

SVM et machines à noyaux

Stéphane Canu
stephane.canu@litislab.eu

Ecole d'été du GRETSI - Peyresq

June 27, 2010

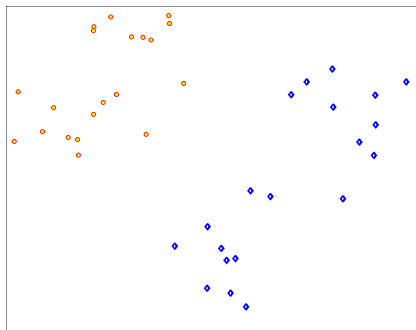
Plan

1 Introduction aux SVM

- Le problème de discrimination linéaire
- La notion de marge
- Le problème des SVM linéaires
- Formulation duale des SVM linéaires
- Méthodes de résolution du problème dual

Le problème de discrimination linéaire

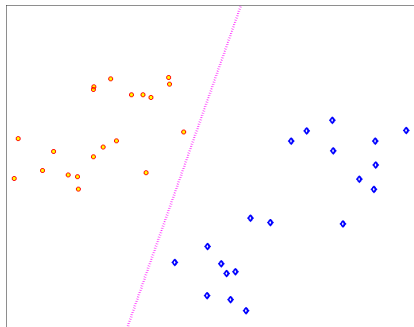
Trouver une droite que sépare les bleus des rouges



$$D(x) = \text{sign}(\mathbf{v}^T \mathbf{x} + a)$$

Le problème de discrimination linéaire

Trouver une droite que sépare les bleus des rouges



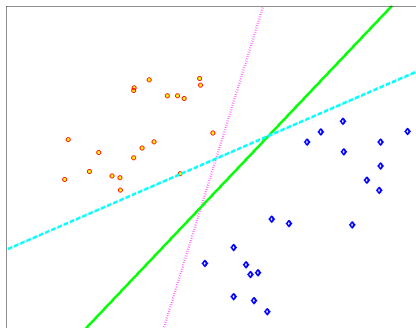
$$D(x) = \text{sign}(\mathbf{v}^T \mathbf{x} + a)$$

la frontière de décision:

$$\mathbf{v}^T \mathbf{x} + a = 0$$

Le problème de discrimination linéaire

Trouver une droite que sépare les bleus des rouges



$$D(x) = \text{sign}(\mathbf{v}^T \mathbf{x} + a)$$

la frontière de décision:

$$\mathbf{v}^T \mathbf{x} + a = 0$$

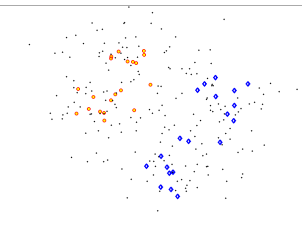
Il y a plusieurs solutions...

Le problème est mal posé

Quelle solution choisir ?

Ce problème n'est pas tout à fait celui qui nous intéresse

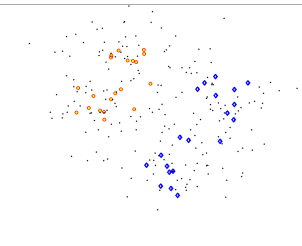
$\{(x_i, y_i); i = 1 : n\}$ un échantillon tiré suivant $\mathbb{P}(x, y)$ **inconnue**



Nous cherchons à bien classer des points issus de cette même distribution : minimiser $\mathbb{P}(\text{erreur})$

Ce problème n'est pas tout à fait celui qui nous intéresse

$\{(x_i, y_i); i = 1 : n\}$ un échantillon tiré suivant $\mathbb{P}(x, y)$ **inconnue**



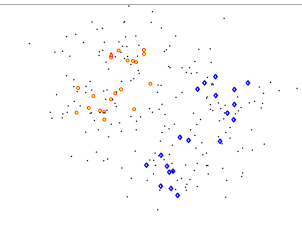
Nous cherchons à bien classer des points issus de cette même distribution : minimiser $\mathbb{P}(\text{erreur})$

A la recherche d'une méthode « universelle »

- être fidèle aux données : (faire **peu** d'erreurs)
- garantir que $\mathbb{P}(\text{erreur})$ ne s'éloigne pas trop
- non linéaire
- passe à l'échelle - complexité algorithmique

Ce problème n'est pas tout à fait celui qui nous intéresse

$\{(x_i, y_i); i = 1 : n\}$ un échantillon tiré suivant $\mathbb{P}(x, y)$ **inconnue**



Nous cherchons à bien classer des points issus de cette même distribution : minimiser $\mathbb{P}(\text{erreur})$

A la recherche d'une méthode « universelle »

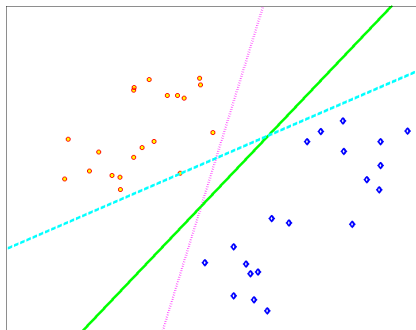
- être fidèle aux données : (faire **peu** d'erreurs)
- garantir que $\mathbb{P}(\text{erreur})$ ne s'éloigne pas trop
- non linéaire
- passe à l'échelle - complexité algorithmique

avec une grande probabilité :

$$\mathbb{P}(\text{erreur}) < \underbrace{\hat{\mathbb{P}}(\text{erreur})}_{=0 \text{ ici}} + \varphi\left(\underbrace{\frac{\|\mathbf{v}\|}{1}}_{\text{marge}}\right)$$

Le problème de discrimination linéaire

Trouver une droite que sépare les bleus des rouges



$$D(x) = \text{sign}(\mathbf{v}^T \mathbf{x} + a)$$

la frontière de décision:

$$\mathbf{v}^T \mathbf{x} + a = 0$$

Il y a plusieurs solutions...

Le problème est mal posé

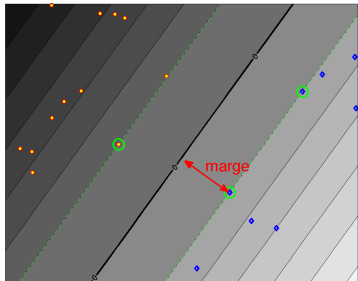
Quelle solution choisir ?

⇒

Celle qui aura la plus grande marge

Maximiser la « confiance » = maximiser la marge

La frontière de décision : $\Delta(\mathbf{v}, a) = \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{v}^\top \mathbf{x} + a = 0\}$



Maximiser la marge

$$\max_{\mathbf{v}, a} \underbrace{\min_{i \in [1, n]} \text{dist}(\mathbf{x}_i, \Delta(\mathbf{v}, a))}_{\text{marge : } m}$$

Maximiser la marge

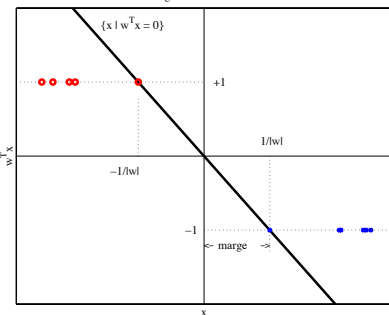
$$\begin{cases} \max_{\mathbf{v}, a} & m \\ \text{avec} & \min_{i=1, n} |\mathbf{v}^\top \mathbf{x}_i + a| \geq m \end{cases}$$

Le problème est toujours mal posé

si (\mathbf{v}, a) est une solution, $\forall k \in \mathbb{R}^+$ $(k\mathbf{v}, ka)$ l'est aussi...

Marge, norme et régularité

Valeur de la marge dans le cas monodimensionnel



Maximiser la marge

$$\left\{ \begin{array}{l} \max_{\mathbf{v}, a} \quad m \\ \text{avec} \quad \min_{i=1, n} |\mathbf{v}^\top \mathbf{x}_i + a| \geq m \\ \|\mathbf{v}\|^2 = 1 \end{array} \right.$$

si le min est plus grand, tout le monde est plus grand ($y_i \in \{-1, 1\}$)

$$\left\{ \begin{array}{l} \max_{\mathbf{v}, a} \quad m \\ \text{avec} \quad y_i(\mathbf{v}^\top \mathbf{x}_i + a) \geq m, \quad i = 1 : n \\ \|\mathbf{v}\|^2 = 1 \end{array} \right.$$

changement de variable : $\mathbf{w} = \frac{\mathbf{v}}{m}$ et $b = \frac{a}{m} \implies \|\mathbf{w}\| = \frac{\|\mathbf{v}\|}{m}$

$$\left\{ \begin{array}{l} \max_{\mathbf{w}, b} \quad m \\ \text{avec} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad ; \quad i = 1, n \\ \text{et} \quad m = \frac{1}{\|\mathbf{w}\|} \end{array} \right.$$

$$\left\{ \begin{array}{l} \min_{\mathbf{w}, b} \quad \|\mathbf{w}\|^2 \\ \text{avec} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ \quad \quad \quad i = 1, n \end{array} \right.$$

Le problème des SVM linéaires

Le problème des SVM linéaires (dans le primal)

soit $\{(\mathbf{x}_i, y_i); i = 1 : n\}$ un ensemble de données étiquetés avec $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{1, -1\}$.

Un séparateur à vaste marge linéaire (SVM) est un discriminateur de la forme : $D(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$ où $\mathbf{w} \in \mathbb{R}^d$ et $b \in \mathbb{R}$ sont donnés par la résolution du problème suivant :

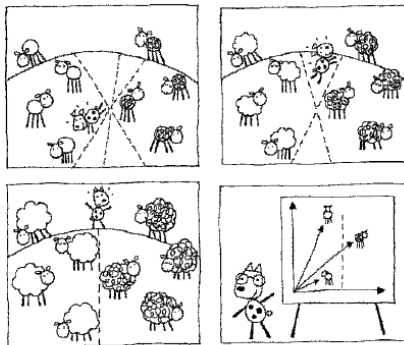
$$\begin{cases} \min_{\mathbf{w}, b} & \frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{avec} & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{cases}$$

C'est un programme quadratique de la forme

$$\begin{cases} \min_{\mathbf{z}} & \frac{1}{2} \mathbf{z}^\top \mathbf{A} \mathbf{z} - \mathbf{d}^\top \mathbf{z} \\ \text{avec} & \mathbf{B} \mathbf{z} \leq \mathbf{e} \end{cases}$$

$$\mathbf{z} = (\mathbf{w}, b)^\top, \mathbf{d} = (0, \dots, 0)^\top, \mathbf{A} = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{B} = -[\mathbf{y} \mathbf{X}, \mathbf{y}] \text{ et } \mathbf{e} = -(1, \dots, 1)^\top$$

Le problème de discrimination linéaire



...the story of the sheep dog who was herding his sheep, and serendipitously invented the large margin classification and Sheep Vectors ...

(drawing by Ana Martin Larranaga)

Formulation duale des SVM linéaires - Le lagrangien

$$\begin{cases} \min_{\mathbf{w}, b} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{avec} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{cases}$$

on recherche un point selle du lagrangien $\max_{\alpha} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha)$ avec des multiplicateurs de lagrange $\alpha_i \geq 0$

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

α_i traduit l'influence de la contrainte et donc l'influence du point (x_i, y_i)

Conditions d'optimalité

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

Les gradients :

$$\begin{cases} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) &= \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial \mathcal{L}(\mathbf{w}, b, \alpha)}{\partial b} &= \sum_{i=1}^n \alpha_i y_i \end{cases}$$

on écrit les conditions d'optimalité :

$$\begin{cases} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) = 0 &\Rightarrow \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial \mathcal{L}(\mathbf{w}, b, \alpha)}{\partial b} = 0 &\Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

Formulation duale des SVM linéaires

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

Optimalité : $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad \sum_{i=1}^n \alpha_i y_i = 0$

$$\begin{aligned} \mathcal{L}(\alpha) &= \frac{1}{2} \underbrace{\sum_{i=1}^n \sum_{j=1}^n \alpha_j \alpha_i y_i y_j \mathbf{x}_j^\top \mathbf{x}_i}_{\mathbf{w}^\top \mathbf{w}} - \sum_{i=1}^n \alpha_i y_i \underbrace{\sum_{j=1}^n \alpha_j y_j \mathbf{x}_j^\top \mathbf{x}_i}_{\mathbf{w}^\top \mathbf{x}_i} - b \underbrace{\sum_{i=1}^n \mathbf{x}_i y_i}_{=0} + \sum_{i=1}^n \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_j \alpha_i y_i y_j \mathbf{x}_j^\top \mathbf{x}_i + \sum_{i=1}^n \alpha_i \end{aligned}$$

Formulation duale des SVM linéaires

$$\begin{cases} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{avec} & \mathbf{y}^\top \alpha = 0 \\ \text{et} & 0 \leq \alpha_i \quad i = 1, n \end{cases}$$

avec G une matrice symétrique $n \times n$ de terme général $G_{ij} = y_i y_j \mathbf{x}_j^\top \mathbf{x}_i$

L'exemple des moindres carrés

Le modèle linéaire

$$y_i = \sum_{j=1}^d \beta_j x_{ij} + \varepsilon_i \quad , \quad i = 1, n$$

n observations et d variables; $d < n$

$$\min_{\beta} = \sum_{i=1}^n \left(\sum_{j=1}^d x_{ij} \beta_j - y_i \right)^2 = \|X\beta - Y\|^2$$

Solution: $\tilde{\beta} = (X^T X)^{-1} X^T Y$

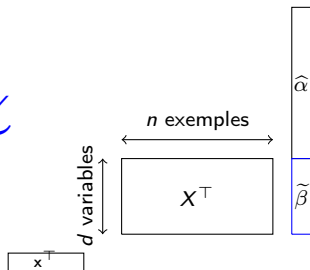
$$f(\mathbf{x}) = \mathbf{x}^T \underbrace{(X^T X)^{-1} X^T Y}_{\tilde{\beta}}$$

Quelle est l'influence de chacun des exemples (les lignes de X) ?

L'influence des exemples

pour une nouvelle observation \mathbf{x}

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{x}^\top (X^\top X)(X^\top X)^{-1} \underbrace{(X^\top X)^{-1} X^\top Y}_{\tilde{\beta}} \\ &= \mathbf{x}^\top X^\top \underbrace{X(X^\top X)^{-1} X^\top Y}_{\hat{\alpha}} \end{aligned}$$

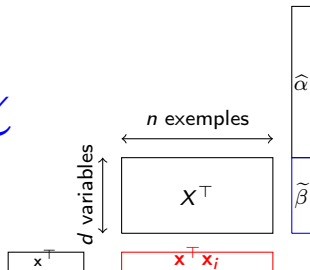


$$f(\mathbf{x}) = \sum_{j=1}^d \tilde{\beta}_j x_j$$

L'influence des exemples

pour une nouvelle observation \mathbf{x}

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{x}^T (X^T X)(X^T X)^{-1} \underbrace{(X^T X)^{-1} X^T Y}_{\tilde{\beta}} \\ &= \mathbf{x}^T X^T \underbrace{X(X^T X)^{-1} (X^T X)^{-1} X^T Y}_{\hat{\alpha}} \end{aligned}$$



$$f(\mathbf{x}) = \sum_{j=1}^d \tilde{\beta}_j x_j = \sum_{i=1}^n \hat{\alpha}_i (\mathbf{x}^T \mathbf{x}_i)$$

des variables aux exemples

$$\underbrace{\hat{\alpha} = X(X^T X)^{-1} \tilde{\beta}}_{n \text{ exemples}}$$

et

$$\underbrace{\tilde{\beta} = X^T \hat{\alpha}}_{d \text{ variables}}$$

et si $d \geq n$!

SVM primal vs. dual

Primal

$$\left\{ \begin{array}{ll} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{avec} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{array} \right.$$

- $d + 1$ inconnues
- n contraintes
- QP classique
- parfait si $d \ll n$

Dual

$$\left\{ \begin{array}{ll} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{avec} & \mathbf{y}^\top \alpha = 0 \\ \text{et} & 0 \leq \alpha_j \quad i = 1, n \end{array} \right.$$

- n inconnues
- G matrice des influences de chaque couple de points
- n contraintes de boîtes
- plus facile à résoudre
- à utiliser si $d > n$

SVM primal vs. dual

Primal

$$\left\{ \begin{array}{ll} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{avec} & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ & i = 1, n \end{array} \right.$$

- $d + 1$ inconnues
- n contraintes
- QP classique
- parfait si $d \ll n$

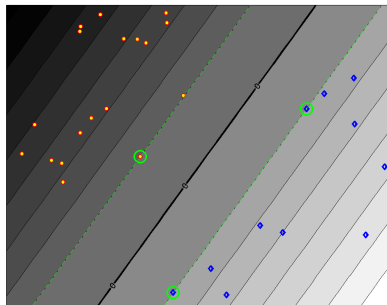
Dual

$$\left\{ \begin{array}{ll} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{avec} & \mathbf{y}^\top \alpha = 0 \\ \text{et} & 0 \leq \alpha_i \quad i = 1, n \end{array} \right.$$

- n inconnues
- G matrice des influences de chaque couple de points
- n contraintes de boîtes
- plus facile à résoudre
- à utiliser si $d > n$

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

Méthodes de résolution du problème dual

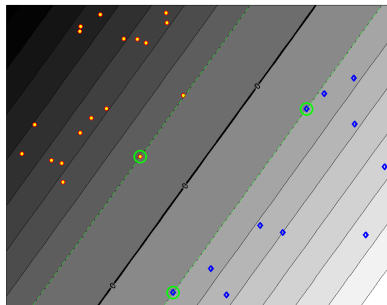


l'influence des points

- $\alpha_i = 0$ le point est inutile
- $\alpha_i \neq 0$ le point est dit **support**

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

Méthodes de résolution du problème dual



l'influence des points

- $\alpha_i = 0$ le point est inutile
- $\alpha_i \neq 0$ le point est dit **support**

$$f(\mathbf{x}) = \sum_{j=1}^d w_j x_j + b = \sum_{i=1}^3 \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

La fonction de décision ne dépend plus que de trois points ($d + 1$)

Supposons que nous connaissions les 3 points en question

Conclusion : variables ou exemples ?

- a la recherche d'une méthode d'apprentissage universelle
 - ▶ pas de modèle
- le cas linéaire - séparable
 - ▶ le cas non séparable est analogue
- double objectif : minimiser les erreur et la régularité de la solution
 - ▶ optimisation multi critère
- dualié : variable – exemple
 - ▶ quand utiliser l'une ou l'autre des formulations ?
- universalité = nonlinéarité
 - ▶ les noyaux