

RÉDUCTION DE DIMENSION

Idee: Passer d'une représentation en p variables à une représentation en d variables
 $d \ll p$

[1] Motivations

– Visualisation (en particulier $d=2$)

– Utiliser moins de ressources

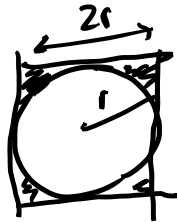
espace de stockage
acquisition des données
tps de calcul des modèles

– Fléau de la dimension
(curse of dimensionality)

Apprentissage supervisé marche
mieux avec peu de variables

En grande dimension, toutes les
distances sont grandes

$d=2$



Proportion du carré en dehors du cercle = $1 - \frac{\pi}{4}$

$d=3$

Proportion du cube en dehors de la sphère = $1 - \pi/6$

$d \rightarrow +\infty$

proportion $\rightarrow 1$

Les intuitions développées en $d=2$ ou $d=3$
ne sont pas toutes valides en grande dimension.

[2] Sélection vs extraction de variables

↓
Éliminer certaines
des variables

↳ Créer de nouvelles variables
- combiner les variables existantes en de nouvelles variables.

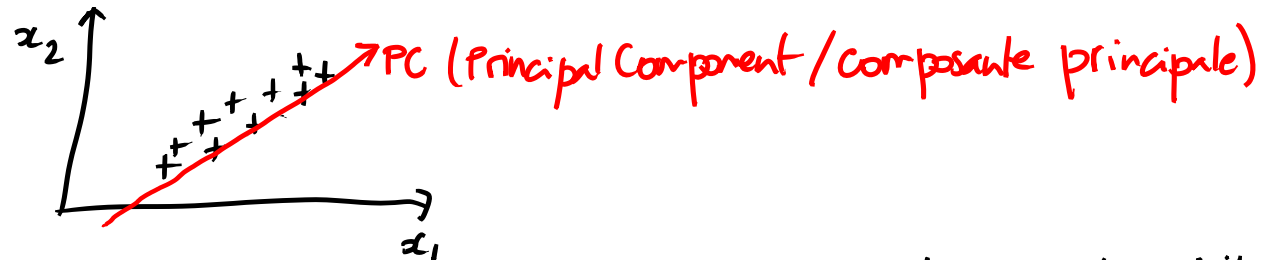
→ ne garder qu'une variable parmi un ensemble de variables corrélées

→ Lasso: éliminer les variables qui ne sont pas utilisées.

[3] Analyse en composantes principales (ACP / PCA)

$$\mathcal{D} = \{ \vec{x}_1, \dots, \vec{x}_n \} \quad \vec{x}_i \in \mathbb{R}^p$$

Idée: trouver un repère orthonormé sur lequel projeter \mathcal{D}
de sorte à maximiser la variance des données projetées



- les nouvelles variables (PCs) sont des combinaisons linéaires des variables initiales
- On connaît facilement la proportion de variance expliquée par chaque composante
→ utile pour choisir le nombre de composantes.

[4] Positionnement multidimensionnel / MDS = Multidimensional Scaling

n observations $\vec{x}_1, \dots, \vec{x}_n \rightarrow$ distances deux à deux $D_{i,l} = d(\vec{x}_i, \vec{x}_l)$

Je cherche $\vec{z}_1, \dots, \vec{z}_n$ $\vec{z}_i \in \mathbb{R}^d$ $d \ll p$

tels que $\min \sum_{i,l} \left(\|\vec{z}_i - \vec{z}_l\|_2 - D_{i,l} \right)^2$ dissimilité

Si $d(\vec{x}_i, \vec{x}_l) = \|\vec{x}_i - \vec{x}_l\|_2$ alors on retrouve l'ACP

[5] kPCA = kernel PCA

version à noyau de l'ACP \Rightarrow Les nouvelles variables sont des fonctions non-linéaires des variables initiales.

[6] t-SNE t-Student Neighborhood Embedding

idée: modéliser la distribution des distances entre observations par une loi t de Student.

très utilisé pour la visualisation

! ne préserve que les structures locales

[7] UMAP