

SVM non-linéaire, et méthodes à noyaux

29 novembre 2019

Rappel SVM linéaire

Le modèle

- ▶ Données d'apprentissage $\{\mathbf{x}_i, y_i\}$
- ▶ Formulation du problème primal

$$\begin{aligned} \min_{w, b, \{\xi_i\}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.c.} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n \\ & \xi_i \geq 0 \quad \forall i = 1, \dots, n \end{aligned}$$

- ▶ Formulation du problème dual

$$\begin{aligned} \max_{\{\alpha_i\}} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.c.} \quad & 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

Handwritten notes:

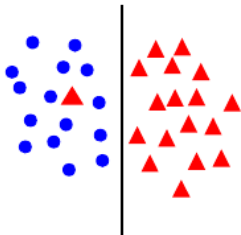
- $\rightarrow f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$
- classifier predicts class 1 if $f(\mathbf{x}) > 0$
- $\omega = \sum \alpha_i y_i \mathbf{x}_i$
- $f(\mathbf{x}) = \omega^\top \mathbf{x} + b$
- $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^\top \mathbf{x} + b$

- ▶ La fonction de décision $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^\top \mathbf{x}$

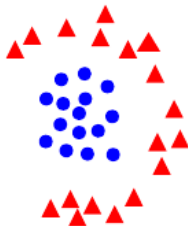
Non-linéarité ?

Exemples de cas d'usage des SVM

- problème linéaire non-séparable



- problème séparable mais non-linéaire ?

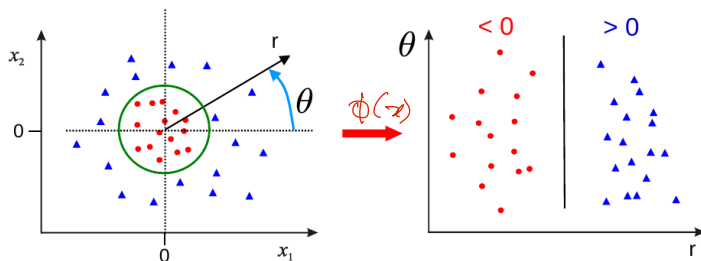


Quelles solutions ?

- transformation non-linéaire

Exemples de transformation non-linéaire

Coordonnées polaires

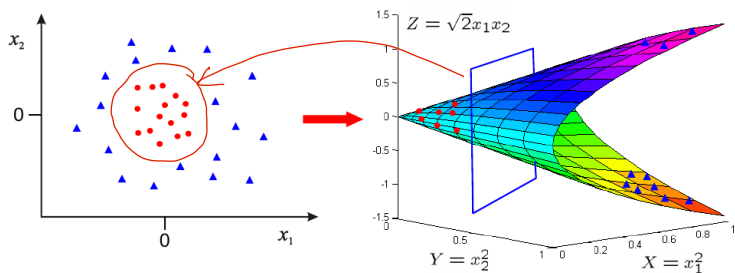


- Les données sont linéairement séparables dans l'espace en coordonnées polaire
- agit comme une transformation non-linéaire de l'espace original

$$\Phi(\mathbf{x}) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} r \\ \theta \end{pmatrix}$$

Exemples de transformation non-linéaire

Projection dans un espace de plus haute dimension



- les données sont séparables dans un espace 3D.

$$\Phi(\mathbf{x}) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}$$

*if I learn
a hyperplane*

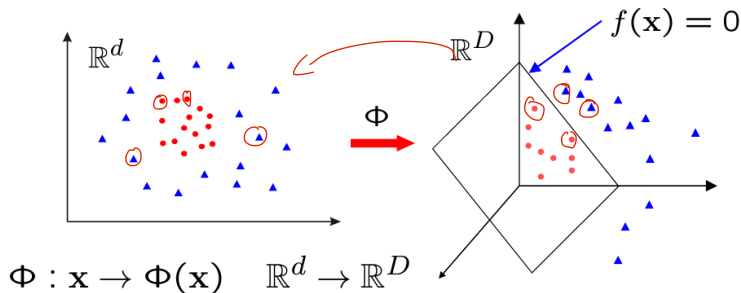
$$\omega^+ \phi(x) + b = 0$$



$$\omega_1 x_1^2 + \omega_2 x_2^2 + \omega_3 \sqrt{2} x_1 x_2 + b = 0$$

- on peut utiliser un SVM linéaire dans un autre espace.

SVM dans un espace transformé



- Apprendre un classifieur linéaire dans le nouvel espace \mathbb{R}^D

$$f(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x}) + b$$

où $\Phi(\mathbf{x})$ est la fonction de transformation des données avec $\Phi : \mathbb{R}^d \mapsto \mathcal{F}$
(dans l'exemple $\mathcal{F} = \mathbb{R}^D$)

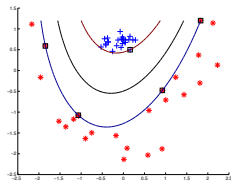
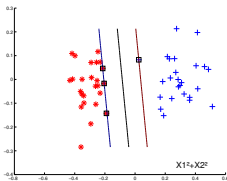
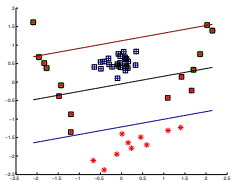
Principes de la non-linéarisation

- ▶ On projette les données \mathbf{x} grâce à une transformation Φ dans un espace \mathcal{F} . La fonction de décision devient $f(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x}) + b$
- ▶ On applique l'algorithme linéaire dans l'espace \mathcal{F} .

$$\begin{array}{ll} \min_{\mathbf{w}, b, \{\xi_i\}} & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.c.} & y_i (\mathbf{w}^\top \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n \\ & \xi_i \geq 0 \quad \forall i = 1, \dots, n \end{array}$$

avec $\mathbf{w} \in \mathcal{F}$.

- ▶ La fonction de transfert obtenue est non-linéaire dans l'espace original.



Problème dual et fonction de décision transformée

Fonction de décision duale

La fonction de décision pour les SVMs

$$f(x) = \sum_{i=1}^n y_i \alpha_i x^\top x_i + b \quad \Rightarrow \quad f(x) = \sum_{i=1}^n y_i \alpha_i \Phi(x)^\top \Phi(x_i) + b$$

Problème dual

$$\begin{array}{ll} \max_{\{\alpha_i\}} & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \Phi(x_i)^\top \Phi(x_j) \\ \text{s.c.} & 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{array}$$

Astuce du noyau

Constat

- ▶ La fonction $\Phi(\mathbf{x})$ intervient toujours sous la forme $\Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$
- ▶ dans le problème dual, une fois que tout les produits scalaires $\Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$ ont été calculés, on n'a besoin que résoudre le problème dual en $\alpha \in \mathbb{R}^n$.
- ▶ On n'a pas besoin de calculer $\mathbf{w} \in \mathbb{R}^D$. c'est avantageux si D est très grand.

Transformation implicite par un noyau $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$

- ▶ Fonction de décision

$$f(x) = \sum_{i=1}^{\ell} \alpha_i k(x, x_i) + b$$

- ▶ problème dual

① the dim
D does not appear

$$\begin{array}{ll} \max_{\{\alpha_i\}} & \\ \text{s.c.} & \end{array}$$

$$\begin{array}{l} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{array}$$

Kernel Ridge Regression

regression with ℓ_1 - ℓ_2 regularization

Formulation regression ridge

- ▶ $\mathcal{D} = \{(x_i, y_i)\} \in \mathcal{X} \times \mathbb{R}$, $i = 1 \dots n$: ensemble de points étiquetés .
- ▶ Modèle linéaire : $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$
- ▶ Coût ℓ_2 pénalité ℓ_2

$$\mathbf{X} = \begin{bmatrix} -x_1 & - \\ & -x_N \end{bmatrix} \in \mathbb{R}^{N \times d}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$$

Apprentissage du modèle

- ▶ Optimisation

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- ▶ Condition d'optimalité $-\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}^*) + \lambda \mathbf{w}^* = 0$

$$-\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X} \mathbf{w}^* + \lambda \mathbf{w}^* = 0$$
$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{w}^* = \mathbf{X}^\top \mathbf{y}$$

- ▶ Solution

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{y})$$

Ici $\mathbf{X}^\top \mathbf{X}$ représente une matrice de covariance $\in \mathbb{R}^{d \times d}$

Kernel Ridge regression : reformulation

- La condition d'optimalité se réécrit :

$$\mathbf{w}^* = \frac{1}{\lambda} \mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{w}^*)$$

replace \mathbf{w}^* by $\mathbf{X}^T \alpha^*$

$$\mathbf{X}^T \alpha^* = \frac{1}{\lambda} \mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{X}^T \alpha^*)$$

- on peut donc dire qu'il existe un vecteur $\alpha \in \mathbb{R}^n$ tel que

$$\mathbf{w}^* = \mathbf{X}^T \alpha = \sum_i \mathbf{x}_i \alpha_i$$

consider the $\hat{\alpha}$ vector st.

$$(\mathbf{I} + \frac{1}{\lambda} \mathbf{X} \mathbf{X}^T) \hat{\alpha} = \frac{1}{\lambda} \mathbf{y}$$

if $\hat{\alpha}$ exists, then $\alpha^* = \hat{\alpha}$

$$(\lambda \mathbf{I} + \mathbf{X} \mathbf{X}^T) \hat{\alpha} = \mathbf{y}$$

- Dans ce contexte, on a $f(\mathbf{x}) = \sum_i \alpha_i \mathbf{x}_i^T \mathbf{x}$

- les variables α s'obtiennent par $\alpha = \frac{1}{\lambda} (\mathbf{y} - \mathbf{X} \mathbf{w}^*) = \frac{1}{\lambda} (\mathbf{y} - \mathbf{X} \mathbf{X}^T \alpha)$ donc

$$\alpha = (\lambda \mathbf{I} + \mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$$

$$\hat{\alpha} = \alpha^* = (\lambda \mathbf{I} + \mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$$

ici, $(\mathbf{X} \mathbf{X}^T)_{i,j} = \mathbf{x}_i^T \mathbf{x}_j$

→ the gram matrix $\in \mathbb{R}^{n \times n}$

Kernel Ridge Regression : reformulation

- la fonction de décision :

$$f(\mathbf{x}) = \sum_i \alpha_i \mathbf{x}_i^T \mathbf{x}$$

- le problème d'apprentissage

$$\min_{\alpha} \frac{1}{2} \|\mathbf{y} - \underbrace{\mathbf{X}\mathbf{X}^T}_{\text{Gram matrix}} \alpha\|_2^2 + \frac{\lambda}{2} \alpha^T \underbrace{\mathbf{X}\mathbf{X}^T}_{\text{K}} \alpha$$

$\omega = \mathbf{X}^T \alpha$
 $\|\omega\|^2$

$$\mathbf{u}^T \mathbf{M} \mathbf{u} = \sum_i \sum_j u_i u_j M_{ij}$$

$$\min_{\alpha} \sum_{i=1}^N \left(y_i - \alpha^T (k(x_i, x_1) \dots k(x_i, x_N)) \right)^2 + \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k(x_i, x_j)$$

Kernelized

- l'ensemble du problème se reformule en fonction des produits scalaires $\mathbf{x}_i^T \mathbf{x}_j$

If we apply the change of representation $x \mapsto \phi(x)$ then, we replace \mathbf{X} by $\begin{bmatrix} -\phi(x_1) - \\ \vdots \\ -\phi(x_N) - \end{bmatrix}$

so \mathbf{K} becomes $\begin{bmatrix} \phi(x_1)^T \phi(x_1) \\ \vdots \\ \phi(x_N)^T \phi(x_N) \end{bmatrix}$ so $\mathbf{K} = [k(x_i, x_j)]_{i,j}$

Kernel Kmeans

Principe

- remplacer la distance euclidienne par la distance dans l'espace transformée

$$d(\mathbf{x}_i, \mu_k)^2 = \|\Phi(\mathbf{x}_i) - \Phi(\mu_k)\|^2 = k(\mathbf{x}_i, \mathbf{x}_i) + k(\mu_k, \mu_k) - 2k(\mathbf{x}_i, \mu_k)$$

où $\Phi(\cdot)$ est la transformation implicite et $k(\cdot, \cdot)$ le produit scalaire dans l'espace transformée

Détail

- μ_k moyenne des $\Phi(\mathbf{x}_i)$ du cluster
- $k(\mathbf{x}_i, \mu_k) = \frac{1}{|J|} \sum_{j \in J} \Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$