

## Recap PCA

Objectif: Passer d'un espace de dimension  $p$  des features à un espace de dimension  $k < p$ .

# Cela peut-être à des fins de visualisation ( $k=2$  ou  $3$ ), computationnelles (complexité) ou encore pour éviter l'overfitting (mais ce n'est pas "une bonne pratique": mieux vaut régulariser).

+ utilisation annexe importante: avoir des nouvelles features décorrélées.

Principe: On va projeter nos données  $x_1, \dots, x_n \in \mathbb{R}^p$  dans un espace de dimension  $k < p$  telle que la "variance totale" résultante est maximisée.

On écrit  $X \in \mathbb{R}^{n \times p}$ ,  $n$  le nombre de données,  $p$  le nombre features:

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} = (x_1 | \dots | x_p)$$

↳ En pratique: on centre et réduit nos données au préalable.

$$x \rightarrow \frac{x - \bar{x}}{s}$$

## Rappel de Maths préliminaires

I/ SVD: Singular Value Decomposition

• Pour  $X \in \mathbb{R}^{n \times p}$ , on peut écrire

$$X = UDV^T \quad U \in O_n(\mathbb{R}) \text{ (i.e. orthogonale)},$$

$$V \in O_p(\mathbb{R}), \quad D = \begin{pmatrix} D_{11} & & 0 \\ & \ddots & \\ 0 & & D_{mm} \end{pmatrix}, \quad m = \min(n, p)$$

$$\text{et } D_{11} \geq \dots \geq D_{mm} \geq 0.$$

⊕ Si  $n > p$ : on peut réécrire cela comme

$$X = UDV^T, \quad U \in \mathbb{R}^{n \times p}, \quad D \in \mathbb{R}^{p \times p}, \quad V \in O_p(\mathbb{R})$$

$$\text{où } U = (U_1 | \dots | U_n) \xrightarrow{\text{devient}} (U_1 | \dots | U_p).$$

(On supprime les dernières colonnes de  $U$ ) et  $D$  diagonale

En effet:

$$\begin{aligned} UDV^T &= (U_1 | \dots | U_n) \begin{pmatrix} D_{11} & & 0 \\ & \ddots & \\ 0 & & D_{pp} \end{pmatrix} \begin{pmatrix} V_1^T \\ \vdots \\ V_p^T \end{pmatrix} \\ &= (U_1 | \dots | U_n) \begin{pmatrix} D_{11} V_1^T \\ \vdots \\ D_{pp} V_p^T \\ 0 \\ \vdots \end{pmatrix} = \sum_{i=1}^p U_i D_{ii} V_i^T \end{aligned}$$

# on peut donc supprimer les dernières colonnes

⊕ Si  $n < p$ : On réécrit  $X = UDV^T$ ,  
 $V^T = \begin{pmatrix} V_1^T \\ \vdots \\ V_n^T \end{pmatrix}$ ,  $U \in O_n(\mathbb{R})$ ,  $D \in \mathbb{R}^{n \times n}$ .

↳ Cette réécriture s'appelle thin SVD (en anglais)

On a alors:  $X^T X = (UDV^T)^T (UDV^T)$

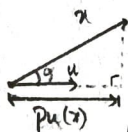
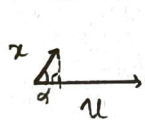
$$X^T X = V D^T D V^T$$

$$\boxed{X^T X = V \Lambda V^T}, \Lambda = D^T D, \\ V \in O_p(\mathbb{R})$$

$\Lambda \rightarrow$  contient les valeurs propres de  $X^T X$  (valeurs singulières) qui sont réelles  $\geq 0$ .

$V \rightarrow$  contient les valeurs propres qui forment une base orthonormée de  $\mathbb{R}^p (v_1, \dots, v_p)$ .

## II / Produit scalaire et projection en 2D



$$\cos(\alpha) = \frac{pu(x)}{\|x\|_2}$$

Projection orthogonale de  $x$  sur la droite avec vecteur directeur  $u$

$$\text{Donc } \langle x, u \rangle = \|x\|_2 \cdot \|u\|_2 \cdot \cos(\alpha) = \|u\|_2 \cdot pu(x)$$

$$\text{Si } \|u\|_2 = 1 \Rightarrow \langle x, u \rangle = \underline{pu(x)}$$

## III / Caractérisation projection orthogonale sur un plan dans $\mathbb{R}^n, n \geq 3$

Rappel:  $\Pi_P(x) = \underset{y \in P}{\operatorname{argmin}} \|y - x\|_2$  pour  $P$  un plan.  
 $\swarrow$   
 projection orthogonale de  $x$  sur  $P$ .

Prenons une base orthonormée de  $P (v_1, v_2)$ .

$$\forall y \in P, y = \lambda_1 v_1 + \lambda_2 v_2$$

$$\|x - y\|_2^2 = \|x\|_2^2 - 2\langle x, y \rangle + \|y\|_2^2.$$

$$\boxed{\|x - y\|_2^2 = \|x\|_2^2 + \lambda_1^2 + \lambda_2^2 - 2\lambda_1 \langle x, v_1 \rangle - 2\lambda_2 \langle x, v_2 \rangle}$$

$$\text{Posons } \varphi(\lambda_1, \lambda_2) = -2\lambda_1 \langle x, v_1 \rangle - 2\lambda_2 \langle x, v_2 \rangle + \lambda_1^2 + \lambda_2^2$$

$$\begin{cases} \frac{\partial \varphi}{\partial \lambda_i} = -2\langle x, v_i \rangle + 2\lambda_i, i=1,2 \\ \frac{\partial^2 \varphi}{\partial \lambda_i^2} = 2, i=1,2 \end{cases}$$

$$\frac{\partial^2 \varphi}{\partial \lambda_1 \partial \lambda_2} = 0$$

$\hookrightarrow \frac{\partial^2 \varphi}{\partial x^2} > 0$  positive définie (Jacobiennes)

donc la fonction est convexe (strictement)

$\Rightarrow$  Minimum atteint pour  $\frac{\partial \varphi}{\partial x} = 0$  ie

$$\boxed{\lambda_i = \langle x, v_i \rangle, i=1,2}$$

$$\hookrightarrow \boxed{y = \langle x, v_1 \rangle v_1 + \langle x, v_2 \rangle v_2 = \Pi_P(x)}$$

On généralise facilement à un espace de dimension  $k$  en prenant  $(v_1, \dots, v_k)$  une base orthonormée de  $E_k$ :

$$\text{On obtient } \Pi_{E_k}(x) = \sum_{i=1}^k \langle x, v_i \rangle v_i$$

$$\Pi_{E_k}(x) = \sum_{i=1}^k (x^T v_i) v_i$$



# PCA - Introduction

La "variance totale" est définie

comme :  $\frac{1}{n} \sum_{i=1}^n \|x_i - \bar{x}\|_2^2 \Leftrightarrow \underbrace{\frac{1}{n} \sum_{i=1}^n \|x_i\|_2^2}_{\text{données centrées}} - \|\bar{x}\|_2^2$

(le facteur  $\frac{1}{n}$  n'est pas important pour la suite).

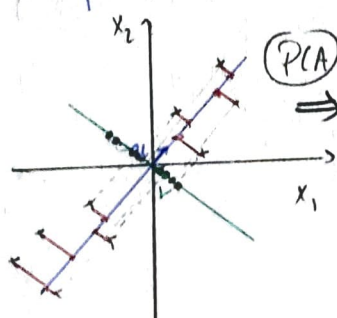
On cherche donc un espace de dimension  $k$  tq la projection des  $x_i$  sur cet espace donne des données (et nouvelles features anciennes) qui gardent le maximum de variance totale possible.

ie : on aura une nouvelle matrice

$$Y = \begin{pmatrix} y_1^T \\ \vdots \\ y_n^T \end{pmatrix}, y_i \in \mathbb{R}^k \text{ telle que}$$

$\sum \|y_i\|_2^2$  est "maximale" dans le sens défini ci-dessus.

Exemple :  $2D \rightarrow 1D$  (ou  $PD \rightarrow 1D$ )



Nouvelle feature  $z$

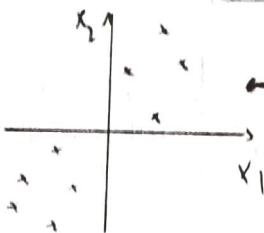
On veut garder une "dispersion" maximale

Bon choix de projection

Nouvelle feature  $z$

\* Mauvais choix de projection

Les nouvelles données sont très rapprochées



Données originales

Bon choix de projection

Mauvais choix

Prendre le cas où on a des données avec 2 features qu'on veut "compresser" en une seule.

On cherche  $\tilde{x}$  ( $\|u\|_2=1$ ) tq les données projetées sur  $\tilde{x}$  ( $x_i^T \tilde{x}$ ,  $i=1, \dots, n$ ) gardent une variance maximale.

ie  $x_u = \begin{pmatrix} x_1^T u \\ \vdots \\ x_n^T u \end{pmatrix}$  et on a vu que

$x_i^T u = p_u(x_i)$  pour  $u$  de norme 1.

On cherche  $u \in \mathbb{R}^2$  tq :  $u = \arg \max_{\substack{u \in \mathbb{R}^2 \\ \|u\|_2=1}} \|x_u\|_2^2$

ou, plus généralement :  $u = \arg \max_{\substack{u \in \mathbb{R}^p \\ \|u\|_2=1}} \|x_u\|_2^2$

Soit part de  $\mathbb{R}^p$ ,  $p \geq 2$ .

On a :  $\|x_u\|_2^2 = u^T X^T X u = u^T V D^T D V u$  (avec  $\|u\|_2=1$ )

En posant  $a = V^T u$ ,  $\|a\|_2 = \|u\|_2 = 1$ ,  $\|x_u\|_2^2 = a^T D^T D a$

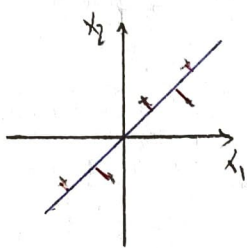
$$\|x_u\|_2^2 = \sum_{i=1}^m a_i^2 D_{ii}^2 \leq D_{11}^2 \sum_{i=1}^m a_i^2 = D_{11}^2$$

Donc  $a = \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix}$  maximise  $\|x_u\|_2^2$  ie  $\tilde{u} = V a = V_1$

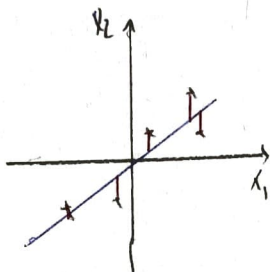
La direction qui maximise  $\|x_u\|_2^2$  est  $V_1$ , valeur propre associée à  $D_{11}^2$  :  $\|x_{V_1}\|_2^2 = D_{11}^2 \|V_1\|_2^2 = D_{11}^2$

Remarque: les données projetées  $x_i \cdot v$  sont centrées ( $\sum x_i \cdot v = \sum (x_i \cdot v) = 0$ ), donc la variance totale est bien  $\sum_{i=1}^n (x_i \cdot v)^2 = \|Xv\|_2^2$  (#facteur  $\frac{1}{n}$  pas important pour trouver le maximum)

• Ne pas confondre PCA et régression linéaire.



PCA



Régression linéaire

↳ Les "traits rouge" (ce qu'on cherche à minimiser) sont différents.

## Formulation du problème

Notre problème peut se résumer à :

$$\arg \max_{(v_1, \dots, v_k) \in (\mathbb{R}^p)^k} \left\{ \sum_{i=1}^k \|Xv_i\|_2^2 \right\} \text{ tels que } v_i \cdot v_j = \delta_{i,j}$$

En effet, on a, pour un espace  $E_n$  de dimension  $k$  et une base orthonormée  $(v_1, \dots, v_k)$  :

$$\Pi_{E_n}(x_i) = \sum_{l=1}^k (x_i \cdot v_l) \cdot v_l \quad (\text{projection orthogonale sur } E_n)$$

N.B.  $\sum_{i=1}^n \|\Pi_{E_n}(x_i)\|_2^2 = \sum_{i=1}^n \sum_{l=1}^k (x_i \cdot v_l)^2$

Dans ce nouvel espace, on peut récrire les coordonnées des  $x_i$  projetées dans la base orthonormée :

$$Y = \begin{pmatrix} y_{11} \\ \vdots \\ y_{n1} \end{pmatrix} = \begin{pmatrix} x_{11} \cdot v_1 & x_{11} \cdot v_2 & \dots & x_{11} \cdot v_k \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} \cdot v_1 & x_{n1} \cdot v_2 & \dots & x_{n1} \cdot v_k \end{pmatrix} = (Xv_1 \dots Xv_k) = (y_{11} \dots y_{1k})$$

La variance totale est :

$$\hookrightarrow \sum_{i=1}^n \sum_{l=1}^k x_{il}^2 = \sum_{l=1}^k \|Xv_l\|_2^2$$

## Résolution

On peut commencer par  $v_1$  pour la 1<sup>ère</sup> feature puis chercher  $u$  tq  $\|Xu\|_2^2$  est maximale et  $\|u\|_2 = 1$  et  $u^T v_1 = 0$  (observable).

On peut montrer facilement que  $u = v_2$ .  
Récursivement, on obtient  $(v_1, \dots, v_k)$  les vecteurs propres de  $X^T X$  associés aux valeurs propres  $D_1^2 \geq D_2^2 \geq \dots \geq D_k^2 \geq 0$ .

on a alors :  $Y = (y_{11} \dots y_{1k}) = (Xv_1 \dots Xv_k)$   
 $Y = (D_1 v_1 \dots D_k v_k)$

↳ la variance totale est alors :  $\sum_{l=1}^k \|D_l v_l\|_2^2 = \sum_{l=1}^k D_l^2$

A noter que la variance totale avant projection est :  $\sum_{l=1}^m D_l^2$

Une façon de choisir  $k$  (dimension de l'espace sur lequel on projette) :

$$\min \left\{ k \mid \frac{\sum_{l=1}^k D_l^2}{\sum_{l=1}^m D_l^2} \geq 0.95 \right\}$$

↳ soit autre : 0.99, 0.90, 0.80, ...

Remarques :

- les données sont centrées

$$XV_1 = \begin{pmatrix} x_{11}TV_1 \\ \vdots \\ x_{n1}TV_1 \end{pmatrix} \rightarrow (\sum x_{i1}TV) = (\sum x_{i1})V = 0$$

- les features créées sont décorréllées:

$$i \neq j, (XV_i)^T (XV_j) = V_i^T X^T X V_j = V_i^T V \Lambda V^T V_j \\ \underline{= 0}$$

⇒ D'ailleurs, on utilise cette méthode pour "décorréler des features" (ie avoir de nouvelles features décorréllées)

- les features créées perdent en interprétabilité/explicabilité

Idée générale: garder le plus de "dispersion" possible entre les données.