

# Références / page pub

---

- **Introduction au Machine Learning**, Chloé-Agathe Azencott, Dunod InfoSup  
Version électronique (PDF) sans exercices disponible gratuitement sur <http://cazencott.info> [lien cliquable]  
(nouvelle édition en préparation)
- **Le Machine Learning avec scikit-learn**, Aurélien Géron, Dunod InfoSup
- Textes et vidéos sur **OpenClassrooms** :  
Textes accessibles librement ; vidéos accessibles avec un compte gratuit.
  - [Parcours Ingénieur Machine Learning](#) [lien cliquable]
  - [Objectif IA \(grand public\)](#) [lien cliquable]

# Ressources complémentaires

---

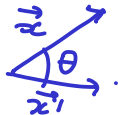
- **Sources de jeux de données**
  - Kaggle [\[lien cliquable\]](#)
  - Le module `datasets` de `scikit-learn` [\[lien cliquable\]](#)
  - Le UCI Machine Learning Repository [\[lien cliquable\]](#)
- **MOOC scikit-learn** [\[lien cliquable\]](#)

# Cours 4 – Méthodes à noyaux

# 1. Noyaux

# 1.1 Produit scalaire

$$- \langle \vec{x}, \vec{x}' \rangle = \sum_{j=1}^p x_j x'_j \quad \vec{x}, \vec{x}' \in \mathbb{R}^p$$



$$- \langle \vec{x}, \vec{x}' \rangle = \|\vec{x}\|_2 \|\vec{x}'\|_2 \cos \theta \quad \|\vec{x}\|_2^2 = \langle \vec{x}, \vec{x} \rangle$$

- Interprétable comme **mesure de similarité** \*  $\theta = \frac{\pi}{2}$ ,  $\cos \theta = 0$ ,  
vecteurs  $\perp$

- **Forme définie positive bilinéaire :**

$$- \langle \vec{x}, \vec{x}' \rangle = \langle \vec{x}', \vec{x} \rangle \text{ pour tout } \vec{x}, \vec{x}' \in \mathcal{X}$$

\*  $\theta = 0$ ,  $\cos \theta = 1$ , vecteurs //

$$- \langle \vec{x} + \vec{z}, \vec{x}' \rangle = \langle \vec{x}, \vec{x}' \rangle + \langle \vec{z}, \vec{x}' \rangle \text{ pour tout } \vec{x}, \vec{x}', \vec{z} \in \mathcal{X}$$

$$- \langle a\vec{x}, \vec{x}' \rangle = a \langle \vec{x}, \vec{x}' \rangle \text{ pour tout } a \in \mathbb{R}$$

$$- \langle \vec{x}, \vec{x} \rangle \geq 0 \text{ et } \langle \vec{x}, \vec{x} \rangle = 0 \text{ ssi } \vec{x} = 0$$

- Apparaît dans de nombreux algorithmes d'apprentissage.

## 1.2 Noyau

$\mathcal{X} = \mathbb{R}^p$  ou un espace  $\oplus$  compliqué (images, textes ...)

- **Généralisation** du produit scalaire :  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- **sémantiquement** une similarité
- **mathématiquement** une forme définie positive
- Aronszajn: Si  $k$  est semi-définie positive<sup>1</sup>, il existe un espace de Hilbert  $\mathcal{H}$  et une application  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$  telle que  $k(\vec{x}, \vec{x}') = \langle \varphi(\vec{x}), \varphi(\vec{x}') \rangle_{\mathcal{H}}$  pour tout  $\vec{x}, \vec{x}', \vec{z} \in \mathcal{X}$  (on ne connaît pas nécessairement  $\varphi$ ).

Ex. régression quadratique  $\varphi : (x_1, \dots, x_p) \mapsto (x_1, \dots, x_p, x_1^2, x_1 x_2, \dots, x_p^2)$

---

<sup>1</sup>Pour tout  $m \in \mathbb{N}^*$ , pour tout  $\{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_m\} \in \mathcal{X}$ , la matrice  $K \in \mathbb{R}^{m \times m}$  telle que  $K_{il} = k(\vec{x}_i, \vec{x}_l)$  est semi-définie positive

## 1.3 Astuce du noyau

Dès que je vois  $\langle \vec{x}, \vec{x}' \rangle$  : ① je peux le remplacer par  $k(\vec{x}, \vec{x}')$

$$k(\vec{x}, \vec{x}') = \langle \varphi(\vec{x}), \varphi(\vec{x}') \rangle_{\mathcal{H}}$$

- Si un algorithme ne fait intervenir les éléments de  $\mathcal{X}$  que dans des produits scalaires, **remplacer ces produits scalaires** par  $k$  est équivalent à appliquer l'algorithme dans  $\mathcal{H}$  après avoir appliqué  $\varphi$
- Utile **si  $k$  est plus simple à calculer que  $\varphi$**

Option ② : je calcule  $\varphi(\vec{x}), \varphi(\vec{x}')$  et j'applique mon  
algorithme aux  $\varphi(\vec{x})$

## 1.3 Astuce du noyau

$$k(\vec{x}, \vec{x}') = \langle \varphi(\vec{x}), \varphi(\vec{x}') \rangle_{\mathcal{H}}$$

- Si un algorithme ne fait intervenir les éléments de  $\mathcal{X}$  que dans des produits scalaires, **remplacer ces produits scalaires** par  $k$  est équivalent à appliquer l'algorithme dans  $\mathcal{H}$  après avoir appliqué  $\varphi$

- Utile **si  $k$  est plus simple à calculer que  $\varphi$**

- Exemple : **noyau quadratique**  $k : (\vec{x}, \vec{x}') \mapsto (1 + \langle \vec{x}, \vec{x}' \rangle)^2 \rightarrow \mathcal{O}(p)$  opérations

Équivaut à  $\varphi : (x_1, x_2, \dots, x_p) \mapsto (1, x_1, \dots, x_p, x_1^2, x_2^2, \dots, x_p^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{p-1}x_p) = \langle \varphi(\vec{x}), \varphi(\vec{x}') \rangle$

ce qu'on a utilisé pour faire des régressions polynomiales de degré  $d=2$ .  
calculer  $\varphi \Rightarrow p + p^2/2$  variables  
 $\hookrightarrow$  calcul de  $\langle \varphi(\vec{x}), \varphi(\vec{x}') \rangle \approx \mathcal{O}(p^2)$  opérations



## 1.4 Régression ridge à noyau

- **Régression ridge** :  $\arg \min_{\vec{\beta} \in \mathbb{R}^{p+1}} \underbrace{\frac{1}{n} (\vec{y} - X\vec{\beta})^\top (\vec{y} - X\vec{\beta})}_{\text{ERM}} + \lambda \underbrace{\|\vec{\beta}\|_2}_{\text{régularisation}}$
- Solution :  $\vec{\beta}^* = (X^\top X + \lambda I_p)^{-1} X^\top \vec{y}$
- Modèle :  $f : \vec{x} \mapsto \langle \vec{\beta}^*, \vec{x} \rangle = \sum \beta_j x_j$
- Reformulation :  $f : \vec{x} \mapsto \vec{x} X^\top (\lambda I_n + X X^\top)^{-1} \vec{y}$

# 1.4 Régression ridge à noyau

- **Régression ridge** :  $\arg \min_{\vec{\beta} \in \mathbb{R}^{p+1}} \frac{1}{n} \left( \vec{y} - X\vec{\beta} \right)^\top \left( \vec{y} - X\vec{\beta} \right) + \lambda \left\| \vec{\beta} \right\|_2$ 
    - Solution :  $\vec{\beta}^* = \left( X^\top X + \lambda I_p \right)^{-1} X^\top \vec{y}$
    - Modèle :  $f : \vec{x} \mapsto \langle \vec{\beta}^*, \vec{x} \rangle$
  - Reformulation :  $f : \vec{x} \mapsto \vec{x} X^\top \left( \lambda I_n + X X^\top \right)^{-1} \vec{y}$ 
    - $\vec{x} X^\top \in \mathbb{R}^n$  a pour  $i$ -ème entrée :  $\langle \vec{x}, \vec{x}_i \rangle$
    - $X X^\top \in \mathbb{R}^{n \times n}$  a pour entrée  $(i, l)$  :  $\langle \vec{x}_i, \vec{x}_l \rangle$
- i-ème observation du jeu d'apprentissage*

# 1.4 Régression ridge à noyau <sup>kernel</sup>

- **Régression ridge** :  $\arg \min_{\vec{\beta} \in \mathbb{R}^{p+1}} \frac{1}{n} (\vec{y} - X\vec{\beta})^\top (\vec{y} - X\vec{\beta}) + \lambda \|\vec{\beta}\|_2^2$   
interprétation de matrice p x p
- Solution :  $\vec{\beta}^* = (X^\top X + \lambda I_p)^{-1} X^\top \vec{y}$
- Modèle :  $f : \vec{x} \mapsto \langle \vec{\beta}^*, \vec{x} \rangle$
- Reformulation :  $f : \vec{x} \mapsto \vec{x} X^\top (\lambda I_n + X X^\top)^{-1} \vec{y}$   
inversion de matrice de taille n x n
- $\vec{x} X^\top \in \mathbb{R}^n$  a pour  $i$ -ème entrée :  $\langle \vec{x}, \vec{x}_i \rangle$
- $X X^\top \in \mathbb{R}^{n \times n}$  a pour entrée  $(i, l)$  :  $\langle \vec{x}_i, \vec{x}_l \rangle$
- **Kernelization** : remplacer  $\vec{x} X^\top$  par  $\kappa \in \mathbb{R}^n$  tel que  $\kappa_i = k(\vec{x}, \vec{x}_i)$  et  $X X^\top$  par  $K \in \mathbb{R}^{n \times n}$  telle que  $K_{il} = k(\vec{x}_i, \vec{x}_l)$

équivalent à transformer les données par  $\varphi$  puis apprendre une régression ridge, **pour environ le même coût algorithmique**

Option 1: Appliquer  $\varphi$  puis une régression ridge  
Option 2: Remplacer  $\langle \cdot, \cdot \rangle$  par  $k \rightarrow$   $k_{\text{np}} \oplus$  efficace

## 1.5 Exemples de noyaux

---

- **Noyau quadratique**  $k : (\vec{x}, \vec{x}') \mapsto (1 + \langle \vec{x}, \vec{x}' \rangle)^2$

$\varphi$  calcule tous les monômes de degré au plus 2 de  $x_1, x_2, \dots, x_p$

# 1.5 Exemples de noyaux

---

- **Noyau quadratique**  $k : (\vec{x}, \vec{x}') \mapsto (1 + \langle \vec{x}, \vec{x}' \rangle)^2$

$\varphi$  calcule tous les monômes de degré au plus 2 de  $x_1, x_2, \dots, x_p$

- **Noyau polynomial**  $k : (\vec{x}, \vec{x}') \mapsto (1 + \langle \vec{x}, \vec{x}' \rangle)^d$

$\varphi$  calcule tous les monômes de degré au plus  $d$  de  $x_1, x_2, \dots, x_p$

# 1.5 Exemples de noyaux

---

- **Noyau quadratique**  $k : (\vec{x}, \vec{x}') \mapsto (1 + \langle \vec{x}, \vec{x}' \rangle)^2$

$\varphi$  calcule tous les monômes de degré au plus 2 de  $x_1, x_2, \dots, x_p$

- **Noyau polynomial**  $k : (\vec{x}, \vec{x}') \mapsto (1 + \langle \vec{x}, \vec{x}' \rangle)^d$

$\varphi$  calcule tous les monômes de degré au plus  $d$  de  $x_1, x_2, \dots, x_p$

- **Noyau RBF gaussien**  $k : (\vec{x}, \vec{x}') \mapsto \exp\left(-\frac{\|\vec{x} - \vec{x}'\|_2^2}{\sigma^2}\right)$

$\varphi$  calcule tous les monômes de  $x_1, x_2, \dots, x_p$  et  $\mathcal{H}$  est de dimension infinie.

# 1.5 Exemples de noyaux

Lister toutes les chaînes de taille  $m$

AAAAA, AAAAS, ..., EXEMP, ..., MPLES, ..., ZZZZZ

"Exemples":  $[0, 0, \dots, 1, \dots, 1, \dots, 0]$

- **Noyau pour chaîne de caractères**  $k : (\vec{x}, \vec{x}') \mapsto \sum_{u \in \mathcal{A}^m} \psi_u(\vec{x}) \psi_u(\vec{x}')$
- $\mathcal{A}^m$  = l'ensemble des chaînes de  $m$  caractères sur l'alphabet  $\mathcal{A}$
- $\psi_u(\vec{x}) = 1$  si  $u$  est une sous-chaîne de  $\vec{x}$  et 0 sinon.

# 1.5 Exemples de noyaux

---

- **Noyau pour chaîne de caractères**  $k : (\vec{x}, \vec{x}') \mapsto \sum_{u \in \mathcal{A}^m} \psi_u(\vec{x}) \psi_u(\vec{x}')$ 
  - $\mathcal{A}^m$  = l'ensemble des chaînes de  $m$  caractères sur l'alphabet  $\mathcal{A}$
  - $\psi_u(\vec{x}) = 1$  si  $u$  est une sous-chaîne de  $\vec{x}$  et 0 sinon.
- $k(\vec{x}, \vec{x}')$  est le nombre de chaînes de  $m$  caractères communes à  $\vec{x}$  et  $\vec{x}'$  et peut être calculé en  $\mathcal{O}(|\vec{x}|)$  au lieu de  $\mathcal{O}(|\mathcal{A}|^m)$ .



# 1.5 Exemples de noyaux

---

- **Noyau pour chaîne de caractères**  $k : (\vec{x}, \vec{x}') \mapsto \sum_{u \in \mathcal{A}^m} \psi_u(\vec{x}) \psi_u(\vec{x}')$ 
  - $\mathcal{A}^m$  = l'ensemble des chaînes de  $m$  caractères sur l'alphabet  $\mathcal{A}$
  - $\psi_u(\vec{x}) = 1$  si  $u$  est une sous-chaîne de  $\vec{x}$  et 0 sinon.
- $k(\vec{x}, \vec{x}')$  est le nombre de chaînes de  $m$  caractères communes à  $\vec{x}$  et  $\vec{x}'$  et peut être calculé en  $\mathcal{O}(|\vec{x}|)$  au lieu de  $\mathcal{O}(|\mathcal{A}|^m)$ .
- Si  $m = 8$ ,  $|\mathcal{A}| = 20$  et en moyenne  $|\vec{x}| = 485$ , on compare 25,6 milliards d'opérations à moins de 500.

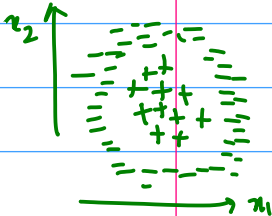
Méthodes à noyaux :

remplacer dans un algo les produits scalaires  
par des noyaux

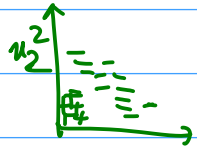
→ apprendre des modèles non-linéaires efficaces

→ appliquer des modèles linéaires à des objets  
dont la représentation en  $p$  dimensions n'est pas  
naturelle.

→ noyaux pour texte, images, molécules...



$$\psi(x_1, x_2) = (x_1^2, x_2^2)$$



linéairement séparable

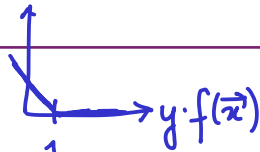
## 2. Machines à vecteurs de support

SVM (Support Vector Machines)

## 2.1 Formulation primale

si  $\vec{x}$  correctement classé:  $y \cdot f(\vec{x}) > 0$ .

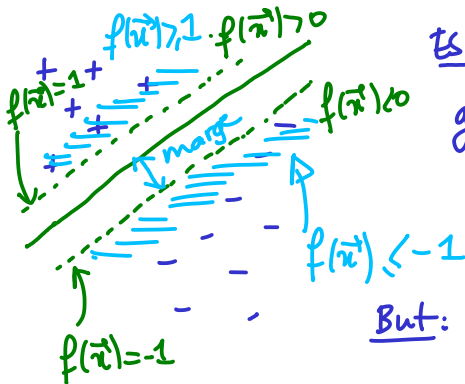
- Classification binaire,  $\mathcal{Y} = \{-1, 1\}$



marge =  $\frac{1}{\|\vec{\beta}\|}$  <sup>régularisation</sup>

perte hinge 1

$$\arg \min_{\vec{\beta} \in \mathbb{R}^p, b \in \mathbb{R}} \frac{1}{2} \|\vec{\beta}\|_2^2 + C \sum_{i=1}^n \max \left( 0, 1 - y_i \left( \langle \vec{\beta}, \vec{x}_i \rangle + b \right) \right)$$

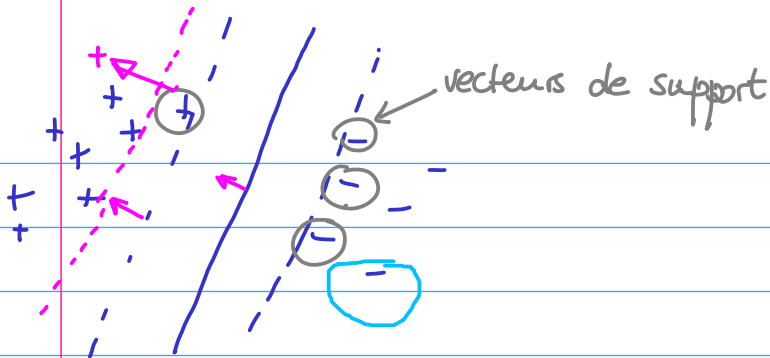


Espace des hypothèses

$$g: \vec{x} \mapsto \langle \vec{\beta}, \vec{x} \rangle + b$$

$$f(\vec{x}) = \begin{cases} 1 & \text{si } g(\vec{x}) > 0 \\ -1 & \text{sinon} \end{cases}$$

But: avoir une grande marge  
et aussi avoir peu d'erreurs.



## 2.2 Formulation duale

on est passés de chercher  $p$  coeffs à chercher  $n$  coeffs

$$\arg \max_{\vec{\alpha} \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{l=1}^n \alpha_i \alpha_l y_i y_l \langle \vec{x}_i, \vec{x}_l \rangle$$

sous les contraintes  $\sum_{i=1}^n \alpha_i y_i = 0$  et  $0 \leq \alpha_i \leq C$

– Modèle :  $f : \vec{x} \mapsto \sum_{i=1}^n \alpha_i y_i \langle \vec{x}_i, \vec{x} \rangle$

$\nearrow_{>0}$   $\alpha_i y_i$   $\nearrow$  étiquette de  $\vec{x}_i$   
similitude entre  $\vec{x}_i$  et  $\vec{x}$

Remplacer  $\langle \vec{x}_i, \vec{x}_l \rangle$  par  $k(\vec{x}_i, \vec{x}_l)$

et  $\langle \vec{x}_i, \vec{x} \rangle$  par  $k(\vec{x}_i, \vec{x})$

Quel que soit  $k$ , on cherche toujours seulement  $n$  coeffs  
( $\alpha_i$ )

$\Rightarrow$  limite le surapprentissage

## 2.3 SVR

## Support Vector Regression

- Perte  $\varepsilon$ -insensible :

$$f(\vec{x}) = \langle \vec{\beta}, \vec{x} \rangle + b$$

$$L(y, f(\vec{x})) = \max(0, |f(\vec{x}) - y| - \varepsilon)$$

