

CLASSIFICATION — LA FIN

PIPELINE D'APPRENTISSAGE ET SVM

Guillaume Wisniewski
guillaume.wisniewski@limsi.fr

Université Paris Sud — LIMSI

Mars 2016

ONE THEOREM TO RULE THEM ALL...

- Tout l'apprentissage repose sur un théorème :

$$\epsilon_{\text{test}} \leq \hat{\epsilon}_{\text{train}} + \sqrt{\frac{\text{complexity}}{n}} \quad (1)$$



- Lien entre l'erreur **estimée** sur un corpus de n exemples et l'erreur de généralisation (sur l'ensemble des données / non estimée)
- dépend de la **complexité** de la classe de fonctions considérée
- Avec des hypothèses très précises

INTERPRÉTATION



- il n'y pas de fonctions qui permet d'obtenir une meilleure erreur en généralisation que celle qui minimise l'erreur d'apprentissage...

INTERPRÉTATION



- il n'y pas de fonctions qui permet d'obtenir une meilleure erreur en généralisation que celle qui minimise l'erreur d'apprentissage...
- ...avec une forte probabilité...

INTERPRÉTATION



- il n'y pas de fonctions qui permet d'obtenir une meilleure erreur en généralisation que celle qui minimise l'erreur d'apprentissage...
- ...avec une forte probabilité...
- ...à l'intérieur d'une classe de fonctions données...

INTERPRÉTATION



- il n'y pas de fonctions qui permet d'obtenir une meilleure erreur en généralisation que celle qui minimise l'erreur d'apprentissage...
- ...avec une forte probabilité...
- ...à l'intérieur d'une classe de fonctions données...
- ...si les exemples d'apprentissage sont représentatifs.

1^{ER} MESSAGE : SMALL IS BEAUTIFUL



Apprendre, c'est minimiser

PARCE QUE LE MONDE EST BIEN FAIT...



- classification binaire avec $\ell^{0/1}$
- minimiser l'erreur sur un ensemble d'apprentissage d'apprentissage \Rightarrow problème NP-difficile
- optimisation directe impossible (ou pas)

<http://jmlr.org/proceedings/papers/v28/nguyen13a.pdf>

Première partie I

RETOUR SUR LE PERCEPTRON

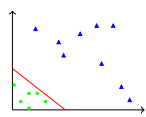
PERCEPTRON



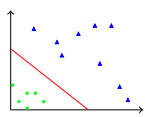
- Méthode **heuristique** d'optimisation de la fonction de $\ell^{0/1}$
- aucune garantie sur la convergence (sauf si les données sont linéairement séparable)
- aucune garantie sur la qualité de la solution trouvée

LES LIMITES DU PERCEPTRON

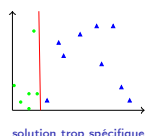
Quelle garantie en généralisation ?



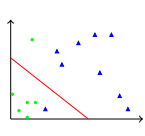
mauvaise solution ?



meilleure solution ?



solution trop spécifique



solution préférable

GRÂCES SOIENT RENDUES AUX MATHÉMATIQUES

- aucun lien théorique entre erreur d'apprentissage et erreur de généralisation
- solution pragmatique : estimer expérimentalement l'erreur en généralisation
- en pratique, 3 ensembles indépendants :
 - 1 corpus d'apprentissage : permet de déterminer les paramètres du classifieur
 - 2 corpus de validation / développement : permet de monitorer l'erreur en généralisation :
 - ★ régulièrement, on calcule l'erreur sur l'ensemble de validation
 - ★ on conserve la valeur des paramètres qui minimise l'erreur de validation
 - 3 corpus de test : permet d'estimer l'erreur de généralisation (1 seule fois)

POURQUOI ÇA MARCHE ?



- estimation **non biaisée** de l'erreur de généralisation (cf. théorème central limite)
- l'erreur est une moyenne \Rightarrow peut-être estimée à partir d'un échantillon suffisamment grand
- mais nécessite plus de données annotées et un surcoût en temps de calcul

PROBLÈME DE STABILITÉ DU PERCEPTRON

LE PROBLÈME

- 10 000 exemples, après le 100^e le perceptron est tellement bon qu'il ne fait pas d'erreurs sur les 9 899 exemples suivants...
- erreur sur le 10 000^e exemple \Rightarrow la mise à jour va détruire le vecteur de poids qui avait parfaitement marcher sur 99.99% des exemples

LA SOLUTION

- `random.shuffle`
- perceptron moyenné

PERCEPTRON MOYENNÉ

- Attention : pendant l'apprentissage on utilise la version « non moyennée » du perceptron

CONCLUSIONS (LOCALES)

- peu de garanties théoriques...
- ...mais on sait faire marcher en pratique

Deuxième partie II

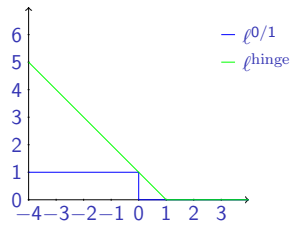
SUPPORT VECTOR MACHINE

NOTRE PROBLÈME

$$\min_{\mathbf{w}} \sum \mathbb{1} \{f(\mathbf{x}_i; \mathbf{w}) \neq y_i\} \quad (2)$$

- problème NP-difficile

2^E MESSAGE : LA NATURE DE L'APPRENTISSAGE STATISTIQUE



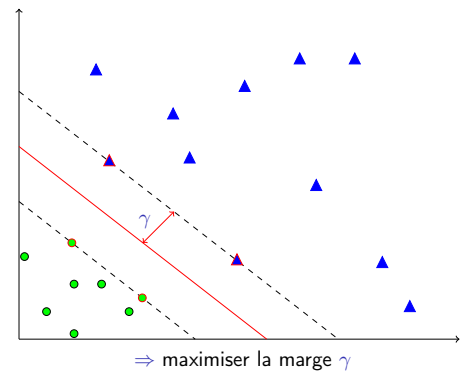
- considérer la hinge loss au lieu de $\ell^{0/1}$:

$$\ell^{\text{hinge}} = \max(0, 1 - t \cdot y) \quad (3)$$

où : t est l'étiquette de référence et y l'étiquette prédite ($y = \mathbf{x} \cdot \mathbf{w}$)

- minimiser la hinge revient à minimiser le $\ell^{0/1}$.

INTUITION DU SVM



MAXIMISER LA MARGE

POURQUOI LA MARGE ?

- trouver l'hyperplan séparateur le « plus loin » des données
- évite que les données inconnues ne soient trop proches de la frontière de décision

FORMALISATION : MARGE

- hyperplan séparateur : $\mathbf{w} \cdot \mathbf{x} = 0$
- hyperplan support : hyperplans parallèles à \mathcal{H} passant par un exemple positif/négatifs et le plus proche de \mathcal{H} :

$$\mathbf{w} \cdot \mathbf{x} = \pm b \Leftrightarrow \mathbf{w} \cdot \mathbf{x} = \pm 1$$
- marge : distance entre \mathcal{H} et l'hyperplan support :

$$\gamma = \frac{2}{\|\mathbf{w}\|}$$

avec de simples considérations géométriques

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :
 - s'il n'y a pas de contraintes : solution triviale

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :
 - ▶ s'il n'y a pas de contraintes : solution triviale
 - ▶ tous les exemples sont du bon côté de leur hyperplan **support**

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :
 - ▶ s'il n'y a pas de contraintes : solution triviale
 - ▶ tous les exemples sont du bon côté de leur hyperplan **support**
 - ▶ pour les exemples positifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :
 - ▶ s'il n'y a pas de contraintes : solution triviale
 - ▶ tous les exemples sont du bon côté de leur hyperplan **support**
 - ▶ pour les exemples positifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$
 - ▶ pour les exemples négatifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \leq -1$

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :
 - ▶ s'il n'y a pas de contraintes : solution triviale
 - ▶ tous les exemples sont du bon côté de leur hyperplan **support**
 - ▶ pour les exemples positifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$
 - ▶ pour les exemples négatifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \leq -1$
 - ▶ de manière compacte : $y^{(i)} \cdot (\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$

CRITÈRE D'APPRENTISSAGE

- 1 maximiser la marge : $\max \frac{2}{\|\mathbf{w}\|} = \min \|\mathbf{w}\|^2$ le carré permet juste de simplifier les équations
- 2 sous les contraintes :
 - ▶ s'il n'y a pas de contraintes : solution triviale
 - ▶ tous les exemples sont du bon côté de leur hyperplan **support**
 - ▶ pour les exemples positifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$
 - ▶ pour les exemples négatifs : $(\mathbf{w} \cdot \mathbf{x}^{(i)}) \leq -1$
 - ▶ de manière compacte : $y^{(i)} \cdot (\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$
 - ▶ intuitivement : garanti que l'exemple est bien classé avec une **marge** de 1 cf. perceptron

BILAN : SUPPORT VECTOR MACHINE

$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2$$

$$\text{s.t. } y^{(i)} \cdot (\mathbf{w} \cdot \mathbf{x}^{(i)}) \geq 1$$

- problème d'optimisation quadratique sous contraintes
- objectif convexe
- problème « simple » d'un point de vue mathématique
- 20 ans de recherche pour avoir une implémentation utilisable
- aujourd'hui : plusieurs bibliothèques proposent des implémentations performantes

ET POUR LES DONNÉES NON SÉPARABLES ?



- données non-séparables : \nexists hyperplan séparateur
 - dans ce cas, le problème précédent n'a pas de solution
 - il y a forcément au moins une contrainte violée
 - critère relâché : maximiser la marge avec le moins de contraintes violées
- ⇒ SVM à marge molle

FORMALISATION

- pour chaque contrainte, on introduit une variable ressort :

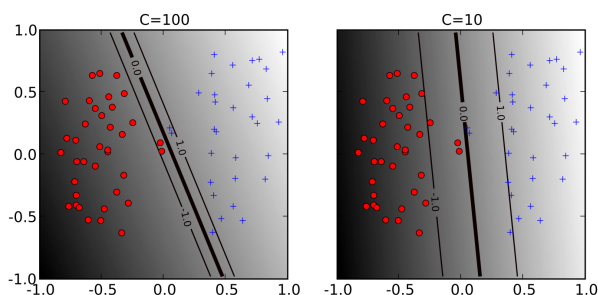
$$y^{(i)} \cdot (w \cdot x^{(i)}) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$

- ξ_i = de combien la i^e contrainte est violée
- nouveau objectif : maximiser la marge en minimisant le non-respect des contraintes :

$$\min ||w||^2 + C \times \sum_{i=1}^n \xi_i$$

- C = compromis entre les deux critères = hyper-paramètres

IMPACT DE C



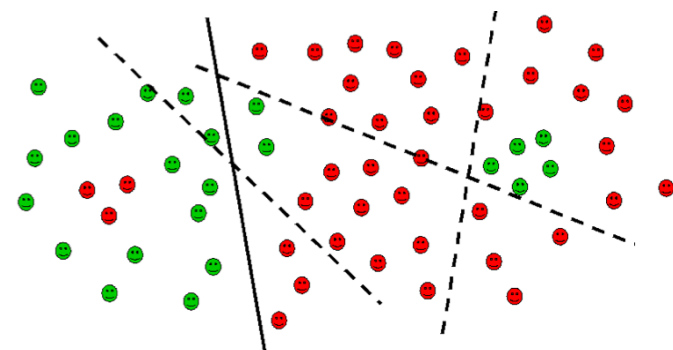
COMMENT CHOISIR C

- question très importante : c'est la différence entre un SVM qui marche et un SVM qui ne marche pas
- sujet encore hautement spéculatif
- en pratique : recherche par force brute :
 - ▶ on considère un ensemble de valeurs (en général 10^i pour $i \in [-5, 3]$)
 - ▶ on garde le C qui permet d'obtenir les meilleures performances sur un corpus de validation

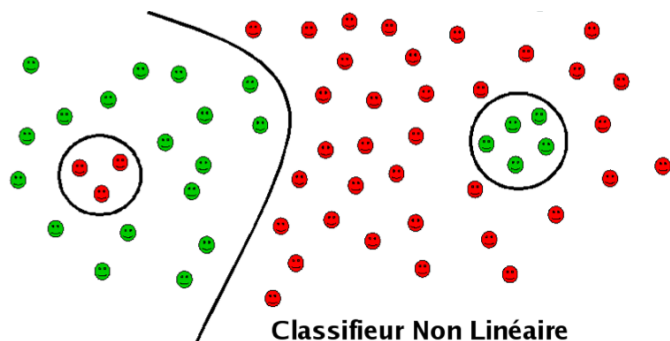
Troisième partie III

NOYAUX

RAPPEL (1) : CLASSIFIEURS LINÉAIRES



RAPPEL (2) : CLASSIFIEURS NON LINÉAIRES



QUESTION LOCALE



- Peut-on généraliser un classifieur linéaire pour apprendre une frontière non linéaire ?
- **oui** → noyaux
- en général associé avec les SVM, mais peut s'appliquer à tous les classifieurs linéaires

RAPPEL (3) : PERCEPTRON

```

w ← 0
while il y a des exemples mal classés do
  choisir un exemple (xi, yi) au hasard
  if yi · w⊤ xi ≤ 0 then
    w ← w + yi · φ(xi)
  end if
end while
    
```

Representer Theorem

Observation : Toute solution trouvée par l'algorithme du perceptron est une combinaison linéaire des exemples mal classés.
(également vrai pour le SVM)

Notations :

- $w = \sum_{i=1}^n \alpha_i \cdot y_i \cdot x_i$
- \Rightarrow représentation alternative du classifieur
- théorie de l'optimisation : α_i = coordonnées duales

PERCEPTRON : FORME DUALE

Nouvelle fonction de décision :

$$f(x) = \text{sign}(\langle x | w \rangle)$$

$$= \text{sign} \left(\left\langle \sum_{j=1}^n \alpha_j \cdot y_j \cdot x_j | x \right\rangle \right)$$

PERCEPTRON : PROBLÈME D'OPTIMISATION DUALE

Require: a linearly separable training set $(x_i, y_i)_{i=1}^n$

```

α ← 0
R = maxi ∈ [1, n] || xi ||
while there are classification errors do
  for i = 1 to n do
    y = ∑j=1n αj · yj · ⟨xj | xi⟩
    if y ≠ yi then
      αi ← αi + 1
    end if
  end for
end while
    
```

INTÉRÊT

- dans la fonction de décision et dans l'algorithme d'optimisation, on accède aux exemples uniquement pour calculer des produits scalaires
- algorithme « rapide » si on sait calculer un produit scalaire rapidement
- et surtout : pas de dépendance quand à la dimension des exemples

LES NOYAUX (1)

Un noyau :

- fonction $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \dots$
- tel qu'il existe un espace de Hilbert \mathcal{H} et une fonction $\phi : \mathcal{X} \mapsto \mathcal{H}$, avec $K(x, x') = \langle \phi(x) | \phi(x') \rangle$,

Exemple : noyau polynomial

$$K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(x, x') \mapsto (\langle x | x' \rangle + 1)^m$$

LES NOYAUX (2)

Supposons : $x = (x_1, x_2)$

$$K(x, x') = (\langle x | x' \rangle + 1)^2 = (x_1 \cdot x'_1 + x_2 \cdot x'_2 + 1)^2$$

$$= x_1^2 \cdot x_1'^2 + (\sqrt{2} \cdot x_1) \cdot (\sqrt{2} \cdot x'_1) + (\sqrt{2} \cdot x_1 \cdot x_2) \cdot (\sqrt{2} \cdot x'_1 \cdot x'_2)$$

$$+ (\sqrt{2} \cdot x_2) \cdot (\sqrt{2} \cdot x'_2) + x_2^2 \cdot x_2'^2 + 1$$

$$= \langle \phi(x) | \phi(x') \rangle$$

avec :

$$\phi : (x_1, x_2) \mapsto \begin{pmatrix} x_1^2 \\ \sqrt{2} \cdot x_1 \\ \sqrt{2} \cdot x_1 \cdot x_2 \\ \sqrt{2} \cdot x_2 \\ x_2^2 \\ 1 \end{pmatrix}$$

BILAN

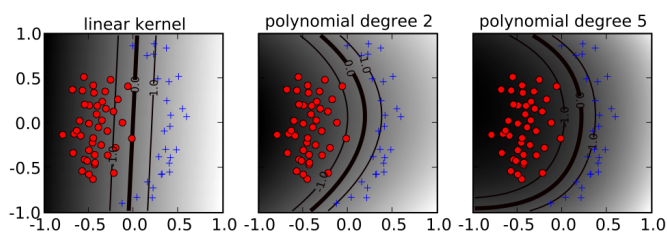
INTÉRÊTS

- 1 réduction de la complexité : avec un noyau $\mathcal{O}(d)$ au lieu de $\mathcal{O}\left(\binom{d+m-1}{m}\right)$
- 2 possibilité d'utiliser des transformations non linéaires :
 - noyaux gaussien $K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2 \cdot \gamma^2}\right)$

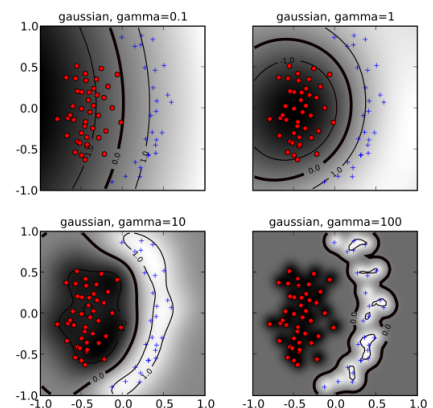
MORALITÉ

Noyaux = moyens de réaliser une transformation non-linéaire **implicite**

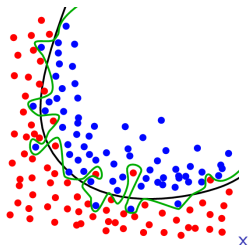
EXEMPLE



EXEMPLE (2)



LES NOYAUX UNE SOLUTION MAGIQUE ?



- les noyaux permettent de modéliser des frontières de décisions complexes
- mais : risque fort de sur-apprentissage
- coût computationnel élevé
- peu utile en grande dimension