.57	1.43	0.14	2.14	0.86	2.29	1.71
4	3	3.71	2.14	2.07	1.57	0	1.29	1	3.14	3.14
5	3.43	1.29	2.86	3.14	2.17	1	1.86	2.86	1.14	0.29
6	3.14	0.86	2.86	3.79	2.57	0.14	1.71	3.29	0.14	0
7	3.14	1.14	2.86	2.86	2	0.43	1.71	1.86	0.14	0
8	2.43	3.71	3.21	1.57	1.71	0	1	0.57	2.57	2.86
9	5.1	2.86	2.86	3.07	1.79	1.71	0.43	1.43	0.57	2.71
10	3.07	3.14	2.57	3	2	0	0.43	1.29	2.57	3.07

Notes obtenues pour 10 cidres selon 10 critères lors d'un concours agricole.

ESPACE DES INDIVIDUS

Métrie

Métrie : L'ACP repose sur les distances entre individus dans \mathcal{E} . Le choix de la métrie a une influence fondamentale sur le résultat de l'analyse.

Définition : Soit \mathbf{M} une matrice définie positive de dimension p . La fonction suivante $d_M : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^+$ définit une métrie

$$d_M^2(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_M^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)$$

Cette distance est appelée distance de Mahalanobis lorsque $\mathbf{M} = \Sigma^{-1}$, où Σ est la matrice de variance-covariance des données.

Produit scalaire : La métrie définie ci-dessus dérive du produit scalaire

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle_M = \mathbf{x}_i^\top \mathbf{M} \mathbf{x}_j$$

On dit que \mathbf{x}_i et \mathbf{x}_j sont M -orthogonaux si $\langle \mathbf{x}_i, \mathbf{x}_j \rangle_M = 0$.

ANALYSE DES DONNÉES

Métrie

Utiliser la métrie $\mathbf{M} = \mathbf{T}^\top \mathbf{T}$ sur le tableau de données \mathbf{X} est équivalent à travailler avec la métrie euclidienne sur le tableau transformé $\mathbf{X}\mathbf{T}^\top$.

Tableau transformé : Lorsqu'on travaille sur le tableau transformé comme ci-dessus, il est alors possible d'utiliser la norme euclidienne. En effet,

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle_M = \mathbf{x}_i^\top (\mathbf{T}^\top \mathbf{T}) \mathbf{x}_j = (\mathbf{T}\mathbf{x}_i)^\top (\mathbf{T}\mathbf{x}_j) = \langle \mathbf{T}\mathbf{x}_i, \mathbf{T}\mathbf{x}_j \rangle_I$$

Réciproque : Pour toute matrice définie positive \mathbf{M} , il existe une matrice définie positive \mathbf{T} telle que $\mathbf{M} = \mathbf{T}^\top \mathbf{T}$. On notera improprement $\mathbf{T} = \mathbf{M}^{\frac{1}{2}}$.

Appliquer préalablement la transformation $\mathbf{X}\mathbf{T}^\top \rightarrow \mathbf{X}$ permet de simplifier les traitements.

ANALYSE DES DONNÉES

Métriques particulières

Métrique euclidienne : Elle est obtenue pour $M = I$.

L'une des difficultés rencontrées avec la métrique euclidienne est qu'elle privilégie les variables les plus dispersées et dépend donc de leur unité de mesure.

Métrique réduite : Elle consiste à prendre $M = D_{1/\sigma^2}$, où D_{1/σ^2} est la matrice diagonale de termes diagonaux les inverses $\frac{1}{\sigma_i^2}$ des variances des variables.

$$D_{1/\sigma^2} = \begin{pmatrix} \frac{1}{\sigma_1^2} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma_p^2} \end{pmatrix}$$

Cette métrique permet de s'affranchir de l'unité de mesure des variables, et de donner la même importance à chaque variable dans le calcul de la distance.

ANALYSE DES DONNÉES

Tableau de données centrées réduites

Utiliser la métrique $M = D_{1/\sigma^2} = D_{1/\sigma}^\top D_{1/\sigma}$ sur le tableau de données X revient à travailler avec la métrique euclidienne sur le tableau transformé $XD_{1/\sigma}^\top$.

En effet :

$$\begin{aligned}d_M^2(\mathbf{x}_i, \mathbf{x}_j) &= \|\mathbf{x}_i - \mathbf{x}_j\|_M^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top D_{1/\sigma}^\top D_{1/\sigma} (\mathbf{x}_i - \mathbf{x}_j) \\ &= (D_{1/\sigma} \mathbf{x}_i - D_{1/\sigma} \mathbf{x}_j)^\top (D_{1/\sigma} \mathbf{x}_i - D_{1/\sigma} \mathbf{x}_j)\end{aligned}$$

Il est équivalent de travailler avec la métrique D_{1/σ^2} sur le tableau X , ou avec la métrique euclidienne I sur le tableau centré réduit Z composé des données :

$$z_i^j = \frac{x_i^j - \bar{x}^j}{\sigma_j}$$

Le tableau de données centré réduit Z se calcule matriciellement ainsi :

$$Z = YD_{1/\sigma} = (X - \mathbf{1}m^\top) D_{1/\sigma}.$$

INERTIES

Inertie par rapport à un point

Définition : L'inertie du nuage de points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ en un point quelconque \mathbf{a} est donnée par

$$I_a = \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{a}\|_M^2$$

Propriété : L'inertie du nuage de points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ en son point moyen \mathbf{m} , ou *centre de gravité*, est

$$\begin{aligned} I_m &= \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{m}\|_M^2 = \frac{1}{2} \sum_{i,j=1}^n p_i p_j \|\mathbf{x}_i - \mathbf{x}_j\|_M^2 \\ &= \text{Trace}(\Sigma \mathbf{M}) \end{aligned}$$

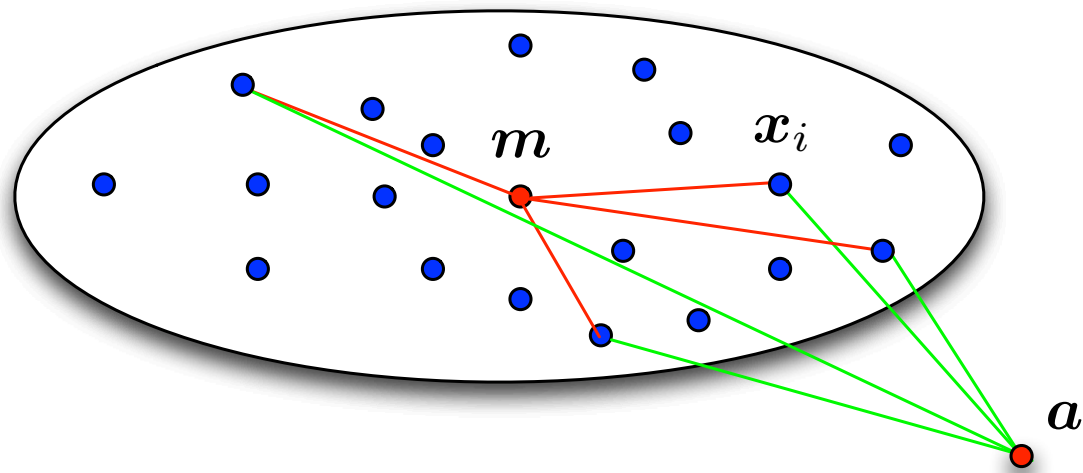
Cette propriété a été démontrée au chapitre précédent.

Propriété : Pour un tableau de données centrées réduites, on a

$$I_m = \text{Trace}(\mathbf{R}) = p$$

INERTIES

Inertie par rapport à un point



INERTIES

Théorème de Huygens

Propriété : Soit \mathbf{m} le centre de gravité du nuage de points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, et \mathbf{a} un point quelconque de \mathbb{R}^p . L'inertie I_a du nuage au point \mathbf{a} est donnée par

$$I_a = I_m + d_M^2(\mathbf{a}, \mathbf{m})$$

En conséquence, I_a est minimum pour $\mathbf{a} = \mathbf{m}$.

Démonstration :

$$\begin{aligned} I_a &= \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{m} - (\mathbf{a} - \mathbf{m})\|_M^2 \\ &= \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{m}\|_M^2 + \sum_{i=1}^n p_i \|\mathbf{a} - \mathbf{m}\|_M^2 - 2 \sum_{i=1}^n p_i \langle \mathbf{x}_i - \mathbf{m}, \mathbf{a} - \mathbf{m} \rangle_M \\ &= I_m + d_M^2(\mathbf{a}, \mathbf{m}) \end{aligned}$$

INERTIES

Inertie par rapport à un axe

Définition : L'inertie du nuage de points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ par rapport à un axe Δ est définie par

$$I_{\Delta} = \sum_{i=1}^n p_i d_M^2(\mathbf{x}_i, \Delta)$$

Cette inertie quantifie la dispersion du nuage des individus autour de Δ .

Rappels : Soit $\Delta(\mathbf{a}, \mathbf{u})$ un axe de vecteur directeur normé \mathbf{u} passant par \mathbf{a} .

On rappelle que

$$d_M^2(\mathbf{x}_i, \Delta) = \min_{v_i} d^2(\mathbf{z}_i, \mathbf{a} + v_i \mathbf{u})$$

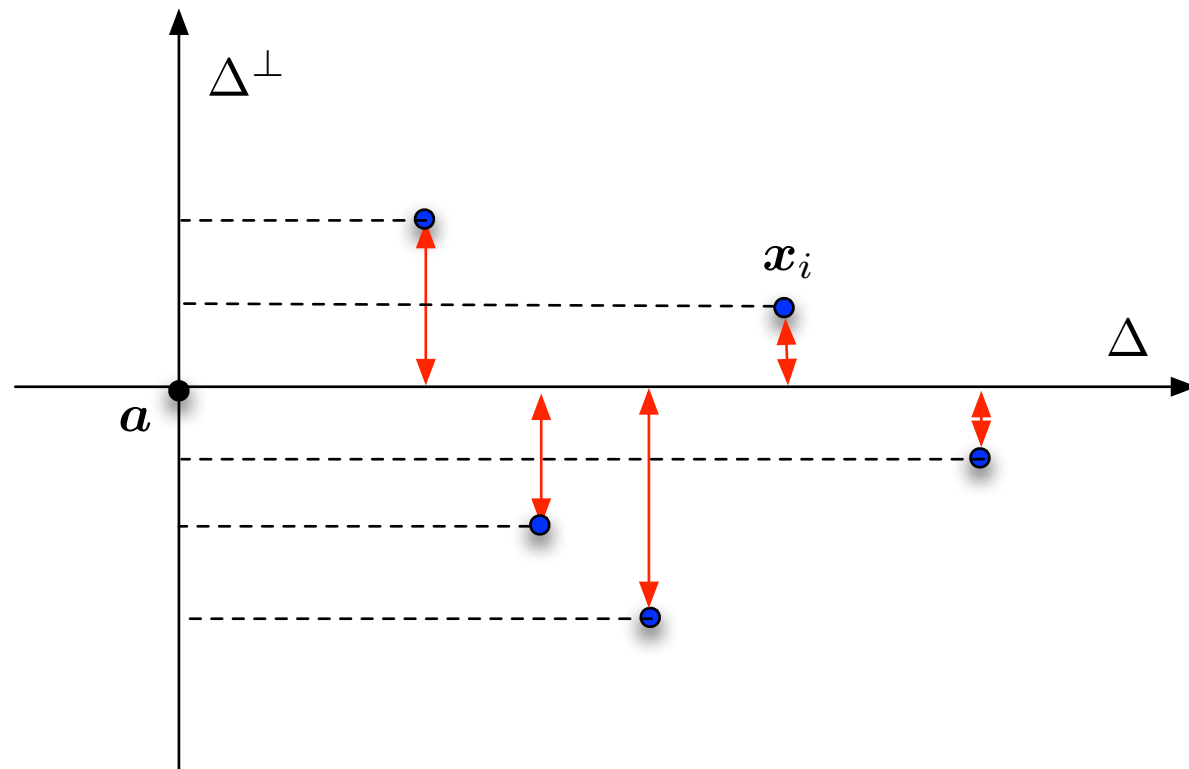
En conséquence, l'inertie du nuage $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ par rapport à $\Delta(\mathbf{a}, \mathbf{u})$ est définie par

$$I_{\Delta(\mathbf{a}, \mathbf{u})} = \min_{\mathbf{v}} \sum_{i=1}^n p_i d_M^2(\mathbf{x}_i, \mathbf{a} + v_i \mathbf{u})$$

avec $\mathbf{v} = [v_1, \dots, v_n]^T$.

INERTIES

Inertie par rapport à un axe



ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du premier axe

Définition du problème : L'axe Δ recherché est celui représentant le mieux les données au sens de l'inertie I_{Δ} minimum. Il est donc obtenu par

$$\min_{\mathbf{u}, \mathbf{a}} I_{\Delta(\mathbf{a}, \mathbf{u})} \quad \text{avec} \quad I_{\Delta(\mathbf{a}, \mathbf{u})} = \min_{\mathbf{v}} \sum_{i=1}^n p_i d_M^2(\mathbf{x}_i, \mathbf{a} + v_i \mathbf{u})$$

où Δ est défini par un point \mathbf{a} et un vecteur directeur normé \mathbf{u} .

Détermination de \mathbf{a} : En appliquant le Théorème de Huygens, on trouve

$$I_{\Delta(\mathbf{a}, \mathbf{u})} = I_{\Delta(\mathbf{m}, \mathbf{u})} + d_M^2(\mathbf{a}, \mathbf{m})$$

En conséquence, $I_{\Delta(\mathbf{a}, \mathbf{u})}$ est minimum pour $\mathbf{a} = \mathbf{m}$.

L'axe recherché passe par \mathbf{m} .

ANALYSE EN COMPOSANTES PRINCIPALES

Centrage et normalisation des données

Nous considérerons les données centrées et normalisées

$$\mathbf{Z} = (\mathbf{X} - \mathbf{1} \mathbf{m}^\top) \mathbf{M}^{\frac{1}{2}}$$

pour :

- ▷ faire correspondre \mathbf{m} avec l'origine $\mathbf{0}$ du repère
- ▷ travailler avec la métrique euclidienne sur le nuage de points $\{z_1, \dots, z_n\}$

ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du premier axe

Redéfinition du problème : L'axe Δ recherché passe par l'origine $\mathbf{0}$ et est solution du problème

$$\min_{\mathbf{u}} I_{\Delta}(\mathbf{u}) \quad \text{avec} \quad I_{\Delta}(\mathbf{u}) = \min_{\mathbf{v}} \sum_{i=1}^n p_i d^2(\mathbf{z}_i, v_i \mathbf{u})$$

Calcul du gradient : En développant la fonction coût, notée $J_{\Delta}(\mathbf{u}, \mathbf{v})$, on obtient

$$J_{\Delta}(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^n p_i [\|\mathbf{z}_i\|^2 - 2v_i \langle \mathbf{z}_i, \mathbf{u} \rangle + v_i^2 \|\mathbf{u}\|^2]$$

En conséquence

$$\nabla_{\mathbf{u}} J_{\Delta}(\mathbf{u}, \mathbf{v}) = -2\mathbf{Z}^{\top} \mathbf{D} \mathbf{v} + 2\|\mathbf{v}\|_D^2 \mathbf{u}$$

$$\nabla_{\mathbf{v}} J_{\Delta}(\mathbf{u}, \mathbf{v}) = -2\mathbf{D} \mathbf{Z} \mathbf{u} + 2\|\mathbf{u}\|^2 \mathbf{D} \mathbf{v}$$

avec \mathbf{D} la matrice diagonal de termes diagonaux p_i .

ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du premier axe

Conditions d'optimalité : Celles-ci s'écrivent

$$\begin{cases} \mathbf{Z}^\top \mathbf{D} \mathbf{v} = \|\mathbf{v}\|_D^2 \mathbf{u} \\ \mathbf{D} \mathbf{Z} \mathbf{u} = \|\mathbf{u}\|^2 \mathbf{D} \mathbf{v} \end{cases}$$

En combinant les 2 équations, on obtient

$$\mathbf{Z}^\top \mathbf{D} \mathbf{Z} \mathbf{u} = \underbrace{\|\mathbf{u}\|^2 \|\mathbf{v}\|_D^2}_{\lambda} \mathbf{u}$$

La solution est un vecteur propre \mathbf{u} de la matrice de covariance $\mathbf{\Sigma} = \mathbf{Z}^\top \mathbf{D} \mathbf{Z}$

Lequel ?

ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du premier axe

Lequel ? En injectant le résultat dans la fonction coût, on obtient

$$\begin{aligned} I_{\Delta}(\mathbf{u}) &= \sum_{i=1}^n p_i [\|\mathbf{z}_i\|^2 - 2v_i \langle \mathbf{z}_i, \mathbf{u} \rangle + v_i^2 \|\mathbf{u}\|^2] \\ &= \text{Trace}(\boldsymbol{\Sigma}) - 2 \mathbf{u}^{\top} \mathbf{Z}^{\top} \mathbf{D} \mathbf{v} + \|\mathbf{v}\|_D^2 \|\mathbf{u}\|^2 \\ &= \text{Trace}(\boldsymbol{\Sigma}) - 2 \|\mathbf{v}\|_D^2 \|\mathbf{u}\|^2 + \|\mathbf{v}\|_D^2 \|\mathbf{u}\|^2 \\ &= \text{Trace}(\boldsymbol{\Sigma}) - \lambda(\mathbf{u}) \end{aligned}$$

où $\lambda(\mathbf{u})$ est la valeur propre associée à \mathbf{u}

Le vecteur propre recherché est donc celui associé à la plus grande valeur propre

ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du premier axe

Soient $\lambda_1 \geq \dots \geq \lambda_p$ les valeurs propres de la matrice de covariance Σ , de vecteurs propres unitaires associés $\mathbf{u}_1, \dots, \mathbf{u}_p$

Choix du premier axe : L'inertie $I_{\Delta(\mathbf{u})} = \text{Trace}(\Sigma) - \lambda(\mathbf{u})$ est minimum pour le couple $(\mathbf{u}_1, \lambda_1)$

Inertie correspondante : La valeur λ_1 est, par définition, l'inertie expliquée par l'axe $\Delta(\mathbf{u}_1)$.

INERTIES

Théorème d'Inclusion

- ▷ Soient \mathcal{E}_1 et \mathcal{E}_1^\perp les sous-espaces affines de directions respectives Δ et Δ^\perp passant par l'origine.
- ▷ L'inertie du nuage par rapport à \mathcal{E}_1 et \mathcal{E}_1^\perp est définie par

$$I_{\mathcal{E}_1} = \sum_{i=1}^n p_i d^2(\mathbf{z}_i, \mathcal{E}_1)$$
$$I_{\mathcal{E}_1^\perp} = \sum_{i=1}^n p_i d^2(\mathbf{z}_i, \mathcal{E}_1^\perp)$$

Voir le schéma ci-après. Par le Théorème de Pythagore, on a

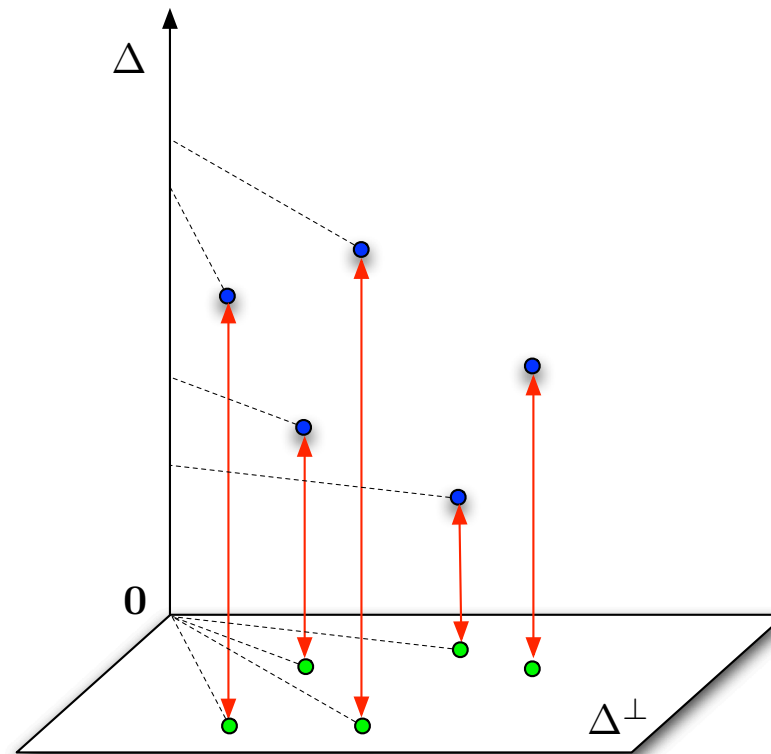
$$I_0 = I_{\mathcal{E}_1} + I_{\mathcal{E}_1^\perp}$$

où I_0 est l'inertie du nuage par rapport à l'origine.

On en déduit le Théorème d'Inclusion ci-après.

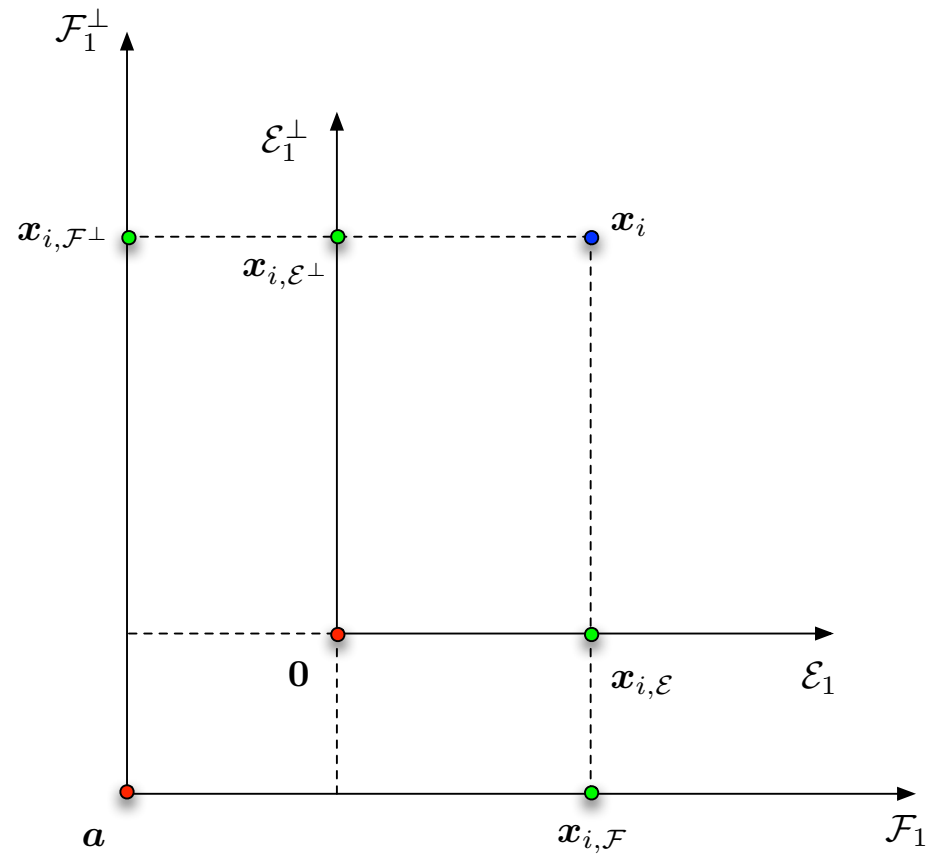
INERTIES

Inertie par rapport à un sous-espace affine



INERTIES

Théorème d'Inclusion



ANALYSE EN COMPOSANTES PRINCIPALES

Théorème d'Inclusion

Théorème d'Inclusion : Soit \mathcal{E}_k un sous-espace, de dimension k , d'inertie expliquée maximum.

Le sous-espace \mathcal{E}_{k+1} de dimension $k + 1$, d'inertie expliquée maximum, est la somme directe de \mathcal{E}_k et du sous-espace de dimension 1 de \mathcal{E}_k^\perp dont l'inertie expliquée y est maximum.

ANALYSE EN COMPOSANTES PRINCIPALES

Théorème d'Inclusion

Recherches des axes : Etant donné le Théorème d'Inclusion, pour déterminer le sous-espace optimum \mathcal{E}_k , il est possible de procéder séquentiellement :

- recherche du premier axe d'inertie portée maximum
- recherche du deuxième axe, orthogonal au premier, d'inertie portée maximum
- etc.

ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du deuxième axe

- ▷ On construit le tableau de données résiduel $\mathbf{Z}_1 = \mathbf{Z} - \mathbf{v}_1 \mathbf{u}_1^\top$. Il est aisé de montrer que les données correspondantes appartiennent à Δ^\perp

$$\begin{aligned}\mathbf{Z}_1 \mathbf{u}_1 &= \mathbf{Z} \mathbf{u}_1 - \|\mathbf{u}_1\|^2 \mathbf{v}_1 \\ &= 0\end{aligned}$$

d'après la deuxième condition d'optimalité, car \mathbf{D} est inversible.

- ▷ La matrice $\mathbf{Z}_1^\top \mathbf{D} \mathbf{Z}_1$ admet les mêmes couples valeurs/vecteurs propres $(\lambda_k, \mathbf{u}_k)$ que $\mathbf{Z}^\top \mathbf{D} \mathbf{Z}$, pour $k \neq 1$,

$$\begin{aligned}(\mathbf{Z}_1^\top \mathbf{D} \mathbf{Z}_1) \mathbf{u}_k &= (\mathbf{Z}^\top \mathbf{D} \mathbf{Z}) \mathbf{u}_k \quad k \neq 1 \\ &= \lambda_k \mathbf{u}_k\end{aligned}$$

car $(\mathbf{u}_1 \perp \mathbf{u}_k)$, excepté λ_1 qui devient nulle

$$\begin{aligned}(\mathbf{Z}_1^\top \mathbf{D} \mathbf{Z}_1) \mathbf{u}_1 &= [\mathbf{Z} - \mathbf{v}_1 \mathbf{u}_1^\top]^\top \mathbf{D} [\mathbf{Z} - \mathbf{v}_1 \mathbf{u}_1^\top] \\ &= (1 - 2 + 1) \|\mathbf{u}_1\|^2 \|\mathbf{v}_1\|_D^2 \mathbf{u}_1 = 0\end{aligned}$$

ANALYSE EN COMPOSANTES PRINCIPALES

Détermination du deuxième axe

Pour déterminer le deuxième axe, on pose le même problème que pour le premier, avec le tableau \mathbf{Z}_1 . On aboutit alors à \mathbf{u}_2 . En effet

Choix du deuxième axe : L'inertie $I_{\Delta(\mathbf{u})} = \text{Trace}(\mathbf{\Sigma}_1) - \lambda(\mathbf{u})$ est minimum pour le couple $(\mathbf{u}_2, \lambda_2)$

Remarque : $\text{Trace}(\mathbf{\Sigma}_1) = \text{Trace}(\mathbf{\Sigma}) - \lambda_1$

Inerties correspondantes : L'inertie expliquée par l'axe $\Delta(\mathbf{u}_2)$ vaut λ_2 .

Pour l'espace affine de direction $\Delta_{\mathbf{u}_1} + \Delta_{\mathbf{u}_2}$, par le Théorème de Pythagore, elle vaut $\lambda_1 + \lambda_2$.

et ainsi de suite...

ANALYSE EN COMPOSANTES PRINCIPALES

Axes principaux et inertie expliquée

- ▷ Les axes $\Delta(\mathbf{u}_k)$ sont appelés **axes factoriels** ou **axes principaux**
- ▷ L'inertie expliquée par l'axe $\Delta(\mathbf{u}_k)$ est la valeur propre λ_k de Σ
- ▷ L'inertie expliquée par le sous-espace factoriel \mathcal{E}_{k_0} engendré par les k_0 premiers axes principaux est

$$\lambda_1 + \dots + \lambda_{k_0}$$

- ▷ Le **taux d'inertie expliquée** par \mathcal{E}_{k_0} est donné par

$$\frac{\lambda_1 + \dots + \lambda_{k_0}}{\sum_{k=1}^p \lambda_k}$$

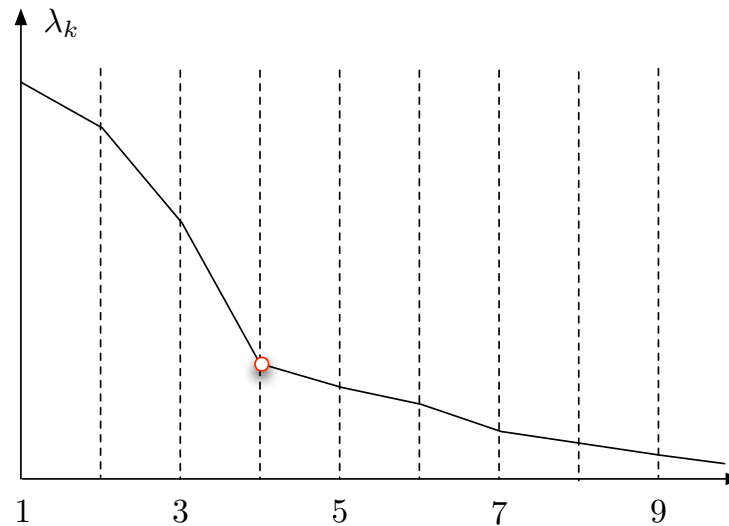
$$\text{car } I_0 = \text{Trace}(\Sigma) = \sum_{k=1}^p \lambda_k$$

ANALYSE EN COMPOSANTES PRINCIPALES

Nombre d'axes à retenir

L'ACP vise à réduire la dimension de l'espace des individus. On souhaite conserver aussi peu d'axes principaux que possible. Des critères empiriques :

- ▷ Retenir les axes offrant une corrélation suffisante entre composantes principales et variables initiales. Voir plus loin.
- ▷ Fixer un taux minimum d'inertie expliquée, par exemple 0.80
- ▷ Rechercher un coude dans l'éboulis des valeurs propres λ_k



ANALYSE EN COMPOSANTES PRINCIPALES

Composantes principales

Coordonnées des individus : La coordonnée $c_{\ell k}$ de l'individu z_{ℓ} sur l'axe principal $\Delta(\mathbf{u}_k)$ est donnée par sa projection sur \mathbf{u}_k

$$c_{\ell k} = \mathbf{z}_{\ell}^{\top} \mathbf{u}_k$$

Composantes principales : Il s'agit des vecteurs \mathbf{c}_k des coordonnées des individus sur l'axe principal $\Delta(\mathbf{u}_k)$, c'est à dire

$$\mathbf{c}_k = \mathbf{Z} \mathbf{u}_k$$

ANALYSE EN COMPOSANTES PRINCIPALES

Propriétés des composantes principales

Moyenne arithmétique :

$$\bar{\mathbf{c}}_k = \mathbf{1}^\top \mathbf{D} \mathbf{c}_k = \mathbf{1}^\top \mathbf{D} \mathbf{Z} \mathbf{u}_k = 0$$

car $\mathbf{1}^\top \mathbf{D} \mathbf{Z} = \mathbf{0}$, le tableau étant centré

Variance :

$$\text{var}(\mathbf{c}_k) = \mathbf{c}_k^\top \mathbf{D} \mathbf{c}_k = \mathbf{u}_k^\top \mathbf{Z}^\top \mathbf{D} \mathbf{Z} \mathbf{u}_k = \lambda_k$$

car \mathbf{u}_k est vecteur propre de $\Sigma = \mathbf{Z}^\top \mathbf{D} \mathbf{Z}$

Covariance : Il s'agit des vecteurs \mathbf{c}_k des coordonnées des individus sur l'axe principal $\Delta(\mathbf{u}_k)$, c'est à dire

$$\text{cov}(\mathbf{c}_k, \mathbf{c}_\ell) = \mathbf{c}_k^\top \mathbf{D} \mathbf{c}_\ell = \mathbf{u}_k^\top \mathbf{Z}^\top \mathbf{D} \mathbf{Z} \mathbf{u}_\ell = 0$$

car $(\mathbf{u}_k \perp \mathbf{u}_\ell)$. Les composantes principales sont donc décorrélées.

ANALYSE EN COMPOSANTES PRINCIPALES

Qualité de la représentation d'un individu

Définition : La qualité de la représentation d'un individu ℓ sur l'axe principal k est donnée par le cosinus de leur angle, soit

$$\cos(\mathbf{z}_\ell, \mathbf{u}_k) = \frac{\mathbf{z}_\ell^\top \mathbf{u}_k}{\|\mathbf{z}_\ell\|}$$

Comme $\mathbf{z}_\ell^\top \mathbf{u}_k = c_{\ell k}$, et $\{\mathbf{u}_k\}_k$ est une base orthonormée de \mathbb{R}^p , on a

$$\cos^2(\mathbf{z}_\ell, \mathbf{u}_k) = \frac{c_{\ell k}^2}{\sum_{k=1}^p c_{\ell k}^2}$$

ANALYSE EN COMPOSANTES PRINCIPALES

Contribution d'un individu à une composante principale

On a vu que $\text{var}(\mathbf{c}_k) = \lambda_k = \sum_{i=1}^n p_i c_{ik}^2$.

Définition : On définit la contribution de l'individu ℓ à un axe principal k par

$$\theta_{\ell k} = \frac{p_{\ell} c_{\ell k}^2}{\lambda_k}$$

Interprétation : La contribution d'un individu ℓ doit être relativisée par rapport à son poids p_{ℓ} . Ainsi, si $\theta_{\ell k} > \alpha p_{\ell}$ avec $\alpha > 2$, on peut juger la contribution de l'individu ℓ importante.

Sur-représentation : L'individu ℓ peut être considéré comme sur-représenté par l'axe principal k si $\theta_{\ell k} > 0.25$. Il a eu une influence trop importante dans la détermination de celui-ci, risquant de perturber la définition des autres axes.

ANALYSE EN COMPOSANTES PRINCIPALES

Espace des variables

On s'intéresse à présent aux colonnes \mathbf{z}^j du tableau \mathbf{Z} , représentant les variables.

Métrique D : On munit l'espace des variables d'une métrique naturelle, D , celle des poids des individus

$$\langle \mathbf{z}^i, \mathbf{z}^j \rangle_D = (\mathbf{z}^i)^\top \mathbf{D} \mathbf{z}^j \quad \|\mathbf{z}^i\|_D^2 = (\mathbf{z}^i)^\top \mathbf{D} \mathbf{z}^i$$

Interprétation : Pour des variables centrées, ce qui est le cas ici, on a

$$\text{cor}(\mathbf{z}^i, \mathbf{z}^j) = \frac{\langle \mathbf{z}^i, \mathbf{z}^j \rangle_D}{\|\mathbf{z}^i\|_D \|\mathbf{z}^j\|_D} = \cos(\mathbf{z}^i, \mathbf{z}^j)$$

Remarque : Les composantes principales normalisée $\mathbf{c}_k / \sqrt{\lambda_k}$ forment une base D -orthonormale :

$$\left\langle \frac{\mathbf{c}_i}{\sqrt{\lambda_i}}, \frac{\mathbf{c}_j}{\sqrt{\lambda_j}} \right\rangle_D = \text{cor} \left(\frac{\mathbf{c}_i}{\sqrt{\lambda_i}}, \frac{\mathbf{c}_j}{\sqrt{\lambda_j}} \right) = \delta_{ij}$$

ANALYSE EN COMPOSANTES PRINCIPALES

Corrélation des composantes principales et variables

- ▷ Le coefficient de corrélation entre \mathbf{c}_k et \mathbf{z}^i est donné par

$$\text{cor}(\mathbf{c}_k, \mathbf{z}^i) = \frac{\langle \mathbf{c}_k, \mathbf{z}^i \rangle_D}{\sqrt{\lambda_k}}$$

- ▷ Les coefficients de corrélation entre \mathbf{c}_k et l'ensemble des variables de \mathbf{Z} peuvent être calculés ainsi

$$\text{cor}(\mathbf{c}_k, \mathbf{Z}) = \frac{\mathbf{Z}^\top \mathbf{D} \mathbf{c}_k}{\sqrt{\lambda_k}}$$

- ▷ Comme $\mathbf{Z}^\top \mathbf{D} \mathbf{c}_k = \mathbf{Z}^\top \mathbf{D} \mathbf{Z} \mathbf{u}_k = \boldsymbol{\Sigma} \mathbf{u}_k$, on a directement

$$\text{cor}(\mathbf{c}_k, \mathbf{Z}) = \sqrt{\lambda_k} \mathbf{u}_k$$

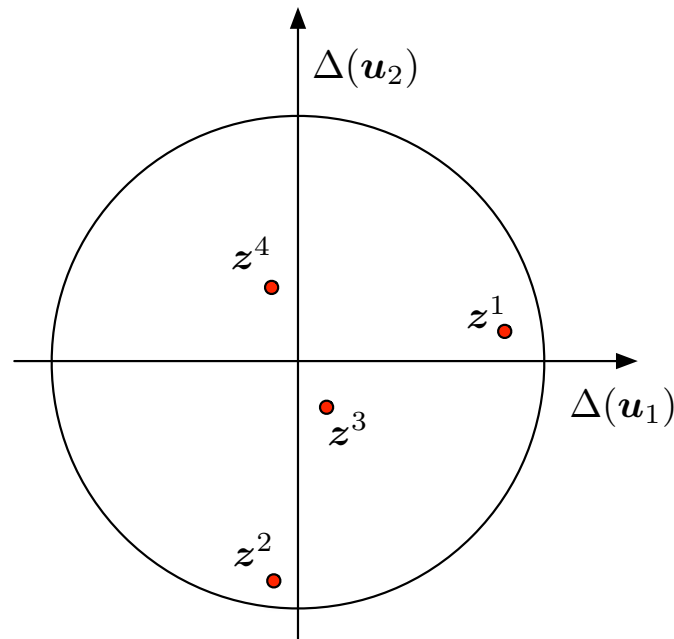
ANALYSE EN COMPOSANTES PRINCIPALES

Cercle des corrélations

On représente traditionnellement chaque variable z^i par un point de coordonnées

$$(\text{cor}(z^i, \mathbf{c}_1), \text{cor}(z^i, \mathbf{c}_2), \text{cor}(z^i, \mathbf{c}_3), \dots)$$

afin d'interpréter les axes principaux obtenus.



ANALYSE EN COMPOSANTES PRINCIPALES

En résumé

- ▷ Centrer et réduire les données $\mathbf{X} \rightarrow \mathbf{Z}$
- ▷ Calculer la matrice de covariance $\mathbf{\Sigma} = \mathbf{Z}^\top \mathbf{D} \mathbf{Z}$
- ▷ Diagonaliser $\mathbf{\Sigma}$ pour obtenir $\{(\mathbf{u}_k, \lambda_k)\}_{k=1, \dots, p}$
- ▷ Choisir le nombre k_0 d'axes à retenir
- ▷ Préciser le taux d'inertie expliqué par chacun des axes retenus
- ▷ Calculer les composantes principales $\{\mathbf{c}_k\}_{k=1, \dots, k_0}$
- ▷ Représenter les individus dans le(s) repère(s) des axes principaux
- ▷ Tracer les cercles de corrélation
- ▷ Interpréter