Réseaux de Neurones, Mélanges de classifieurs, SVM avancé

### Réseau Artificiel de Neurones

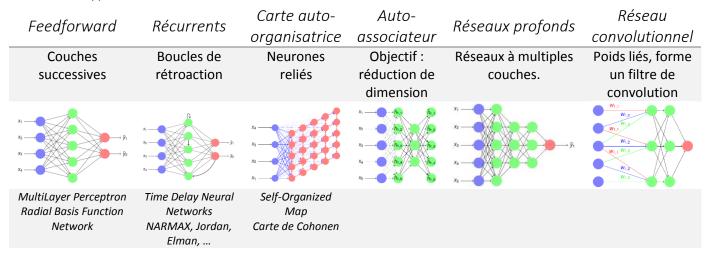
#### 1. Neurone

Exemples de choix pour la fonction f:

$$\hat{y} = f(w^{\mathsf{T}}x + b)$$

- $sigm(x) \in [0,1]$
- $tanh(x) \in [-1; 1]$
- $\operatorname{softmax}(x) \in [-1; 1]$  suivant une loi de probabilité

#### 2. Type de réseaux



Les réseaux profonds nécessitent qu'on initialise les poids pour stabiliser les calculs. Pour cela, on préapprend les couches basses (vertes) avec des auto-associateurs (cf ex d'auto-associateur).

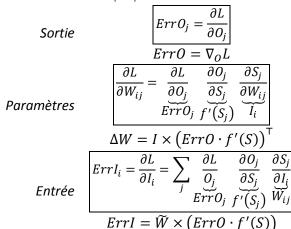
#### 3. Perceptron Multi-Couches

#### a. Définition d'une couche

Pour une couche, on a n entrées notées i, p sorties notées j.

$$I = \begin{bmatrix} I_0 = 1 \\ I_1 \\ \vdots \\ I_i \\ \vdots \\ I_n \end{bmatrix} \qquad W = \begin{bmatrix} w_1 & \dots & w_j & \dots & w_p \\ \end{bmatrix}$$
 
$$w_j = \begin{bmatrix} w_{0j} & \dots & w_{nj} \end{bmatrix}^{\mathsf{T}}$$
 
$$\widetilde{W} = W(1:n,1:p)$$
 Poids

### b. Principe pour une couche



$$\begin{bmatrix} W = \begin{bmatrix} w_1 & \dots & w_j & \dots & w_p \end{bmatrix} \\ w_j = \begin{bmatrix} w_{0j} & \dots & w_{nj} \end{bmatrix}^{\mathsf{T}} \\ \widetilde{W} = W(1:n,1:p) \\ Poids \end{bmatrix} S = \begin{bmatrix} S_1 \\ \vdots \\ S_j = w_j^{\mathsf{T}}I \\ \vdots \\ S_p \end{bmatrix} = W^{\mathsf{T}}I \\ \vdots \\ S_p \end{bmatrix} = W^{\mathsf{T}}I$$

$$O = \begin{bmatrix} O_1 \\ \vdots \\ O_j = f(S_j) \\ \vdots \\ O_p \end{bmatrix} = f(S)$$

$$Sommage$$

$$Sortie$$

#### c. Calcul pour *m* couches et rétropropagation

La technique consiste à choisir une fonction L pour la sortie finale (ex :  $L(O,Y) = ||O-Y||^2$  où Y est l'objectif de sortie) et de calculer  $ErrO_j^{(m)} = \frac{\partial L}{\partial O_i^{(m)}}$ 

Pour les autres couches (l < m), on applique la rétropropagation du gradient, soit  $\overline{\mathit{ErrO}_i^{(l)} = \mathit{ErrI}_i^{(l+1)}}$ 

Réseaux de Neurones, Mélanges de classifieurs, SVM avancé

## II. Mélange de classifieurs

#### 1. Généralités sur les mélanges

**Principe** Consiste à combiner les décisions de plusieurs classifieurs.

#### Limitations d'un classifieur unique

Limitation	roblème du classifieur seul Avantage du mélange	
	Il existe plusieurs classifieurs avec les mêmes performances en test.	Moyenner plusieurs classifieurs aussi performants.
Représentation (biais)	On peut créer un nombre limité de classifieurs.	Obtenir un classifieur qu'on ne pourrait pas construire directement.
Computationnelle	On tombe souvent dans des minimas locaux.	S'approcher du minimum global.

**Types de mélanges** Hétérogènes (classifieurs de types ≠) / Homogènes (classifieurs de même type)

Approches Injecter de l'aléatoire / Manipuler les données (apprentissage, entrée, sortie) / ...

#### 2. Bootstrap

Méthode de ré-échantillonnage permettant d'estimer une grandeur s(X) d'un jeu initial X de n valeurs en créant B jeux de données  $\tilde{X}_k$  de m valeurs par tirage aléatoire avec remise sur le jeu.

$$X = \{x_1, \dots, x_i, \dots x_n\} \rightarrow \begin{cases} \tilde{X}_1 = \{x_4, x_6, x_4, x_1, x_3, \dots\} \\ \vdots \\ \tilde{X}_B = \{x_2, x_1, x_n, x_2, x_2, \dots\} \end{cases}$$

$$s(X) = \sum_{k=1}^{B} s(\tilde{X}_k) / B \qquad \sigma_s = \sum_{k=1}^{B} \left( s(\tilde{X}_k) - s(X) \right)^2 / (B - 1)$$

#### 3. Bagging

#### a. Principe

Le bagging applique l'idée du bootstrap à la classification. A partir du jeu initial X, on crée B (souvent 50) jeux par tirage bootstrap (souvent m=n) sur lesquels on apprends des classifieurs que l'on combine finalement (souvent vote majoritaire).

#### b. Out-Of-Bag

L'OOB d'un jeu  $\tilde{X}_k$  est l'ensemble de valeurs  $x_i$  de X qui ne sont pas dans  $\tilde{X}_k$  (en moyenne 37%).

Cet ensemble peut servir comme jeu de test lors de l'apprentissage sur le jeu  $\tilde{X}_k$ .

Réseaux de Neurones, Mélanges de classifieurs, SVM avancé

#### 4. Boosting

#### a. Principe

Apprentissage itératif/adaptatif d'un ensemble de classifieurs par pondération de chaque donnée en fonction de son erreur par les classifieurs précédents.

### b. Algorithme d'apprentissage AdaBoost

$$\begin{split} D_1(x_i) &= \frac{1}{m} \quad \forall \ i \\ \text{pour } t = 1, \dots, T \text{ faire} \\ h_t &= \mathcal{L}(X, D_t) \\ \widehat{e_t} &= \sum_i D_t(x_i) \\ \alpha_t &= \frac{1}{2} \ln \frac{1 - \widehat{e_t}}{\widehat{e_t}} \\ D_{t+1}(x_i) &= \frac{D_t(x_i) \exp{-\alpha_t y_i h_t(x_i)}}{\sum_i D_{t+1}(x_i)} \quad \forall \ i \quad \text{Mise à jour des poids des données} \end{split}$$

fin pour

#### c. Classification

On fait de la classification par vote pondéré, c'est-à-dire :

$$\hat{y} = signe\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

#### 5. Foret aléatoires

#### a. Arbres de décision

- 1. Créer un premier groupe avec l'ensemble de données
- 2. Apprendre une règle pour séparer le groupe en 2 sous-groupes les plus purs possibles
- 3. Répéter l'opération 2 sur chaque sous-groupe, sauf si on a atteint le critère d'arrêt

On peut prendre par exemple comme critères d'arrêts une pureté minimale à atteindre et/ou une hauteur maximale de l'arbre.

#### b. Principe de la forêt aléatoire

Apprendre un ensemble d'arbres de décision à partir de données différentes  $\tilde{X}_k$  mais suivant la même distribution statistique et combiner leurs décisions.

#### c. Construction des jeux de données

Pour augmenter les performances des forêts, on ajoute de l'aléatoire dans la création du jeu d'apprentissage de chaque arbre  $\tilde{X}_k$  à partir des données X.

- Création des jeux  $\tilde{X}_k$  par bootstrap
- Construction des arbres par rapport à une seule variable (Random Feature Selection TODO a revoir...)
- Sélection d'un sous-ensemble de  $\tilde{p}$  variables parmi les p variables de X

#### d. Décision de la forêt

- Vote majoritaire simple des décisions des arbres
- Pondération des arbres par leur performance sur l'OOB

Réseaux de Neurones, Mélanges de classifieurs, SVM avancé

### III. SVM et noyaux

#### 1. Rappel : le SVM linéaire

a. Problème primal

$$X_{app} \in \mathbb{R}^{n \times p}$$

$$y_{app} \in \{1; -1\}^{n} \Rightarrow \begin{bmatrix} \min_{w, \xi_{i}} \frac{1}{2} ||w||^{2} + C \sum_{i=1}^{n} \xi_{i} \\ y_{i}(w^{\mathsf{T}}x_{i} + b) \geq 1 - \xi_{i} \ \forall i \end{bmatrix} \Rightarrow \begin{bmatrix} f(x) = w^{\mathsf{T}}x + b \\ \hat{y} = sign(f(x)) \end{bmatrix}$$

Données initiales

#### b. Problème dual

Lagrangien du problème primal

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i(w^\top x + b) - 1 + \xi_i) - \sum_{i=1}^n \beta_i \xi_i$$

Focus sur la complémentarité : cas possibles de classement des points

#### Problème dual

$$\begin{cases} \min_{\alpha} \frac{1}{2} \alpha^{\mathsf{T}} G \alpha - \alpha^{\mathsf{T}} e \\ 0 \le \alpha_i \le C \\ \alpha^{\mathsf{T}} y = 0 \end{cases} \text{ avec } \begin{cases} G_{ij} = x_i^{\mathsf{T}} x_j y_i y_j \\ \widehat{y} = sign(f(x)) \end{cases} \Rightarrow \begin{cases} w = \sum_{i=1}^n \alpha_i y_i x_i \\ f(x) = \sum_{i=1}^n \alpha_i y_i x_i^{\mathsf{T}} x + b \\ \widehat{y} = sign(f(x)) \end{cases}$$

Problème SVM linéaire dual

Fonction de décision

#### 2. Passage au problème non-linéaire : novau

On se propose de rendre linéaire un problème qui ne l'est pas en changeant d'espace via une feature map  $\phi$  passant des  $x_i$  aux  $t_i$ , on utilisera donc les  $t_i$  dans le problème.

L'espace d'arrivée est très grand. Dans le primal, le nombre de variable explose. Dans le dual, on a toujours autant de variables mais il faut calculer les produits scalaires entre les  $t_i$ .

On définit le kernel permettant de calculer  $t_i^{\mathsf{T}}t_i$  à partir de  $x_i, x_i$ .

En pratique donc, on ne calcule pas les  $t_i$ , on intègre la fonction kernel au problème quand on a besoin des produits scalaires des  $t_i$ .

Le problème et la fonction de décision deviennent donc :

$$\begin{cases} \min \frac{1}{f, b, \xi_i} \frac{1}{2} \|f\|^2 + C \sum_{i=1}^n \xi_i \\ y_i(f(x_i) + b) \ge 1 - \xi_i \\ \xi_i \ge 0 \end{cases} \Rightarrow \begin{cases} \min \frac{1}{2} \alpha^\mathsf{T} G \alpha - \alpha^\mathsf{T} e \\ 0 \le \alpha_i \le C \\ \alpha^\mathsf{T} y = 0 \end{cases} \Rightarrow \begin{cases} w = \sum_{i=1}^n \alpha_i y_i t_i \\ f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

$$\Rightarrow \begin{cases} f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b \\ \hat{y} = sign(f(x)) \end{cases}$$

Exemple de  $\phi$  polynomiale :

$$b: \mathbb{R}^2 \to \mathbb{R}^p$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \to t = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ \vdots \end{bmatrix}$$

#### Exemple de kernels :

$$k_{polynomial}(x_i, x_j) = (x_i^{\mathsf{T}} x_j + 1)^m$$
  
 $k_{gaussien}(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$ 

$$w = \sum_{i=1}^{n} \alpha_i y_i t_i$$
$$f(x) = \sum_{i=1}^{n} \alpha_i y_i k(x_i, x) + b$$
$$\hat{y} = sign(f(x))$$

Fonction de décision

# DATA MINING 2

Réseaux de Neurones, Mélanges de classifieurs, SVM avancé

## IV. Choix des paramètres par validation croisée

1. Estimation de l'erreur par validation croisée

$$X \rightarrow \{X_k\} \rightarrow \forall k, \begin{cases} X_k \text{ app } \overline{X_k} \text{ test} \\ \overline{X_k} \text{ app } X_k \text{ test} \end{cases} \rightarrow e_k \rightarrow e_k \rightarrow e_k$$

$$Erreur: e \pm t_{\alpha/2,K-1} \sqrt{\frac{\sigma^2}{K}}$$

### 2. Optimisation des hyperparamètres

Pour optimiser P paramètres, on génère une grille de dimension p et on calcule l'erreur pour chaque point de la grille. On peut répéter l'opération sur plusieurs grilles de plus en plus précises.

### V. Support Vector Data Description

1. Minimum enclosing ball problem

$$\begin{cases} \min_{c,R,\xi} R^2 + C \sum_{i=1}^n \xi_i \\ \|x_i - c\|^2 \le R^2 + \xi_i \\ \xi_i \ge 0 \end{cases}$$

Problème primal

$$\begin{bmatrix} \min_{\alpha} \alpha^{\mathsf{T}} G \alpha - \alpha^{\mathsf{T}} \operatorname{diag} G \\ \alpha^{\mathsf{T}} e = 1 \\ 0 \le \alpha_i \le C \end{bmatrix} \quad \text{avec } G_{ij} = x_i^{\mathsf{T}} x_j$$

$$R^2 = \mu + \alpha^{\mathsf{T}} G \alpha = \mu + ||c||^2$$
Problème dual

2. Ajout d'un kernel

$$\begin{cases} \min_{c,R,\xi} R^2 + C \sum_{i=1}^n \xi_i \\ \|k(\cdot,x_i) - c(\cdot)\|^2 \le R^2 + \xi_i \\ \xi_i \ge 0 \end{cases}$$

Problème SVDD primal

$$\begin{bmatrix} \min_{\alpha} \alpha^{\mathsf{T}} G \alpha - \alpha^{\mathsf{T}} \operatorname{diag} G \\ \alpha^{\mathsf{T}} e = 1 \\ 0 \le \alpha_i \le C \end{bmatrix} \quad \text{avec } G_{ij} = k(x_i, x_j)$$

Problème SVDD dual

$$c(\cdot) = \sum_{i=1}^{n} \alpha_i k(\cdot, x_i) \qquad \qquad f(x) = \|k(x, \cdot) - c(\cdot)\|^2 - R^2 = -2 \sum_{i=1}^{n} \alpha_i k(x, x_i) + k(x, x) - \mu$$
 Fonction de décision

#### VI.SVM multi-classe

Approd	che	Problem size	Nb of sub-problems	Total size	Discrimination	Rejection
All togeti	her	$n \times c$	1	$n \times c$	+	
1 vs.	all	n	С	$n \times c$	++	-
Décomposition 1 vs	s. 1	2n/c	c(c-1)/2	$n \times (c-1)$	++	-
c SV	DD	n/c	С	n	-	++
Couplage	СН	$\mid n \mid$	1	n	+	+

$$\begin{array}{l} \text{All together:} & \begin{cases} \min\limits_{w,\xi_i} \frac{1}{2} \sum_{\ell=1}^c \|w_\ell\|^2 \\ x_i^\top \big(w_{y_i} - w_\ell\big) + b_{y_i} - b_\ell \geq 1 \quad \forall \ i,\ell \neq y_i \end{cases} \end{array}$$

Réseaux de Neurones, Mélanges de classifieurs, SVM avancé

#### Multi kernel SVM VII.

#### 1. Principe

$$\boxed{K(\cdot,x_i) = \sum_{m=1}^M d_m k_m(\cdot,x_i)} \quad \text{avec} \quad \sum_{m=1}^M d_m = 1 \quad \text{et} \quad 0 \leq d_m \qquad \qquad f(\cdot) = \sum_m \underbrace{d_m \sum_i \alpha_i y_i k_m(\cdot,x_i)}_{f_m(\cdot)}$$

Mélange de kernels : combinaison linéaire de kernels

Fonction de décision

### 2. Problème d'optimisation

$$\begin{cases} \min_{f_m, b, \xi, d} \frac{1}{2} \sum_{m} \frac{1}{d_m} \|f_m\|^2 + C \sum_{i=1}^n \xi_i \\ y_i \left( \sum_{m} f_m(x_i) + b \right) \ge 1 - \xi_i \\ \xi_i \ge 0 \quad ; \quad \sum_{m} d_m = 1 \quad ; \quad d_m \ge 0 \end{cases}$$

$$\begin{cases} \min_{f_m,b,\xi,d} \frac{1}{2} \sum_{m} \frac{1}{d_m} \|f_m\|^2 + C \sum_{i=1}^n \xi_i \\ y_i \left( \sum_{m} f_m(x_i) + b \right) \ge 1 - \xi_i \\ \xi_i \ge 0 \quad ; \quad \sum_{m} d_m = 1 \quad ; \quad d_m \ge 0 \end{cases}$$

$$\begin{cases} \min_{d} J(d) = \begin{cases} \min_{f_m,b,\xi} \frac{1}{2} \sum_{m} \frac{1}{d_m} \|f_m\|^2 + C \sum_{i=1}^n \xi_i \\ y_i \left( \sum_{m} f_m(x_i) + b \right) \ge 1 - \xi_i \\ \xi_i \ge 0 \end{cases}$$

$$\sum_{m} d_m = 1 \quad ; \quad d_m \ge 0$$

Problème MKL primal bi-niveaux

#### 3. Résolution

On résout le problème par descente de gradient. A chaque itération on résout I(d) et on calcule son gradient  $\nabla_d I$ , on en déduit la direction de descente D et le pas optimal dans cette direction  $\gamma$ .

$$G = \sum_{m} d_{m}G_{m} \qquad G_{m,ij} = k_{m}(x_{i}, x_{j}) \qquad \Rightarrow \qquad \nabla_{d_{m}}J = \frac{1}{2}\alpha^{\top}G_{m}\alpha \qquad \Rightarrow \qquad \sum_{m} d_{m} = 1 \Rightarrow \sum_{m} D_{m} = 0$$

$$\Rightarrow \begin{bmatrix} D_{m} = \nabla_{d_{m}}J - \nabla_{d_{1}}J \\ D_{1} = -\sum_{m=2}^{M}D_{m} \end{bmatrix}$$

$$Cout à l'optimum \qquad Gradient \qquad Direction de descente$$

Gradient