

APPRENTISSAGE PAR COMBINAISON DE CLASSIFIEURS ELEMENTAIRES (« dopage » ou « Boosting »)

Pr. Fabien Moutarde
Centre de Robotique (CAOR)
MINES Paris Tech (ENSMP)
PSL Research University

Fabien.Moutarde@mines-paristech.fr
<http://people.mines-paristech.fr/fabien.moutarde>

Apprentissage par dopage (« boosting ») Pr. Fabien Moutarde, CAOR, MINES ParisTech, PSL Fév.2017 1



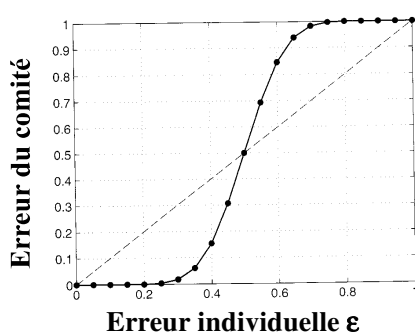
Principe essentiel : l'union fait la force

Réunir un « comité d'experts »
chacun peut se tromper, mais en combinant les avis,
on a plus de chance d'avoir la bonne prédiction !!...

Justification théorique :

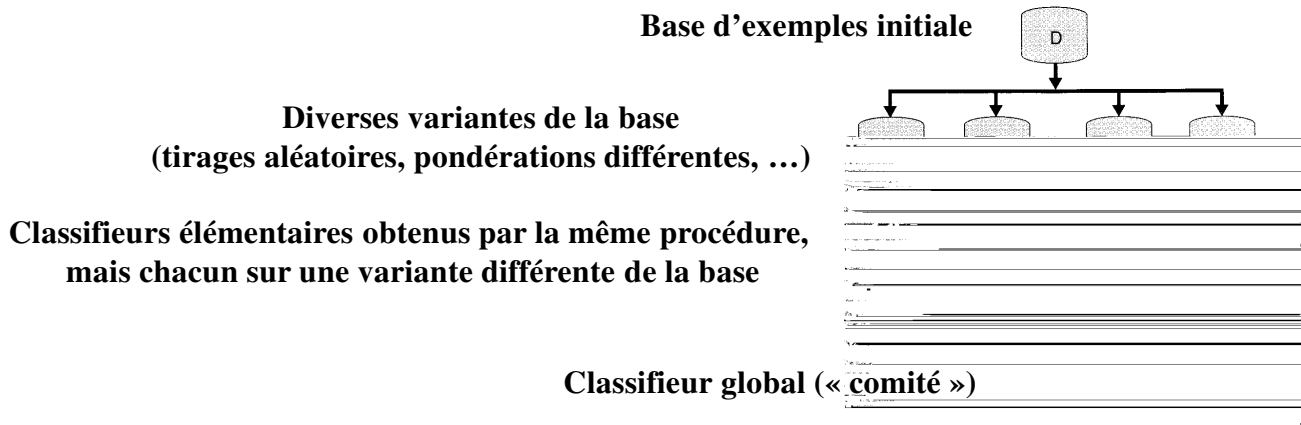
- supposons N classifieurs *indépendants*, faisant chacun une erreur $E_{gen} = \epsilon$
- si on décide « à la majorité », alors on se trompe si et seulement si plus de la moitié du « comité » se trompe

$$\rightarrow Erreur_{ensemble} = \sum_{k=N/2}^N C_k^N \epsilon^k (1-\epsilon)^{N-k}$$



Décision fortement améliorée,
(sous réserve que $\epsilon < 0.5$!!)...
...et ce d'autant plus qu'on augmente
le nombre N d'experts

- Utiliser des algos totalement différents
- Même algo mais avec des paramètres et/ou initialisations différent(e)s
- **Faire varier l'ensemble d'apprentissage**



→ Méthodes très générales, applicables
avec n'importe quel algorithme « élémentaire »

Le « bagging »

Variantes de la base d'apprentissage obtenues par tirages aléatoires avec remise depuis la base initiale

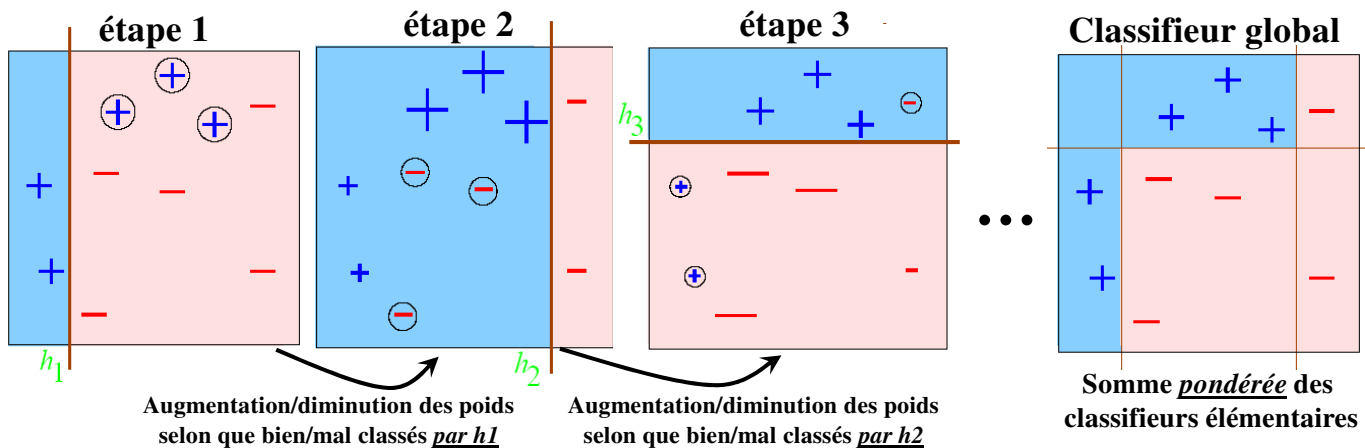
(sorte de « bootstrap » → duplication/disparition aléatoires de certains exemples selon variantes)

- Utile et efficace en particulier si l'algo de base utilisé est « instable » (au sens « sensible au bruit des données ») car différences entre variantes de base → classifieurs élémentaires très différents
- Evite « over-fitting » (sur-apprentissage), car on « moyenne » des classifieurs construits avec différentes réalisations aléatoires des mêmes données

Méthode itérative pour l'ajout des classifieurs :

variantes de la base d'apprentissage obtenues par des pondérations successives des mêmes exemples

(calculées pour « se focaliser » sur les exemples « difficiles », i.e. mal classés par plusieurs classifieurs déjà assemblés)



L'algorithme adaBoost

adaBoost (« adaptive Boosting »)

- Base d'exemples initiale : $S = \{ (x_1, u_1), \dots, (x_k, u_k) \}$, avec $u_i \in \{+1, -1\}$, $i=1, k$
- Poids initiaux : $w_0(x_i) = 1/m$ pour tout $i=1, k$ (ou $1/2p$ pr pos et $1/2n$ pr neg)
- Pour chaque étape (« round ») t de 1 à T , faire :
 1. apprendre/choisir 1 règle de classification h_t sur (S, w_t) par l'algorithme A
 2. calculer l'erreur pondérée de h_t sur (S, w_t) : $\varepsilon_t = \sum_{i=1}^k w_t(x_i) \times \|h_t(x_i) - u_i\|$
 3. en déduire la « fiabilité » de h_t : $\alpha_t = \frac{1}{2} \ln\left(\frac{1-\varepsilon_t}{\varepsilon_t}\right)$ [$\alpha_t > 0$ si $\varepsilon_t < 0.5$, et $\rightarrow +\infty$ si $\varepsilon_t \rightarrow 0$]
 4. modifier les poids, i.e. pour i de 1 à k faire :

$$w_{t+1}(x_i) = \frac{w_t(x_i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{si } h_t(x_i) = u_i \text{ (i.e. } x_i \text{ bien classé)} \\ e^{+\alpha_t} & \text{si } h_t(x_i) \neq u_i \text{ (i.e. } x_i \text{ mal classé)} \end{cases}$$

- Fournir en sortie le classifieur global : $H(x) = \text{signe}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$

- L'écart entre l'erreur empirique la généralisation peut être borné par :

$$E_{gen}(H_T) < E_{emp}(H_T) + O\left(\sqrt{\frac{T\delta}{n}}\right)$$

(où T est le nombre d'étapes de boosting, n le nombre d'exemples, et δ la VC-dimension de l'espace des h_t)

⇒ Risque over-fitting ?? (cf aussi survalorisation des exemples « ambigus »)

- Heureusement une borne plus intéressante est :

$$E_{gen}(H_T) < \Pr(m(H_T, x) \leq \theta) + O\left(\sqrt{\frac{\delta}{n\theta^2}}\right)$$

indépendant de T !!

→ si $\Pr(m(H_T, x) < \theta)$ très faible pour un θ assez grand, alors on a une garantie de bonne généralisation...

Résumé sur le boosting

- Trois principes essentiels :

1. Combiner les avis/estimations de différents experts
 2. Modifier, avant chaque ajout d'un expert, la distribution des exemples en surpondérant au fur et à mesure les exemples mal classés
 3. Au final, utiliser une moyenne des « votes » des experts, pondérée par leurs fiabilités respectives
- Permet d'obtenir un très bon classifieur en associant des classifieurs très « faibles »
 - Problème essentiel : déterminer le « weak learner », i.e. l'algo employé à chaque étape (« round ») de boosting pour choisir/entraîner le classifieur faible h_t