

Master IA2S (Intelligence Artificielle, Science des données et Systèmes Cyber-Physiques)

Apprentissage Automatique

Ferhat ATTAL

ferhat.attal@u-pec.fr

Université Paris Est Créteil (UPEC)

Novembre 2021

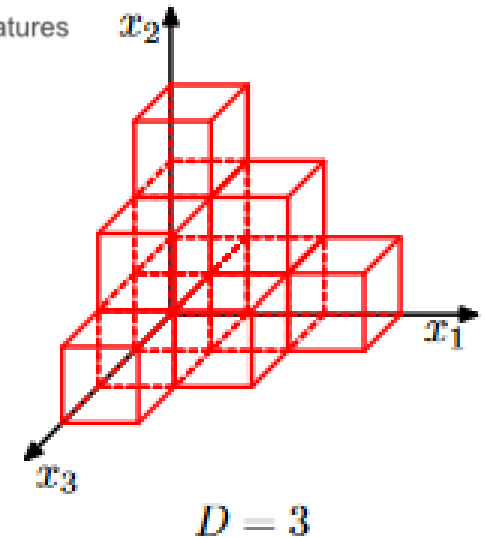
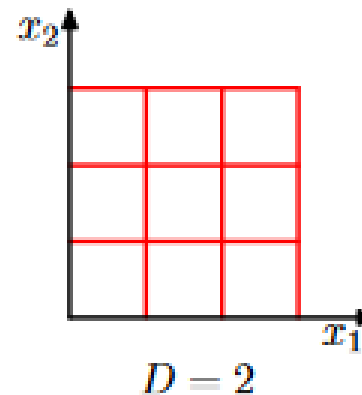
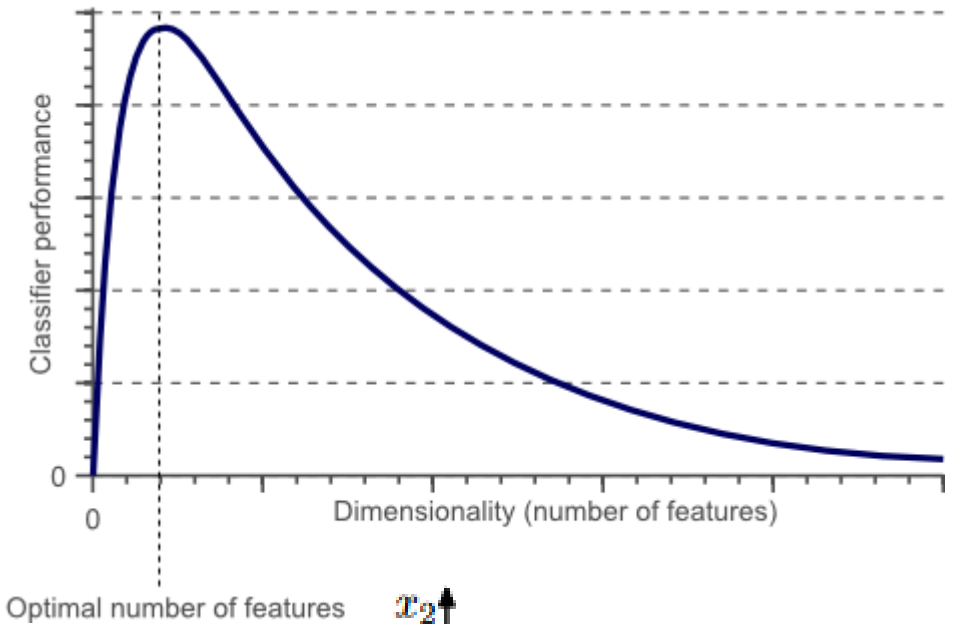
Cours 5 : Réduction de la dimension

Réduction de la dimension

- Réduction de la dimension
- Analyse en composantes principales
- Analyse factorielle

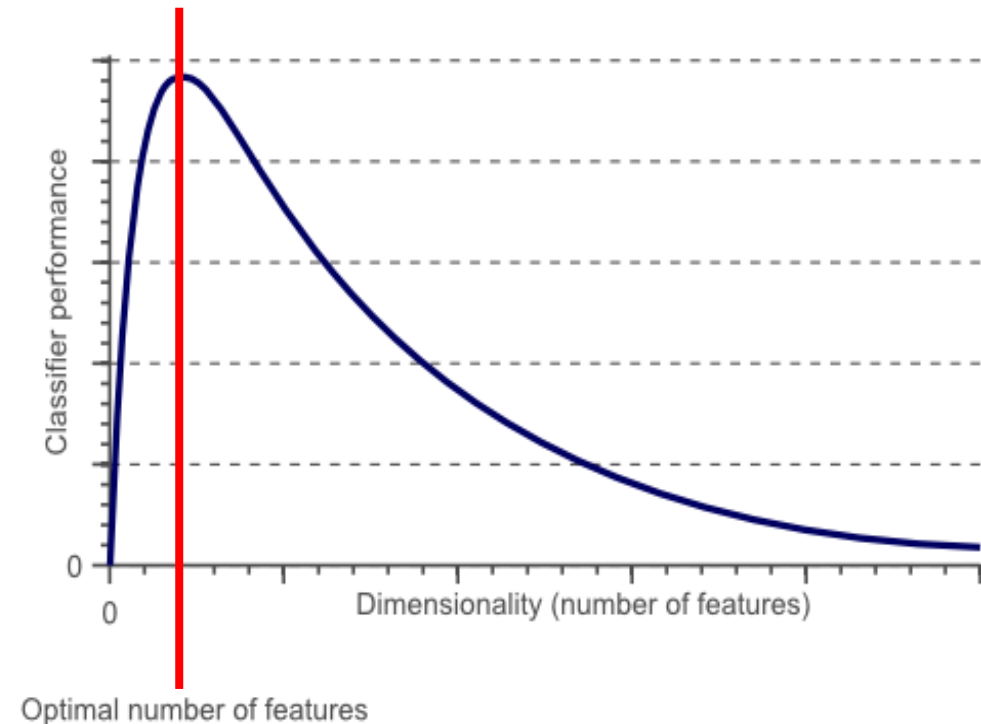
Malédiction de la dimensionnalité

- L'augmentation du nombre de variables (caractéristiques) n'améliorera pas toujours les performances de la classification.
- Le nombre d'exemples (individus) nécessaires pour l'apprentissage augmente de manière exponentielle avec la dimension d .



Réduction de la dimension

- Choisir un ensemble optimal de variables de dimension inférieure pour améliorer la capacité discriminative des algorithmes
- Réduire le temps de calcul des algorithmes
- Réduit la complexité de l'espace des variables
- faciliter l'interprétation des données
- Augmenter la capacité de généralisation des algorithmes
- Réduire la redondance



Réduction de la dimension

Extraction de caractéristiques (Feature extraction) : permet de trouver un ensemble de **nouvelles** caractéristiques (c'est-à-dire par le biais d'une fonction de projection $f()$) à partir des caractéristiques existantes.

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_j \\ \vdots \\ x_d \end{bmatrix} \xrightarrow{f()} z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_K \end{bmatrix}$$

$K \ll N$

Sélection de caractéristiques (Feature selection) : choisir un sous-ensemble de variables à partir de l'ensemble de variables d'originales.

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_j \\ \vdots \\ x_d \end{bmatrix} \rightarrow z = \begin{bmatrix} x_1 \\ x_5 \\ \vdots \\ x_K \end{bmatrix}$$

$K \ll N$

Réduction de la dimension

D'un point de vue mathématique, trouver une projection optimale $z = f(x)$ équivaut à optimiser une fonction objective.

Exemple de fonctions objectives:

- Minimiser la perte d'information : le but est de représenter les données aussi précisément que possible (c'est-à-dire sans perte d'information) dans l'espace de dimension inférieure.
- Augmenter la capacité discriminatoire : l'objectif est de trouver une projection permettant d'améliorer la capacité discriminatoires dans l'espace de dimension inférieure.

Analyse en composantes principales (ACP) transformation de Karhunen–Loève (KLT)

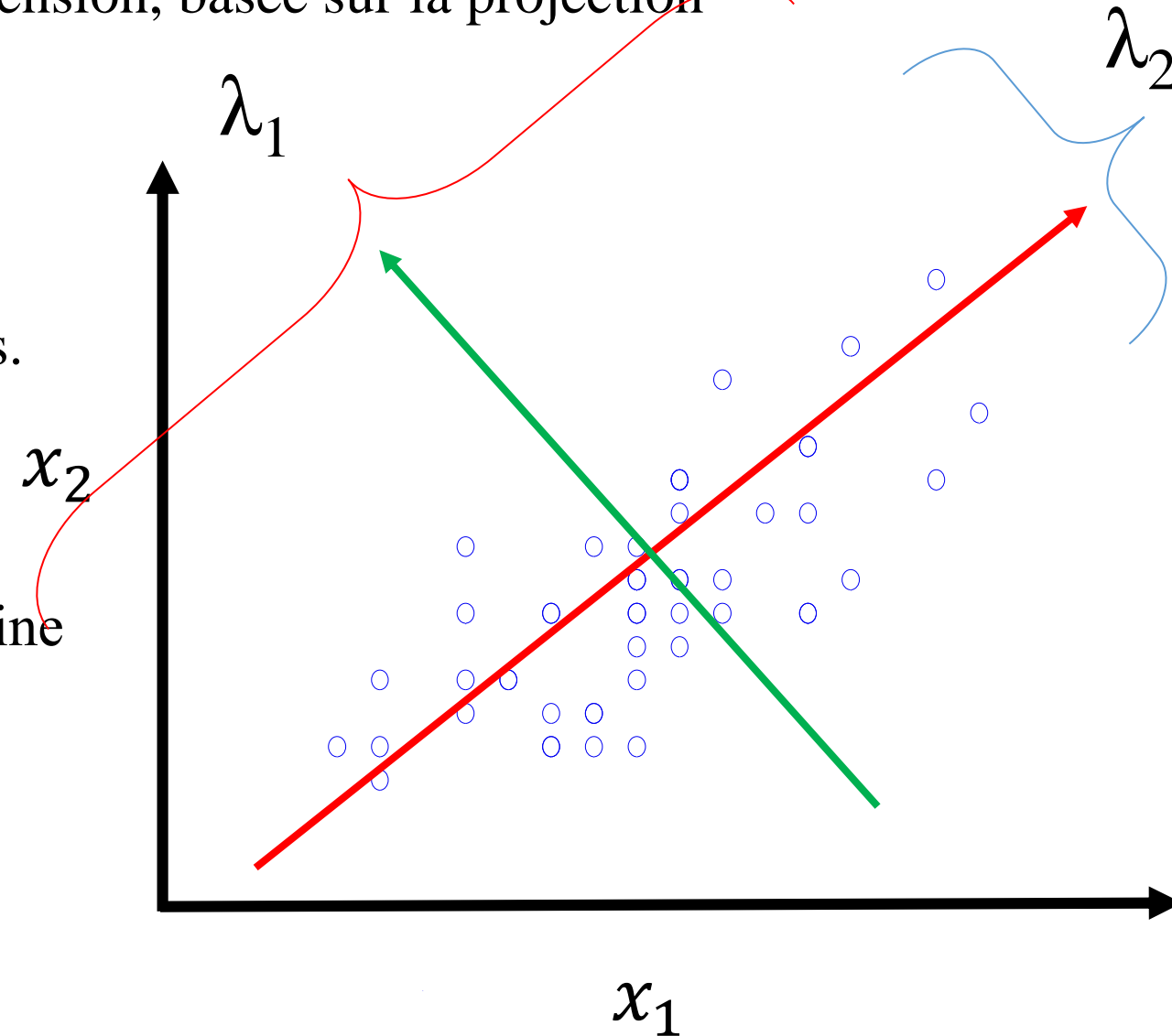
L'ACP est une technique de réduction de la dimension, basée sur la projection linéaire des données.

Soit en

- Maximisant la variance des données projetées.

Ou bien en

- Minimisant l'erreur entre les données d'origine et les données projetées



Analyse en composantes principales

Soit un échantillon $S = \{x_i\}_{i \leq n}$ où x_i appartient à \mathbb{R}^d .

	<i>1</i>	<i>...</i>	<i>j</i>	<i>...</i>	<i>d</i>
<i>1</i>	x_{11}	$...$	x_{1j}	$...$	x_{1d}
\vdots	\vdots		\vdots		\vdots
<i>i</i>	x_{i1}	$...$	x_{ij}	$...$	x_{id}
\vdots	\vdots		\vdots		\vdots
<i>n</i>	x_{n1}	$...$	x_{nj}	$...$	x_{nd}

L'objectif est de projeter les données sur un espace ayant une dimension $K < d$ tout en maximisant la variance des données projetées

Analyse en composantes principales

Soit le vecteur moyen $\bar{\mathbf{x}}$ et la matrice de covariance Σ

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

$$\Sigma = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Avec

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ij}, \dots, x_{id})^T$$

$$\bar{\mathbf{x}} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_j, \dots, \bar{x}_d)^T$$

Analyse en composantes principales

Considérons \mathbf{z}_i la projection de \mathbf{x}_i sur un espace unidimensionnel ($K=1$) muni d'un vecteur unitaire $\mathbf{u} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d)^T$

$$\mathbf{u}^T \mathbf{u} = 1$$

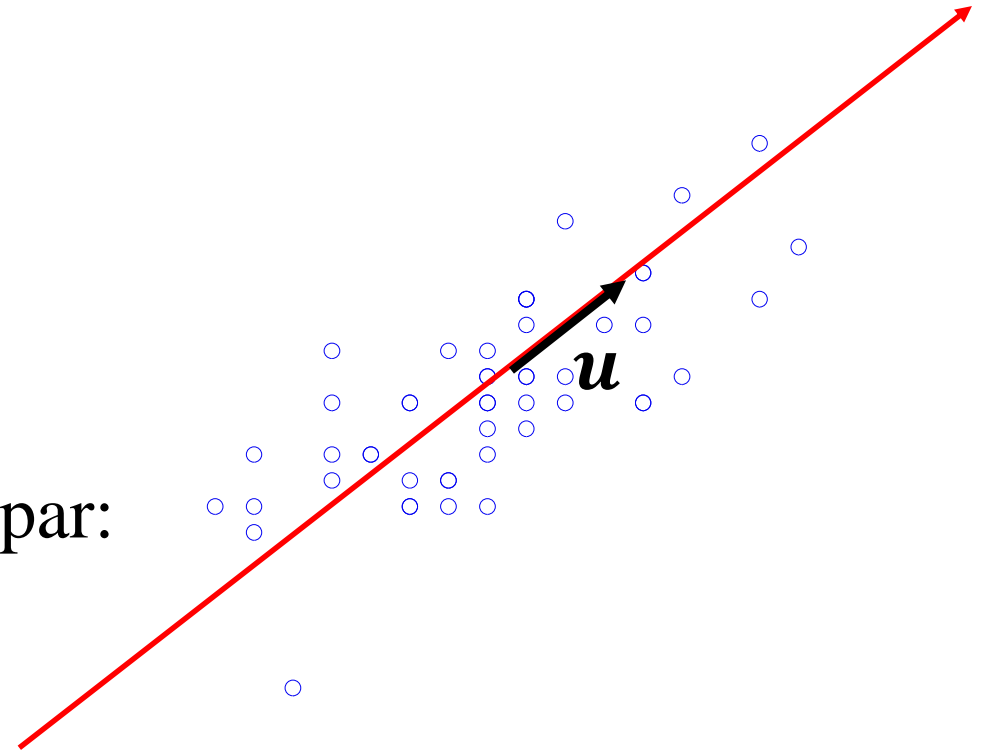
$$\mathbf{z}_i = \mathbf{u}^T \mathbf{x}_i$$

Ainsi

La variance des points projetés est donnée par:

$$\text{Var}(\mathbf{z}) = \text{var}(\mathbf{u}^T \mathbf{x})$$

$$= \mathbf{u}^T \Sigma \mathbf{u}$$



Analyse en composantes principales

Maximisation de la variance des données projetées

$$\begin{cases} \text{Maximiser } \mathbf{u}^T \Sigma \mathbf{u} \\ \text{s. c} \\ \mathbf{u}^T \mathbf{u} = 1 \end{cases}$$

Il s'agit d'un problème de maximisation sous contraintes \Rightarrow Utilisation du multiplicateur de Lagrange

Analyse en composantes principales

Maximisation de la variance des données projetées

$$\mathcal{L}(\mathbf{u}, \lambda) = \mathbf{u}^T \Sigma \mathbf{u} - \lambda (\mathbf{u}^T \mathbf{u} - 1)$$

$$\frac{\delta \mathcal{L}(\mathbf{u}, \lambda)}{\delta \mathbf{u}} = 0 \Rightarrow 2 \Sigma \mathbf{u} - 2 \lambda \mathbf{u} = \mathbf{0}$$

$$\Rightarrow \Sigma \mathbf{u} = \lambda \mathbf{u}$$

$\lambda \geq 0$ multiplicateur de Lagrange

Analyse en composantes principales

Maximisation de la variance des données projetées

$$\Sigma \mathbf{u} = \lambda \mathbf{u}$$

Avec

λ : valeurs propre de la matrice de covariance Σ

\mathbf{u} : vecteurs propre de la matrice de covariance Σ

On a

$$\begin{aligned}\text{Var}(\mathbf{z}) &= \mathbf{u}^T \Sigma \mathbf{u} \\ &= \mathbf{u}^T \lambda \mathbf{u} \\ &= \lambda \mathbf{u}^T \mathbf{u} \\ &= \lambda\end{aligned}$$

Analyse en composantes principales

La matrice de covariance admet d valeurs propres $(\lambda_1, \lambda_2, \dots, \lambda_d)$ et d vecteurs propres $(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d)$

Avec

$$\lambda_1 > \lambda_2 > \dots > \lambda_d$$

- Maximiser la variance des données projetées $\mathbf{u}^T \Sigma \mathbf{u}$ revient à trouver la valeur propre maximale λ c.à.d λ_1
- Le vecteur propre \mathbf{u}_1 correspondant à la valeur propre maximale λ_1 représente le vecteur directeur du premier axe principal (première composante principale).

Analyse en composantes principales

Nous pouvons définir des composantes principales supplémentaires de manière progressive en choisissant chaque nouvelle direction comme étant celle qui maximise la variance projetée parmi toutes les directions possibles orthogonales à celles déjà considérées.

Nous pouvons démontrer facilement que les vecteurs directeurs des composantes principales correspondent aux vecteurs propres $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$.

Remarque :

Σ est symétrique $\Rightarrow \mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ forment une base orthogonale

Analyse en composantes principales

L'espace de projection (l'espace des axes principaux) est alors obtenu en prenant les K premiers vecteurs propres $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K$

Comment choisir K

Choisir les K composantes permettant d'avoir une proportion de la variance expliquée supérieur à un certain seuil (0.9, 0.95 , etc.)

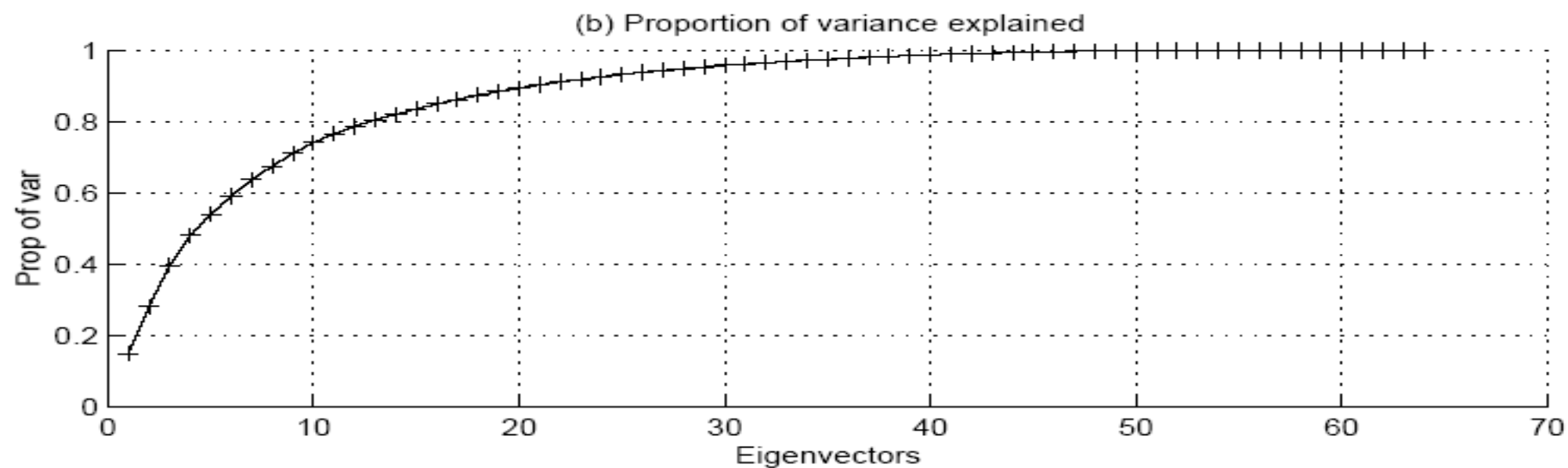
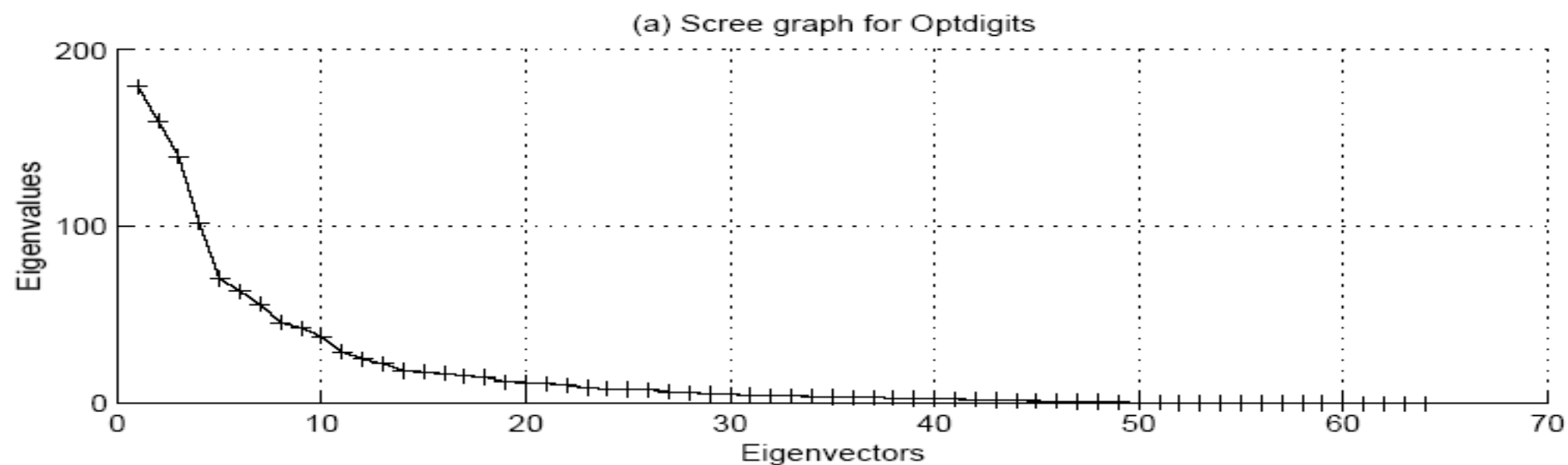
$$\text{Proportion de la variance expliquée} = \frac{\sum_{j=1}^K \lambda_j}{\sum_{j=1}^d \lambda_j}$$

Avec

Variance totale (inertie totale) des données projetés $\sum_{j=1}^d \lambda_j$

Variance des données projetées sur K axes principaux $\sum_{j=1}^K \lambda_j$

Analyse en composantes principales



La mise en œuvre de l'ACP peut être divisée en 6 étapes principales :

1. Calculer les centres des données

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

2. Calculer la matrice de covariance

$$\Sigma = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

3. Extraire les valeurs et vecteurs propres de Σ

4. Classer les vecteurs propres dans l'ordre décroissant des valeurs propres associées

$$(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d) \text{ Avec } \lambda_1 > \lambda_2 > \dots > \lambda_d$$

5. Choisir K composantes principales vérifiant $\frac{\sum_{j=1}^K \lambda_j}{\sum_{j=1}^d \lambda_j} \geq \text{seuil}$

6. Projeter les données sur le nouveaux espace $Z = \mathbf{U}^T X$

Avec $X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n)$ et $\mathbf{U} = (\mathbf{u}_1^T, \mathbf{u}_2^T, \dots, \mathbf{u}_K^T)$

Analyse factorielle

L'**analyse factorielle** est un outils très pertinent dans la mise en évidence des relations entre variables et, plus généralement, dans la compréhension des données à modéliser. Cette méthode vise à réduire la **dimension** des données (le nombre de variables) en conservant au mieux l'information utile.

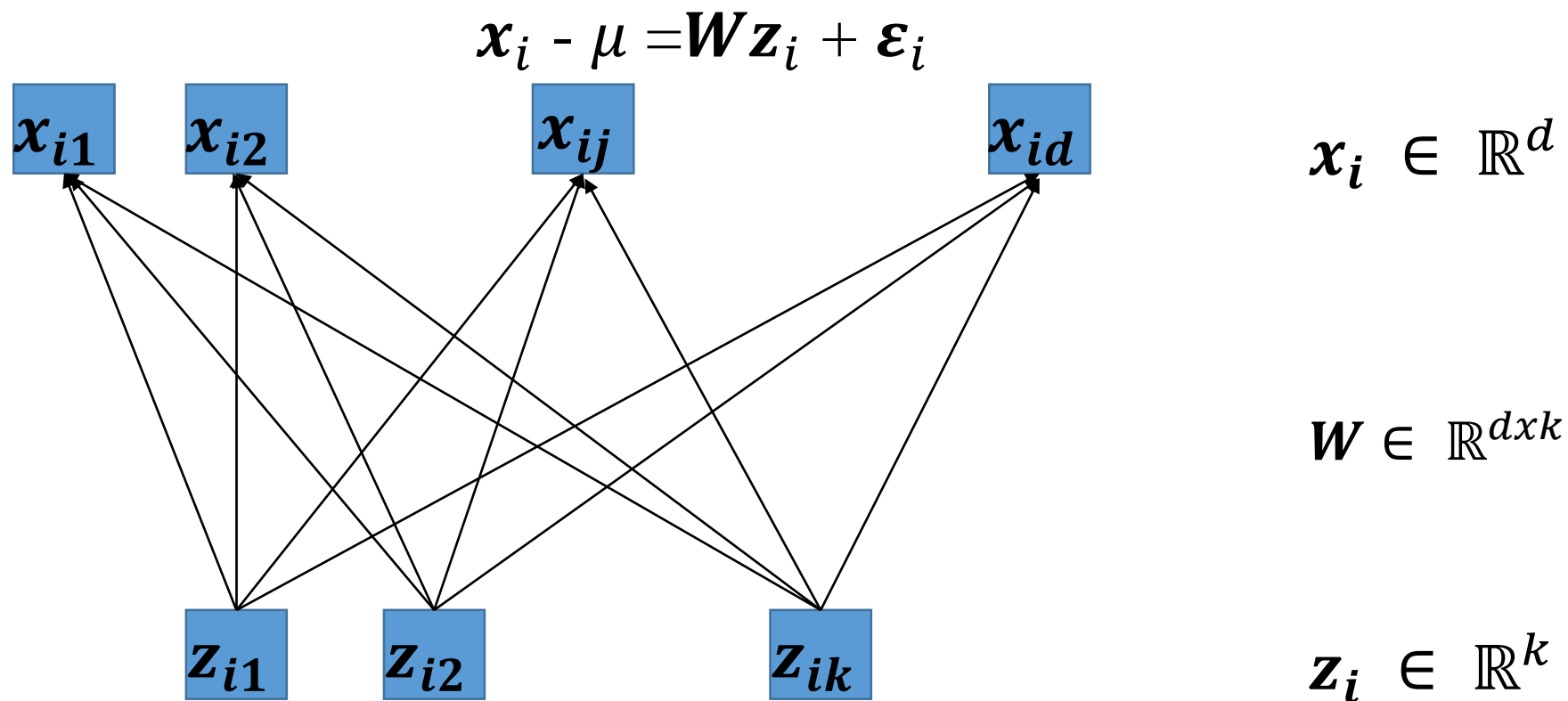
Objectif général : recherche de « facteurs » (variables dérivées) permettant de résumer les (caractéristiques) des données

- Améliorer la « lisibilité » des données
- Réduire le nombre de variables en conservant au mieux l'information utile

Analyse factorielle

- Soit $X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n)$ avec $\mathbf{x}_i \in \mathbb{R}^d$
- On fait l'hypothèse que ces variables dépendent (linéairement) en partie de k variables non observables, ou variables latentes ou facteurs $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$.

On cherche donc à décomposer les variables observées \mathbf{x}_i (supposées centrées) de la façon suivante



Analyse factorielle

Avec

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ij}, \dots, x_{id})^T$$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_j, \dots, \mu_d)^T$$

$$\mathbf{z}_i = (z_{i1}, z_{i2}, \dots, z_{ik})^T$$

$$\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \varepsilon_{i2}, \dots, \varepsilon_{ij}, \dots, \varepsilon_{id})^T$$

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_{11} & \mathbf{w}_{12} & \dots & \mathbf{w}_{1k} \\ \mathbf{w}_{21} & \mathbf{w}_{22} & \dots & \mathbf{w}_{2k} \\ & \vdots & & \\ \mathbf{w}_{j1} & \mathbf{w}_{j2} & \dots & \mathbf{w}_{jk} \\ & \vdots & & \\ \mathbf{w}_{d1} & \mathbf{w}_{d2} & \dots & \mathbf{w}_{dk} \end{bmatrix}$$

Analyse factorielle

$$\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\varepsilon}_i$$

$$\mathbf{z}_i \sim \mathcal{N}(0, \mathbf{I})$$

$$\boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, \boldsymbol{\psi})$$

$$\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\psi})$$

$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi})$$

$\boldsymbol{\psi}$ est une matrice diagonale

Analyse factorielle

Estimation des paramètre du modèle (W, ψ)

- Méthode basée sur l'ACP
- Méthode basée sur le maximum de vraisemblance
- Méthode basée sur les facteurs principaux

Estimation des paramètre du modèle (W, ψ)

➤ Méthode basée sur l'ACP

On sait que

$$\text{Var}(\mathbf{x}_i) = \Sigma = \mathbf{W}\mathbf{W}^T + \psi$$

$$\Sigma = \sum_{j=1}^d \lambda_j \mathbf{u}_j \mathbf{u}_j^T$$

$$= \underbrace{[\sqrt{\lambda_1} \mathbf{u}_1, \sqrt{\lambda_2} \mathbf{u}_2, \dots, \sqrt{\lambda_k} \mathbf{u}_k, \dots, \sqrt{\lambda_d} \mathbf{u}_d]}_{\mathbf{W}} \underbrace{\begin{bmatrix} \sqrt{\lambda_1} \mathbf{u}_1^T \\ \sqrt{\lambda_2} \mathbf{u}_2^T \\ \vdots \\ \sqrt{\lambda_k} \mathbf{u}_k^T \\ \vdots \\ \sqrt{\lambda_d} \mathbf{u}_d^T \end{bmatrix}}_{\mathbf{W}^T}$$

$\lambda_1, \lambda_2, \dots, \lambda_d$: valeurs propres de la matrice de covariance Σ

$\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$: vecteurs propres de la matrice de covariance Σ

Estimation des paramètre du modèle (W, ψ)

➤ Méthode basée sur l'ACP

$$\Sigma = WW^T + [\sqrt{\lambda_{k+1}}\mathbf{u}_{k+1}, \dots, \sqrt{\lambda_d}\mathbf{u}_d] \begin{bmatrix} \sqrt{\lambda_{k+1}}\mathbf{u}_{k+1} \\ \vdots \\ \sqrt{\lambda_d}\mathbf{u}_d^T \end{bmatrix}$$

ainsi

$$\psi = \text{diag}\{[\sqrt{\lambda_{k+1}}\mathbf{u}_{k+1}, \dots, \sqrt{\lambda_d}\mathbf{u}_d] \begin{bmatrix} \sqrt{\lambda_{k+1}}\mathbf{u}_{k+1} \\ \vdots \\ \sqrt{\lambda_d}\mathbf{u}_d^T \end{bmatrix}\}$$

Ou bien

$$\psi = \text{diag}\{\Sigma - WW^T\}$$

Estimation des paramètre du modèle (W, ψ)

➤ Méthode basée sur le maximum de vraisemblance

$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi})$$

Log- vraisemblance de \mathbf{x}

$$\mathcal{L}(\theta; \mathbf{x}) = \log \prod_{i=1}^n p(\mathbf{x}_i | \theta)$$

$$\mathcal{L}(\theta; \mathbf{x}) = \log \prod_{i=1}^n \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi})$$

$$\mathcal{L}(\theta; \mathbf{x}) = \frac{nd}{2} \log(2\pi) - \frac{n}{2} \log |\mathbf{W}\mathbf{W}^T + \boldsymbol{\psi}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T (\mathbf{W}\mathbf{W}^T + \boldsymbol{\psi})^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

Estimation des paramètre du modèle (W, ψ)

➤ Méthode basée sur le maximum de vraisemblance

Log- vraisemblance des données complétées

$$\mathcal{L}_c(\theta; \mathbf{x}, \mathbf{z}) = \log \prod_{i=1}^n p(x_i | z_i; \theta) p(z_i; \theta)$$

$$= \log \prod_{i=1}^n \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\psi}) \mathcal{N}(0, \mathbf{I})$$

$$= \sum_{i=1}^n \log \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\psi}) + \log \mathcal{N}(0, \mathbf{I})$$

Estimation des paramètre du modèle (W, ψ)

➤ Méthode basée sur le maximum de vraisemblance

Algorithme EM

Etape E

$$E[\mathbf{z}_i] = (\mathbf{I} + \mathbf{W}\psi^{-1}\mathbf{W}^T)^{-1}\mathbf{W}^T(\psi^{-1}(\mathbf{x}_i - \boldsymbol{\mu}))$$

$$E[\mathbf{z}_i\mathbf{z}_i^T] = (\mathbf{I} + \mathbf{W}\psi^{-1}\mathbf{W}^T)^{-1} + E[\mathbf{z}_i]E[\mathbf{z}_i]^T$$

Etape M

$$\mathbf{W}_q = \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})E[\mathbf{z}_i]^T \sum_{i=1}^n E[\mathbf{z}_i\mathbf{z}_i^T]$$

$$\psi_q = \text{diag} \left\{ \Sigma - \mathbf{W}_q \frac{1}{n} \sum_{i=1}^n E[\mathbf{z}_i](\mathbf{x}_i - \boldsymbol{\mu})^T \right\}$$