

SVM

Vincent Lefieux



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Plan

Introduction

Cas linéairement séparable

Cas non-séparable

Astuce du noyau

Cas de la régression

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Plan

Introduction

Cas linéairement séparable

Cas non-séparable

Astuce du noyau

Cas de la régression

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Généralités I

- ▶ Les **SVM** (**Support Vector Machine**) (en français : *séparateurs à vaste marge* ou *machines à vecteurs supports*) sont issus de la théorie de Vapnik-Tchervonenkis (dénommée théorie VC) : (Cortes et Vapnik, 1995), (Vapnik, 1995).
- ▶ L'objectif historique des SVM est de **classifier une variable binaire** via un **hyperplan de marge maximale**, les SVM constituent une généralisation des **classifieurs linéaires**.
- ▶ Les SVM intègrent le **contrôle de la complexité**, ce qu'on peut appréhender via la dimension de Vapnik-Tchervonenkis qui est un indicateur du pouvoir séparateur d'une famille de fonctions.
- ▶ C'est une méthode souvent utilisée en pratique au vu des bons résultats obtenus.

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

- ▶ On parle de **marge** (*hard margin*) lorsque les données sont linéairement séparables et de **marge souple** (*soft margin*) lorsque les données ne le sont pas.
- ▶ Dans le cas où les données ne sont pas linéairement séparables, on utilise ce qu'on appelle l'**astuce du noyau** (*kernel trick*).
- ▶ Il existe également les **SVR** dans le cadre de la régression.
- ▶ Il faut **normaliser les covariables**.

Données considérées

- ▶ On dispose d'un échantillon de (X, Y) :

$$\mathcal{D}_n = (X_i, Y_i)_{i \in \{1, \dots, n\}} .$$

On note :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}} .$$

- ▶ On considère dans la suite que :
 - ▶ $X \in \mathbb{R}^p$:
*Toutes les covariables sont considérés quantitatifs.
Mais il est également possible de considérer des
covariables qualitatives.*
 - ▶ $Y \in \{-1, 1\}$:
*On se place dans le cadre d'une classification supervisée
binaire*

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Généralités sur les hyperplans I

- Dans \mathbb{R}^p , un hyperplan \mathcal{H} admet comme équation :

$$\omega_0 + \omega_1 x_1 + \dots + \omega_p x_p = 0 ,$$

ce qu'on peut noter également :

$$\omega_0 + \langle \omega, x \rangle = 0$$

ou encore :

$$\omega_0 + \omega^\top x = 0$$

où $\omega = (\omega_1, \dots, \omega_p)^\top \in \mathbb{R}^p$ et $x = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$.

- ω est le vecteur normal de l'hyperplan \mathcal{H} .
- Par exemple : un hyperplan dans \mathbb{R}^2 est une droite, un hyperplan dans \mathbb{R}^3 est un plan.

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Généralités sur les hyperplans II

$$w_0 + w^\top x = 0$$

\mathcal{H}

w

x_i

$$d(x_i, \mathcal{H}) = \frac{|w_0 + w^\top x_i|}{\|w\|}$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Plan

Introduction

Cas linéairement séparable

Cas non-séparable

Astuce du noyau

Cas de la régression

Introduction

**Cas linéairement
séparable**

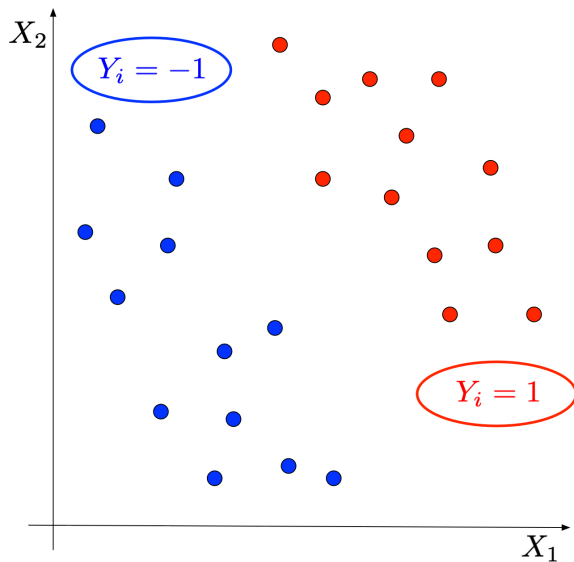
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Données linéairement séparables I



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Données linéairement séparables II

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

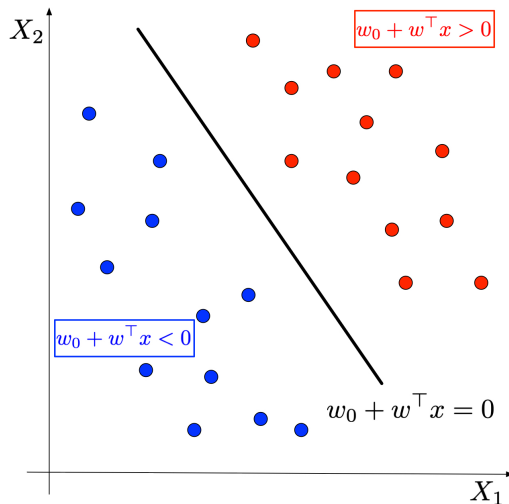
- ▶ On dit que $(x_1, y_1), \dots, (x_n, y_n)$ sont **linéairement séparables** s'il existe $(\omega_0, \omega) \in \mathbb{R} \times \mathbb{R}^d$ tels que :

$$\forall i \in \{1, \dots, n\} : y_i = \begin{cases} 1 & \text{si } \omega_0 + \omega^\top x_i > 0 \\ -1 & \text{si } \omega_0 + \omega^\top x_i < 0 \end{cases} .$$

- ▶ Cette propriété est équivalente à :

$$y_i (\omega_0 + \omega^\top x_i) > 0 .$$

Données linéairement séparables III



Introduction

Cas linéairement
séparable

Cas
non-séparable

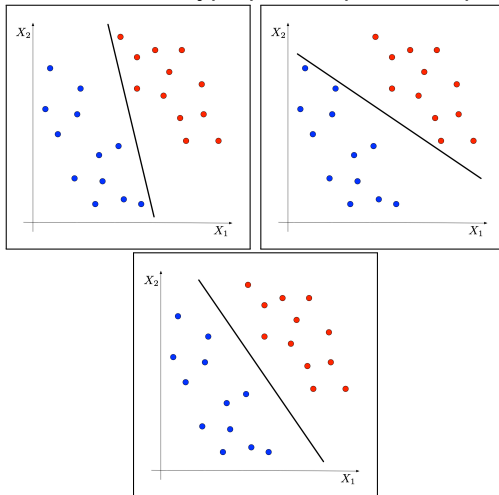
Astuce du noyau

Cas de la
régression

Références

Le choix de l'hyperplan séparateur

- Il existe une infinité d'hyperplans séparateurs possibles :



- Vapnik a proposé de **maximiser la marge**, soit la distance minimale entre les 2 classes déterminées par l'hyperplan séparateur.

Introduction

Cas linéairement
séparable

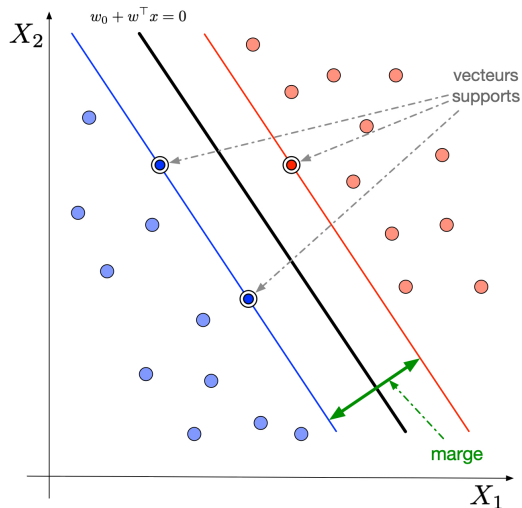
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Marge et vecteurs supports



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Formalisation du problème I

- On pose comme contrainte que les vecteurs supports sont situés sur les hyperplans canoniques d'équations :

$$\begin{cases} \omega_0 + \omega^\top x = -1 \\ \omega_0 + \omega^\top x = 1 \end{cases} .$$

- La marge vaut dans ce cas :

$$\frac{2}{\|\omega\|} .$$

Introduction

Cas linéairement
séparable

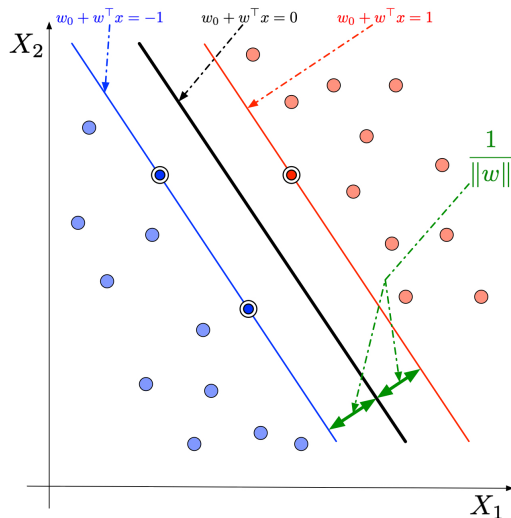
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Formalisation du problème II



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Formalisation du problème III

- On obtient donc le problème suivant :

$$\begin{aligned} & \max_{\omega_0, \omega} \frac{2}{\|\omega\|} \\ \text{sc} \quad & \forall i \in \{1, \dots, n\} : y_i \left(\omega_0 + \omega^\top x_i \right) \geq 1 . \end{aligned}$$

- Dans la suite, on considère le **problème primal** équivalent :

$$\begin{aligned} & \min_{\omega_0, \omega} \frac{1}{2} \|\omega\|^2 \\ \text{sc} \quad & \forall i \in \{1, \dots, n\} : y_i \left(\omega_0 + \omega^\top x_i \right) \geq 1 . \end{aligned}$$

- Le carré et la division par 2 ont comme seul objectif d'améliorer la lisibilité des résultats obtenus.
- Il s'agit d'un programme d'**optimisation quadratique** classique.

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Règle de classification

La règle de classification obtenue est :

$$\forall x \in \mathbb{R}^p : g(x) = \begin{cases} 1 & \text{si } \omega_0^* + \omega^{*\top} x > 0 \\ -1 & \text{si } \omega_0^* + \omega^{*\top} x < 0 \end{cases} .$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Plan

Introduction

Cas linéairement séparable

Cas non-séparable

Astuce du noyau

Cas de la régression

Introduction

Cas linéairement
séparable

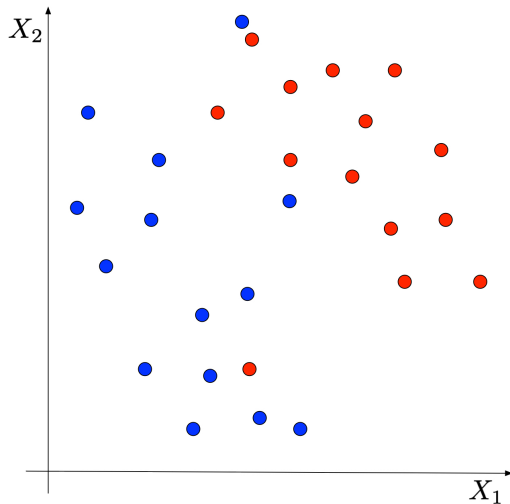
**Cas
non-séparable**

Astuce du noyau

Cas de la
régression

Références

Exemple non-séparable



Lever les contraintes I

- ▶ Il est rare d'être confronté à un problème linéairement séparable.
- ▶ On **lève la contrainte** en tolérant que :
 - ▶ certains **points** soient **bien classés** mais à l'intérieur de la zone définie par la marge,
 - ▶ certains **points** soient **mal classés**.

Introduction

Cas linéairement
séparable

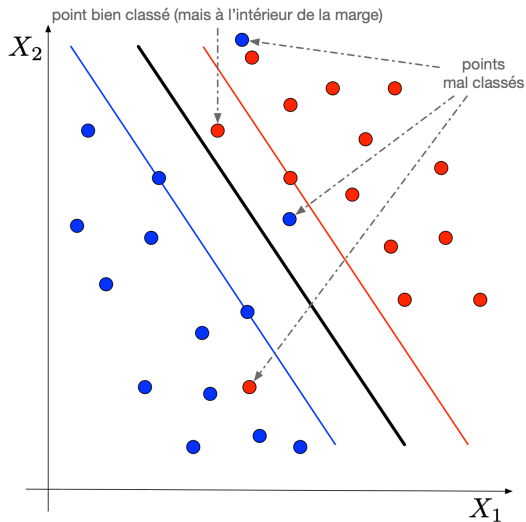
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Lever les contraintes II



Introduction

Cas linéairement
séparable

**Cas
non-séparable**

Astuce du noyau

Cas de la
régression

Références

Un outil : les variables ressorts I

- ▶ On crée des variables ressorts (*slack variables*) (ξ_1, \dots, ξ_n) telles que :

$$y_i \left(\omega_0 + \omega^\top x_i \right) \geq 1 - \xi_i .$$

- ▶ On peut distinguer les cas suivants :
 - ▶ $\xi_i \in]0, 1]$: les points sont bien classés mais à l'intérieur (strictement) de la zone définie par la marge.
 - ▶ $\xi_i > 1$: les points sont mal classés.
 - ▶ $\xi_i = 0$: les points sont bien classés et à l'extérieur de la zone définie par la marge.
- ▶ L'enjeu est de ne pas avoir trop de variables ressorts non nulles (et lorsqu'elles le sont, qu'elles soient les plus faibles possibles).

Introduction

Cas linéairement
séparable

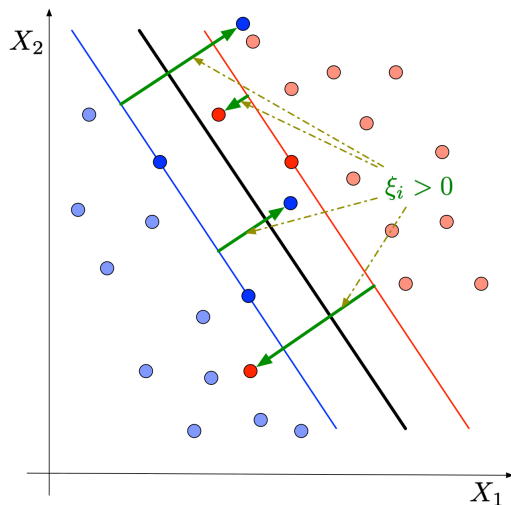
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Un outil : les variables ressorts II



Introduction

Cas linéairement
séparable

**Cas
non-séparable**

Astuce du noyau

Cas de la
régression

Références

Un nouveau problème I

Introduction

Cas linéairement
séparable

**Cas
non-séparable**

Astuce du noyau

Cas de la
régression

Références

On considèrerait le problème suivant :

$$\begin{array}{ll} \min_{\omega_0, \omega} & \frac{1}{2} \|\omega\|^2 \\ \text{sc} & \forall i \in \{1, \dots, n\} : y_i (\omega_0 + \omega^\top x) \geq 1 . \end{array}$$

Un nouveau problème II

Introduction

Cas linéairement
séparable

**Cas
non-séparable**

Astuce du noyau

Cas de la
régression

Références

On considère maintenant le problème suivant, avec $\xi = (\xi_1, \dots, \xi_n)^\top$:

$$\begin{aligned} \min_{\omega_0, \omega, \xi} \quad & \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i \\ \text{sc} \quad & \forall i \in \{1, \dots, n\} : y_i (\omega_0 + \omega^\top x) \geq 1 - \xi_i, \\ & \forall i \in \{1, \dots, n\} : \xi_i \geq 0. \end{aligned}$$

Choix de l'hyper-paramètre C

- ▶ L'hyper-paramètre C contrôle de le compromis entre le nombre d'erreurs de classification et le niveau de la marge.
- ▶ Le cas linéairement séparable correspond à une valeur C infinie.
- ▶ On choisit l'hyper-paramètre C par **validation croisée**.

Introduction

Cas linéairement
séparable

**Cas
non-séparable**

Astuce du noyau

Cas de la
régression

Références

Plan

Introduction

Cas linéairement séparable

Cas non-séparable

Astuce du noyau

Cas de la régression

Introduction

Cas linéairement
séparable

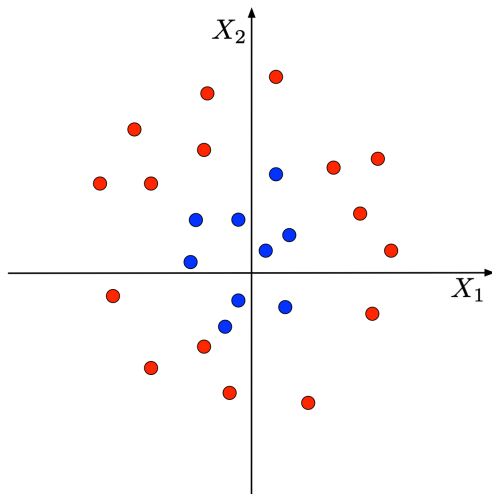
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Changer la dimension I



Introduction

Cas linéairement
séparable

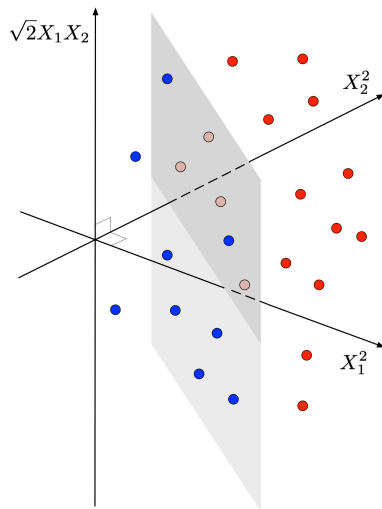
Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Changer la dimension II



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Astuce du noyau

- ▶ Déterminer un classifieur linéaire dans l'espace des observations n'est pas toujours opportun. On espère que la séparation linéaire sera plus simple dans un nouvel espace.
- ▶ On « envoie » les observations (dans l'espace \mathcal{X}) dans un nouvel espace \mathcal{X}' : l'espace de représentation (*feature space*).
- ▶ On considère pour cela une fonction Φ définie sur \mathcal{X} et à valeurs dans \mathcal{X}' .
- ▶ Dans le problème d'optimisation des SVM, on retrouve les produits $x_i^\top x_j$ dans l'espace des observations, donc des produits $\Phi(x_i)^\top \Phi(x_j)$ dans l'espace de représentation.
- ▶ Il n'est pas nécessaire de déterminer Φ , on utilisera des noyaux K tels que :

$$K(x_i, x_j) = \Phi(x_i)^\top \Phi(x_j)$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Retour au problème d'optimisation I

Dans le cas linéairement séparable, on devait résoudre le problème dual suivant dans l'espace des observations :

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^\top x_j \\ \text{sc} \quad & \forall i \in \{1, \dots, n\} : \alpha_i \geq 0, \\ & \sum_{i=1}^n \alpha_i y_i = 0. \end{aligned}$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Retour au problème d'optimisation II

Dans le cas linéairement séparable, on doit maintenant résoudre le problème dual suivant dans l'espace de représentation :

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \Phi(x_i)^{\top} \Phi(x_j)$$

$$\text{sc} \quad \forall i \in \{1, \dots, n\} : \alpha_i \geq 0 ,$$

$$\sum_{i=1}^n \alpha_i y_i = 0 .$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Retour au problème d'optimisation III

Dans le cas linéairement séparable, on doit résoudre le problème dual suivant dans l'espace de représentation :

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{sc } \forall i \in \{1, \dots, n\} : \alpha_i \geq 0 ,$$

$$\sum_{i=1}^n \alpha_i y_i = 0 .$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Une fonction $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ est un **noyau** si et seulement si :

- K est une fonction **symétrique** :

$$\forall (x, x') \in \mathcal{X} \times \mathcal{X} : K(x, x') = K(x', x) .$$

- K est une fonction **semi-définie positive** :

$$\forall n \in \mathbb{N}^*, \forall (x_1, \dots, x_n) \in \mathcal{X}^n, \forall (a_1, \dots, a_n) \in \mathbb{R}^n :$$

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j K(x_i, x_j) \geq 0 .$$

Un exemple de noyau

- Pour une observation $x_i = (x_{i1}, x_{i2})^\top$, on considère la fonction suivante :

$$\begin{aligned}\Phi : \quad \mathbb{R}^2 &\rightarrow \mathbb{R}^3 \\ (x_{i1}, x_{i2})^\top &\mapsto (x_{i1}^2, \sqrt{2} x_{i1} x_{i2}, x_{i2}^2)^\top\end{aligned}$$

- On peut montrer que pour 2 observations x_i et x_j :

$$\begin{aligned}K(x_i, x_j) &= \Phi(x_i)^\top \Phi(x_j) \\ &= (x_{i1}x_{j1})^2 + 2(x_{i1}x_{j1})(x_{i2}x_{j2}) + (x_{i2}x_{j2})^2 \\ &= (x_{i1}x_{j1} + x_{i2}x_{j2})^2 \\ &= \left(x_i^\top x_j\right)^2.\end{aligned}$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Quelques noyaux (parmi bien d'autres)

- Noyau **affine** :

$$K(x_i, x_j) = x_i^\top x_j + c .$$

- Noyau **polynomial** :

$$K(x_i, x_j) = \left(x_i^\top x_j + c \right)^d .$$

- Noyau **laplacien** :

$$K(x_i, x_j) = \exp \left(-\frac{\|x_i - x_j\|}{\sigma} \right) .$$

- Noyau **gaussien** (ou **RBF** : Radial Basis Function) :

$$K(x_i, x_j) = \exp \left(-\frac{\|x_i - x_j\|^2}{2\sigma^2} \right) .$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

En pratique

On choisit :

- ▶ l'hyper-paramètre C ,
- ▶ le noyau K ,

par validation croisée.

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

- ▶ Dans le cas où on dispose de $K > 2$ classes, on peut par exemple considérer K discriminations binaires « classe k » contre « classe autre que k » pour $k \in \{1, \dots, K\}$.
- ▶ Il est également possible d'utiliser ces méthodes pour la régression : on parle alors de **SVR** : (Drucker et collab., 1997), (Vapnik et collab., 1997).

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Plan

Introduction

Cas linéairement séparable

Cas non-séparable

Astuce du noyau

Cas de la régression

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

**Cas de la
régression**

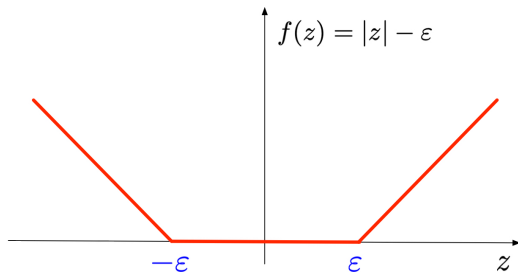
Références

Fonction de perte

Vapnik a introduit la fonction de perte suivante (ε -insensitive loss function) pour mesurer la qualité de l'ajustement de la fonction de régression m :

$$\ell(m(x), y) = \begin{cases} |m(x) - y| - \varepsilon & \text{si } |m(x) - y| > \varepsilon \\ 0 & \text{sinon} \end{cases}$$

où $\varepsilon > 0$.



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Risque empirique

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Le risque empirique vaut :

$$R_n(m) = \sum_{i=1}^n (|m(x_i) - y_i| - \varepsilon) = \sum_{i=1}^n (\xi_i + \xi_i^*)$$

où :

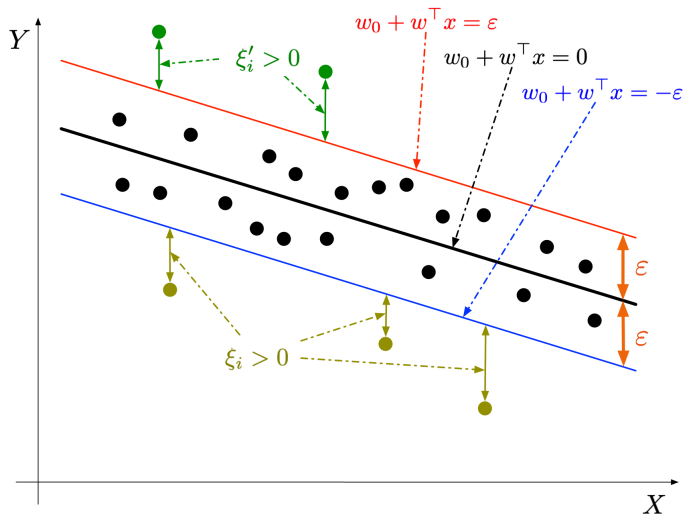
$$\begin{cases} \xi_i = m(x_i) - \varepsilon - y_i & \text{si } y_i < m(x_i) - \varepsilon \\ 0 & \text{sinon} \end{cases}$$

et :

$$\begin{cases} \xi_i^* = y_i - m(x_i) - \varepsilon & \text{si } y_i > m(x_i) + \varepsilon \\ 0 & \text{sinon} \end{cases}$$

.

Cas linéaire I



Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Cas linéaire II

- On considère la fonction de régression :

$$\forall x \in \mathbb{R}^p : m_{\omega_0, \omega}(x) = \omega_0 + \omega^\top x .$$

- On cherche ω_0 et ω de manière à minimiser la somme de la perte qui traduit l'ajustement et d'un terme de régularisation (assurant la parcimonie) $\|\omega\|^2$.
- On considère le problème suivant :

$$\begin{aligned} \min_{\omega_0, \omega} \quad & \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \\ \text{sc} \quad & \forall i \in \{1, \dots, n\} : m_{\omega_0, \omega}(x_i) - y_i \leq \varepsilon + \xi_i , \\ & \forall i \in \{1, \dots, n\} : y_i - m_{\omega_0, \omega}(x_i) \leq \varepsilon + \xi_i^* , \\ & \forall i \in \{1, \dots, n\} : \xi_i \geq 0 , \\ & \forall i \in \{1, \dots, n\} : \xi_i^* \geq 0 . \end{aligned}$$

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Choix des hyper-paramètres ε et C

- ▶ L'hyper-paramètre ε contrôle la largeur du « tube » : plus ε est important, moins on a de vecteurs support et plus lisse est l'estimation.
- ▶ L'hyper-paramètre C contrôle de le compromis entre l'erreur d'ajustement et le niveau de la marge. On le choisit par **validation croisée**.

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références

Références

- Boyd, S. et L. Vandenberghe. 2003, *Convex optimization*, Cambridge University Press.
- Cortes, C. et V. N. Vapnik. 1995, «Support-vector networks», *Machine Learning*, vol. 20, n° 3, p. 273–297.
- Drucker, H., C. J. Burges, L. Kaufman, A. Smola et V. N. Vapnik. 1997, *Advances in Neural Information Processing Systems*, vol. 9, chap. Support vector regression machines, MIT Press, p. 155–161.
- Schölkopf, B. et A. J. Smola. 2001, *Learning with Kernels. Support vector machines, regularization, optimization, and beyond*, MIT Press.
- Vapnik, V. N. 1995, *The nature of statistical learning theory*, Springer.
- Vapnik, V. N., S. E. Golowich et A. Smola. 1997, *Advances in Neural Information Processing Systems*, vol. 9, chap. Support vector method for function approximation, regression estimation, and signal processing, MIT Press, p. 281–287.

Introduction

Cas linéairement
séparable

Cas
non-séparable

Astuce du noyau

Cas de la
régression

Références