

FONDEMENTS DU MACHINE LEARNING

1^{er} octobre 2024

Aujourd'hui :

Cours (13^h45 - 15^h15) : ACP

TD (15^h30 - 18^h45) : Fin exercices SVD

TD/TP ACP

Analyse en composantes principales (ACP) pour la réduction de dimension

Problème: $X \in \mathbb{R}^{m \times n}$ matrice de données

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}, \quad x_i \in \mathbb{R}^n \text{ vecteur représentatif de l'individu } i, \text{ à } n \text{ attributs}$$

BUT: Réduire X à une matrice de taille $m \times k$ avec $k \leq m$ en conservant le maximum d'information pour chacun des individus x_1, \dots, x_m

\Rightarrow Analyse en k composantes principales

① Analyse en 1 composante principale

\hookrightarrow Objectif: Remplacer les n attributs par une nouvelle variable que l'on appelle composante principale.

\Rightarrow Géométriquement, cela revient à représenter le nuage de points $\{x_1, \dots, x_m\} \subseteq \mathbb{R}^n$ par $\{y_1, \dots, y_m\} \subseteq \mathbb{R}^k$ tels que les $\{y_i\}$ appartiennent à une même droite

Deux points de vue: Géométrique: On reste dans \mathbb{R}^n et on passe d'un nuage de points quelconque à un nuage de points sur une droite

Algébrique $X \in \mathbb{R}^{m \times n} \rightarrow c \in \mathbb{R}^{m \times 1}$

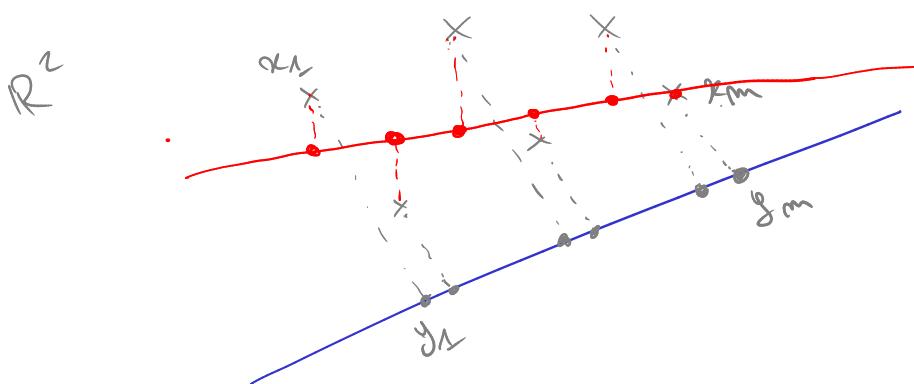
\Rightarrow La meilleure façon de réduire les n attributs en 1 constitue à déterminer la meilleure projection des $|x_i\rangle$ sur une droite dans \mathbb{R}^m

\Rightarrow La "meilleure droite possible" (au sens de fournir la représentation des données la plus proche des points de départ) est celle qui maximise l'inertie du nuage de points projetés

$X = \{x_1, \dots, x_m\}$: nuage de points de départ

On cherche $\mathcal{Y} = \{y_1, \dots, y_m\}$: nuage des points projetés, contenus dans une droite / séries.

$$(\text{Inertie de } \mathcal{Y}: \frac{1}{m-1} \sum_{i=1}^m \|y_i - \bar{y}\|^2 \quad \bar{y} = \frac{1}{m} \sum_{i=1}^m y_i \in \mathbb{R})$$



\hookrightarrow Comment calcule-t-on les y_i / la droite en pratique ?

- Une "bonne" droite devrait passer par l'individu moyen $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$ afin de conserver l'information moyenne du nuage de points

\Rightarrow Les coordonnées des projections du nuage de points ne prennent pas en compte la tendance centrale des données (mais les y_i , si !)

• les $\{y_i\}$ appartiennent à une droite d'équation

$$\{y = \bar{x} + \gamma u \mid \gamma \in \mathbb{R}\}$$

où u est le vecteur directeur de la droite à déterminer

On a alors : $y_i = \bar{x} + c_i u$ avec $c_i \in \mathbb{R}$

\uparrow
Composante principale de x_i

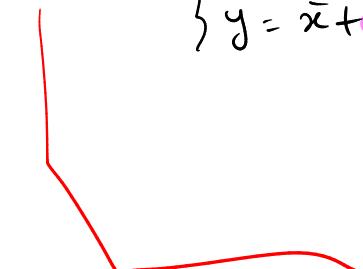
$$X \in \mathbb{R}^{m \times n} \rightarrow c \in \mathbb{R}^{m \times 1}, c = \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix} \text{ vecteur des premières composantes principales}$$

Théorème :

Pour tout $x \in \mathbb{R}^n$, la projection de x sur la droite

$$\{y = \bar{x} + \gamma u\}$$
 est donnée par

$$\bar{x} + \bar{u}^T(x - \bar{x}) u$$



Corollaire : Les projections $\{y_1, \dots, y_m\}$ du nuage de points

$$\{x_1, \dots, x_m\}$$
 sont données par

$$y_i = \bar{x} + \bar{u}^T(x_i - \bar{x}) u \quad \forall i=1..m$$

↪ On cherche le vecteur u tel que l'inertie du nuage de points projetés $\mathcal{Y} = \{y_1, \dots, y_m\}$ soit maximale

$$I(\mathcal{Y}) = \frac{1}{m-1} \sum_{i=1}^m \|y_i - \bar{y}\|^2$$

$$= \frac{1}{m-1} \sum_{i=1}^m \|(\bar{x} + \bar{u}^T(x_i - \bar{x}) u) - \bar{y}\|^2$$

$$\begin{aligned}
 \text{On } \bar{y} &= \frac{1}{m} \sum_{i=1}^m y_i \\
 &= \frac{1}{m} \sum_{i=1}^m \left(\bar{x} + u^\top (x_i - \bar{x}) u \right) \\
 &= \left(\frac{1}{m} \sum_{i=1}^m \bar{x} \right) + \left(\frac{1}{m} \sum_{i=1}^m u^\top (x_i - \bar{x}) u \right) \\
 &= \bar{x} + \frac{1}{m} u^\top \left(\sum_{i=1}^m (x_i - \bar{x}) \right) u
 \end{aligned}$$

$$\begin{aligned}
 &\sum_{i=1}^m u^\top (x_i - \bar{x}) u \\
 &= \left[\sum_{i=1}^m u^\top (x_i - \bar{x}) \right] u \\
 &= \left[u^\top \left(\sum_{i=1}^m x_i - \bar{x} \right) \right] u
 \end{aligned}$$

Done

$$\begin{aligned}
 I(\bar{y}) &= \frac{1}{m-1} \sum_{i=1}^m \| \bar{x} + u^\top (x_i - \bar{x}) u - \bar{y} \|^2 \\
 &= \frac{1}{m-1} \sum_{i=1}^m \| \bar{x} + u^\top (x_i - \bar{x}) u - \bar{x} - \frac{1}{m} u^\top \sum_{i=1}^m (x_i - \bar{x}) u \|^2 \\
 &= \frac{1}{m-1} \sum_{i=1}^m \| u^\top (x_i - \bar{x}) u - \frac{1}{m} u^\top \sum_{i=1}^m (x_i - \bar{x}) u \|^2 \\
 &= \frac{1}{m-1} \sum_{i=1}^m \| u^\top \left[(x_i - \bar{x}) - \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x}) \right] u \|^2
 \end{aligned}$$

$$x_i - \bar{x} - \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x}) = x_i - \bar{x} - \frac{1}{m} \sum_{i=1}^m x_i + \bar{x} = x_i - \frac{1}{m} \sum_{i=1}^m x_i = x_i - \bar{x}$$

$$\begin{aligned}
 I(\bar{y}) &= \frac{1}{m-1} \sum_{i=1}^m \| u^\top (x_i - \bar{x}) u \|^2 \quad \| v \|^2 = v^\top v \\
 &= \frac{1}{m-1} \sum_{i=1}^m u^\top (x_i - \bar{x})^\top u u^\top (x_i - \bar{x}) u \\
 &= \frac{1}{m-1} \sum_{i=1}^m [(x_i - \bar{x})^\top u]^2 u^\top u \\
 &= \frac{1}{m-1} \sum_{i=1}^m [u^\top (x_i - \bar{x})]^2 \| u \|^2
 \end{aligned}$$

$$= \frac{1}{m-1} \sum_{i=1}^m u^\top (x_i - \bar{x}) (x_i - \bar{x})^\top u \times \|u\|^2$$

$$(a^\top b)^2 : (a^\top b)(a^\top b) = (a^\top b)(b^\top a) = a^\top (b b^\top) a$$

$$\begin{aligned} I(u) &= \frac{1}{m-1} \|u\|^2 \sum_{i=1}^m u^\top (x_i - \bar{x}) (x_i - \bar{x})^\top u \\ &= \|u\|^2 u^\top \left(\frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^\top \right) u \\ &= \|u\|^2 u^\top \Sigma u \end{aligned}$$

\uparrow
 Σ : matrice de covariance
de $\{x_1, \dots, x_m\}$

→ Pour maximiser cette matrice par rapport à u , on a 2 cas :

- Soit $x_i = \bar{x} + t_i$ et dans ce cas $u = 0$
 $\hookrightarrow \Sigma = [0]$
- Sinon, on prend u de norme 1 et égal à un vecteur propre de Σ associé à la plus grande valeur propre de Σ

Formule à retenir

ACP avec 1 composante principale

$$\{x_1, \dots, x_m\} \subseteq \mathbb{R}^n \rightarrow \{y_1, \dots, y_m\} \subseteq \mathbb{R}^n$$

avec $y_i = \bar{x} + u^\top (x_i - \bar{x}) u$, u vecteur propre de norme 1 associé à la plus grande valeur propre de $\Sigma = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^\top$

$i = 1 \dots m$

c_i : (1^{re}) composante principale de x_i
 u : (1^{er}) vecteur principal
 $\{\bar{x} + y_u\}$: (1^{er}) axe principal

Toutes ces notions sont déterminées par le nuage de points X

En pratique, on calcule les points projetés via

$$Y = \begin{bmatrix} y_1^T \\ \vdots \\ y_m^T \end{bmatrix}_{m \times m} = \underbrace{\left(\bar{X} - \bar{1}_m \bar{x}^T \right) u u^T}_{m \times m} + \underbrace{\bar{1}_m \bar{x}^T}_{m \times 1}$$

(2) Analyse en k composantes principales, $k \geq 1$

Balr. $X \in \mathbb{R}^{m \times n} \rightarrow C = \begin{bmatrix} c_1^T \\ \vdots \\ c_m^T \end{bmatrix} \in \mathbb{R}^{m \times k}$

$$X = \{x_1, \dots, x_m\} \subseteq \mathbb{R}^n \rightarrow Y = \{y_1, \dots, y_m\} \subseteq \mathbb{R}^m$$

avec $\{y_i\} \subseteq \{\bar{x} + U_k \gamma | \gamma \in \mathbb{R}^k\}$

sous-espace affine de dimension k

et $U_k = [u_1 \dots u_k] \in \mathbb{R}^{m \times k}$ formé par k vecteurs linéairement indépendants et de norme 1

\Rightarrow On écrit: $y_i = \bar{x} + U_k c_i$

$c_i \in \mathbb{R}^k$ est le vecteur des composantes principales

↪ Le meilleur choix possible pour les U_k consiste à prendre u_1, \dots, u_k orthogonaux ($u_i^T u_j = 0 \text{ si } i \neq j$) et à les prendre comme k vecteurs propres de Σ associés aux k plus grandes valeurs propres.

NB: On exclut le cas $\Sigma = [0]$, pour lequel on peut prendre $U_k = [0]$

↪ On a alors $y_i = \bar{x} + U_k c_i$, avec $c_i = \begin{bmatrix} u_1^T(x_i - \bar{x}) \\ \vdots \\ u_k^T(x_i - \bar{x}) \end{bmatrix}$

c_i : vecteur des k premières composantes principales de x_i

U_k : k premiers vecteurs principaux

$\{\bar{x} + u_i \delta | \delta \in \mathbb{R}^k\}$: sous-espace associé aux k premiers axes principaux

$\{\bar{x} + y_{i,1} | \delta \in \mathbb{R}\}, \{\bar{x} + y_{i,2} | \delta \in \mathbb{R}\}, \dots, \{\bar{x} + y_{i,k} | \delta \in \mathbb{R}\}$

Propriété importante: Calculer les composantes principales peut se faire en calculant les composantes principales une à une (notamment parce que Σ possède une base orthonormée de vecteurs propres)

(3) Le cas $k=m$

$$X \in \mathbb{R}^{m \times n} \rightarrow C \in \mathbb{R}^{m \times m}$$

$$x_i \in \mathbb{R}^n \rightarrow y_i = \bar{x} + U_m c_i \\ = x_i$$

$$c_i \in \mathbb{R}^m$$

$[u_1, \dots, u_m]$ base
orthonormée de vecteurs
propres de Σ

c_i : coordonnées de $x_i - \bar{x}$ dans la base (e_1, \dots, e_m)

Avec $k=m$, on fait un changement de base (+ une translation par rapport à \bar{x}) relativement à notre jeu de données x_1, \dots, x_m

\Rightarrow Dans cette base, les coordonnées des $x_i - \bar{x}$ sont décomposées

| Représentation naturelle de X | Représentation avec l'ACP |
|---|--|
| Base canonique $\{e_1, \dots, e_m\}$ | Base de vecteurs propres de Σ |
| $X^c = X - \lambda_m \bar{x}^T$ (Données centrées) | $C^c = \underbrace{X^c}_{m \times m} \underbrace{U_m}_{m \times m}$ |
| $\Sigma = (\underbrace{X^c X^c}_{{m-1}})^T$ | $\Sigma' = \underbrace{(C^c)^T (C^c)}_{m-1} = U_m^T \Sigma U_m$ $= \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_m \end{bmatrix}$ où $\lambda_1 \geq \dots \geq \lambda_m$ sont les valeurs propres de Σ |

\Rightarrow Représentation pour laquelle les m attributs sont décomposés