

---

# GLOBOOST

## Combinaisons de Moindres Généralisés

### Produire des moindres généralisés et les combiner au sein d'un vote pondéré

**Fabien Torre**

GRAppA - Université Charles de Gaulle - Lille 3  
Domaine universitaire du « Pont de bois »  
B.P. 149  
59653 Villeneuve d'Ascq Cedex  
torre@univ-lille3.fr

---

*RÉSUMÉ.* Nous explorons dans cet article l'utilisation de moindres généralisés corrects comme apprenant dans des techniques de boosting [FRE 96]. Les premières expérimentations sur des problèmes classiques de l'UCI [BLA 98] montrent qu'ADABOOST instancié avec un apprenant à base de moindres généralisés obtient des taux d'erreur plus faibles que C4.5, GLOBO [TOR 99] et ADABOOST muni d'un apprenant plus classique. Constatant que notre nouvel apprenant peut être coûteux en temps de calcul, nous proposons un mode de génération des hypothèses qui peut être distribué sur différentes machines et une définition du poids des hypothèses a posteriori. Cela aboutit au nouvel algorithme GLOBOOST et, à nouveau, les expérimentations montrent que GLOBOOST obtient des performances comparables à celles d'ADABOOST.

*ABSTRACT.* The primary goal of this paper is to propose a new learner for boosting algorithms, namely least general generalization. First experiments conducted on benchmarks from the UCI Repository [BLA 98] show that ADABOOST boosting least general generalization obtains smaller error than reference systems. The computation time needed by these experiments leads us to define a method that could be easily distributed. This method, called GloBoost, is able to produce hypotheses independently of one another and then give a weight to each produced hypothesis. New experiments are then conducted and show low error rates for both ADABOOST and GLOBOOST. Moreover, GLOBOOST has the advantage to be naturally distributable to different computers.

*MOTS-CLÉS :* classification, boosting, ADABOOST, méthodes d'ensemble, moindres généralisés, algorithmes distribués.

*KEYWORDS:* classification, boosting, ADABOOST, ensemble methods, leveraging, least general generalization, distributed algorithms.

---

## Introduction

Le *boosting* est une technique qui a des origines théoriques ancrées dans le modèle PAC [VAL 84]. Dans ce cadre, un résultat fort a été montré : tout algorithme d'apprentissage produisant des hypothèses faisant mieux qu'une procédure de choix aléatoire peut être *boostée* jusqu'à atteindre une erreur aussi faible que voulue et ce avec une probabilité aussi grande que désirée et en temps polynomial.

Deux preuves de ce résultat ont été proposées [SCH 90, FRE 95a] ; elles ont pour point commun de définir chacune un algorithme réalisant le *boosting* d'un apprenant faible, toujours dans le cadre PAC. L'idée générale de ces algorithmes est d'utiliser plusieurs fois l'apprenant faible sur des distributions différentes, puis de combiner les différentes hypothèses obtenues au sein d'un vote.

L'algorithme utilisé dans [FRE 95a] a de plus la particularité de se rapprocher d'un cadre plus réaliste que le cadre PAC (par exemple, en constituant son échantillon au tout début de l'algorithme). Cette démarche a finalement conduit à l'algorithme ADABOOST [FRE 95b, FRE 97]. Son principe est d'associer un poids à chaque exemple, puis de solliciter l'apprenant faible sur ces exemples et ces poids pour obtenir une hypothèse qui classe correctement les exemples de poids forts ; ensuite, les poids des exemples sont modifiés en fonction de cette hypothèse (les poids des bien classés sont diminués et ceux des mal classés augmentés) et le processus itéré.

L'algorithme ADABOOST a permis de tester les idées du *boosting* sur des données réelles [FRE 96], tout comme BOOSTEXTER la version dédiée à la classification de textes [SCH 00]. Ces expérimentations, et bien d'autres ensuite, démontrent qu'ADABOOST est capable d'atteindre des erreurs en généralisation très faibles.

Dans cet article, nous nous intéressons à un constructeur d'hypothèses qui, à notre connaissance, n'a jamais été utilisé comme apprenant faible dans un algorithme de *boosting* : le calcul de *moindres généralisés corrects*. Ces moindres généralisés ont montré leur capacité à capturer des sous-concepts, en particulier sur des problèmes hautement disjonctifs [DOM 96, TOR 99]. Outre l'aspect prédictif, les moindres généralisés permettent la production de théories compréhensibles pour un expert humain.

Leur défaut majeur est le temps de calcul qu'ils nécessitent et la question se pose alors de pouvoir distribuer le calcul effectué par ADABOOST sur plusieurs machines. Or, comme nous l'avons signalé, ADABOOST produit une hypothèse en fonction des précédentes ce qui, *a priori*, empêche sa parallélisation.

Ce sont ces questions que nous traitons dans cet article dont le plan est le suivant. Nous revenons à la section 1 sur la technique du *boosting*, les questions qu'elle soulève et les éléments de réponse disponibles. À la section 2, nous rappelons ce qu'est un *moindre généralisé correct* [TOR 99] et discutons de son utilisation comme apprenant faible d'ADABOOST. La section 3 présente les résultats expérimentaux obtenus en suivant cette voie ainsi que les performances sur les mêmes données de systèmes plus classiques. Nous évoquons ensuite, section 4, la possibilité de remplacer la génération des hypothèses d'ADABOOST par une méthode de production plus facilement

parallélisable et une pondération des hypothèses qui soit détachée de cette génération. À la section 5, nous montrons à nouveau des résultats expérimentaux. Enfin, à la section 6 nous dressons un bilan de ce travail relativement aux travaux précédents et une liste des pistes ouvertes.

## 1. Quelques éléments de *boosting*

Il ne s'agit pas ici de donner une vision exhaustive des travaux sur le *boosting*, mais simplement de montrer que certaines questions relatives aux conditions et aux modes d'application de cette technique restent ouvertes.

Dans le cadre PAC [VAL 84], une classe de concepts est dite *apprenable* s'il existe un algorithme capable, pour tout concept de cette classe et toute distribution des exemples, de fournir une hypothèse dont l'erreur par rapport au concept cible peut être aussi proche de 0 que voulu et ce avec une probabilité aussi proche de 1 que désiré. Précisons qu'en plus l'algorithme doit fournir l'hypothèse en temps polynomial dans les inverses de l'erreur et de la confiance choisies.

Cette notion d'apprenabilité est dite *forte*, par opposition à la version dite *faible* qui apparaît dans [KEA 89] : une classe de concepts est dite *faiblement apprenable* s'il existe un algorithme capable, pour tout concept de cette classe et toute distribution des exemples, de fournir une hypothèse dont l'erreur est strictement inférieure à  $\frac{1}{2}$  avec une confiance non nulle. Autrement dit, il est simplement exigé que l'apprenant faible fournisse une hypothèse qui fasse mieux qu'un classifieur aléatoire.

Vient ensuite la question du lien entre ces deux notions d'apprenabilité et la réponse est surprenante : elles sont équivalentes. Une première preuve de ce résultat est donnée dans [SCH 90] (voir [KEA 94] pour une présentation limpide de cette preuve) puis une seconde dans [FRE 95a]. Chacune de ces deux preuves repose sur un algorithme auquel on fournit un apprenant faible et qui fait appel un nombre polynomial de fois à cet apprenant faible pour finalement produire une hypothèse forte. Dans les deux cas, cette hypothèse forte prend la forme d'un vote qui combine les hypothèses faibles produites par l'apprenant faible.

À la lumière de ces preuves, l'équivalence des deux notions d'apprenabilité revient à dire que, dans le cadre PAC, il est possible de *booster* un apprenant faisant à peine mieux que la hasard pour atteindre une erreur aussi petite que voulue.

Un algorithme pratique nommé ADABOOST et inspiré de la preuve de [FRE 95a] a finalement été proposé [FRE 95b, FRE 97]. ADABOOST est présenté à l'algorithme 1 dans sa version à deux classes et avec un apprenant faible pouvant fournir des hypothèses qui s'abstiennent. Cette version est reprise de [SCH 99] et définit pour chaque hypothèse trois quantités :  $W_+$  la somme des poids des exemples bien classés,  $W_-$  la somme des poids des exemples mal classés et  $W_0$  la somme des poids des exemples sur lesquels l'hypothèse s'abstient.

---

**Algorithm 1** ADABOOST

---

**Entrées :**  $n$  exemples  $x_i$  et leurs classes  $y_i \in \{-1, +1\}$ ;  $T$  un nombre d'itérations à effectuer;  $A$  un apprenant faible acceptant en entrée un échantillon d'exemples étiquetés et pondérés.

**Sortie :**  $H$  le classifieur final.

```

for  $i = 1$  to  $n$  do
   $w_i = 1/n$  {initialisation des poids des exemples}
end for
for  $t = 1$  to  $T$  do
   $h_t = A(\{(x_i, y_i, w_i)\})$ 
  for  $b \in \{-, 0, +\}$  do
     $W_b = \sum_{i: \text{sign}(y_i h_t(x_i))=b} w_i$ 
  end for
   $\alpha_t = \frac{1}{2} \log \left( \frac{W_+ + \frac{1}{2}W_0}{W_- + \frac{1}{2}W_0} \right)$ 
   $Z_t = \sum_i [w_i \cdot \exp(-\alpha_t y_i h_t(x_i))]$ 
  for  $i = 1$  to  $n$  do
     $w_i = \frac{w_i \cdot \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$ 
  end for
end for
return  $H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$ 

```

---

Dans la réalité, il est difficile de réunir les conditions théoriques énoncées dans le cadre PAC pour que le *boosting* soit effectivement garanti :

- on n'est jamais certain de disposer d'assez d'exemples et dans la majorité des cas, on ne peut pas comme dans le cadre PAC en réclamer plus à un oracle ;
- on ignore si le nombre d'étapes de *boosting* choisi est suffisant ;
- on est rarement sûr de disposer d'un apprenant faible, c'est-à-dire d'un algorithme fournissant pour toute distribution une hypothèse faisant mieux que l'aléatoire.

De plus, dans la pratique, on pourrait s'attendre à une dégradation de l'erreur en généralisation au cours du *boosting*, un trop grand nombre d'étapes ayant pour effet de sur-spécialiser le système de vote.

Malgré toutes ces réserves, l'observation de très bonnes performances obtenues par ADABOOST sur des données réelles est souvent rapportée. Les conditions imposées par le cadre PAC et les théorèmes de *boosting* n'étant clairement pas remplies, il faut trouver une autre explication à ces bons résultats.

Dans [BRE 96a], l'erreur d'un classifieur est décomposée en biais et variance puis il est avancé que, comme toute méthode de type *arcing* (*adaptively resample and combine*) et comme le *bagging* [BRE 96b], ADABOOST diminuerait la composante due à la variance. Cet argument est battu en brèche dans [FRE 98] où sont présentées des expérimentations où l'on voit ADABOOST faire chuter à la fois les deux composantes de l'erreur mais toujours sans explication du phénomène.

Une observation sur les algorithmes de *boosting* est souvent rapportée : l'erreur en généralisation décroît encore après que l'erreur en apprentissage soit stable ou même nulle. L'explication est que même si tous les exemples d'apprentissage sont déjà bien classés, la poursuite du *boosting* tend à maximiser les *marges* [SCH 97]. À la suite de cette explication, certains ont cherché à maximiser explicitement la marge [RäT 02, RäT 03, GRO 98] et ont dû constater que leurs systèmes étaient moins performants qu'ADABOOST. Sur cette thématique, citons encore [HAR 99] qui parvient à construire un classifieur donnant une marge de 1 à chacun des exemples d'apprentissage mais qui a des performances toujours moins bonnes qu'ADABOOST. Cette dernière expérience permet de rejeter définitivement l'explication des performances d'ADABOOST par les marges.

Nous pouvons également nous interroger sur la signification des poids  $\alpha_t$  qu'ADABOOST associe aux hypothèses produites. À la vue de l'algorithme, il s'agit d'une valeur déterminée sur une distribution artificielle, elle-même fonction des échecs et réussites des hypothèses précédentes. Cependant, nous avons noté lors d'expériences sur des données très simples que l'erreur en généralisation diminuait encore alors que l'apprenant faible avait déjà fourni toutes les hypothèses possibles. Autrement dit, lorsqu'une hypothèse apparaît plusieurs fois, elle vote finalement avec un poids, cumul de tous ses  $\alpha_t$ , qui a peut-être un caractère plus absolu. Si c'est bien le cas, on peut espérer approcher ces valeurs par un processus non adaptatif.

Enfin, une dernière question est celle de l'apprenant faible : comment bien le choisir ? À nouveau les éléments de réponse sont limités. Les auteurs d'ADABOOST indiquent que la définition simple et claire d'un apprenant faible est perdue avec le passage du cadre PAC aux problèmes réels et que la seule caractérisation d'un bon apprenant dont nous disposons est : *an algorithm that gives small errors when run under Adaboost* [FRE 98]. En particulier, faut-il autoriser ou même encourager l'apprenant faible à se tromper sur une partie des exemples d'apprentissage (les exemples de poids faibles, bien classés par les hypothèses produites précédemment) ? La réponse n'est pas évidente, surtout si l'on prend en compte qu'ADABOOST traite de manière privilégiée des données non bruitées (son comportement naturel sur des données bruitées est de s'acharner sur les exemples difficiles à classer, donc sur le bruit).

Dans cet article, nous proposons justement d'utiliser un apprenant faible qui ne produit que des hypothèses correctes : celles-ci peuvent donc s'abstenir sur une partie des exemples mais en aucun cas se tromper sur un exemple. Il s'agit de *moindres généralisés maximalement corrects*.

## 2. Booster des moindres généralisés avec ADABOOST

### 2.1. Moindres généralisés corrects

Comme son nom l'indique, le moindre généralisé de deux exemples est l'hypothèse la plus spécifique possible qui couvre ces deux exemples. En attribut-valeur, cette hypothèse est unique et peut être simplement définie comme suit :

- si les deux exemples partagent la même valeur d'un attribut catégoriel, celle-ci est conservée dans le moindre généralisé ; si ces valeurs sont distinctes, le moindre généralisé fait apparaître une valeur indéterminée pour cet attribut ;
- pour un attribut continu, on construit le plus petit intervalle incluant les deux valeurs présentes dans les exemples à généraliser.

Cette procédure apparaît clairement sur l'exemple suivant :

	âge	fumeur	sexe	classe
e <sub>1</sub>	25	non	homme	positif
e <sub>2</sub>	35	non	femme	positif
g <sub>1</sub>	[25, 35]	non	?	positif

Naturellement, cette opération peut se généraliser pour s'appliquer à deux hypothèses ou, cas qui nous intéresse le plus, à une hypothèse et un exemple :

	âge	fumeur	sexe	classe
g <sub>1</sub>	[25, 35]	non	?	positif
e <sub>3</sub>	30	oui	homme	positif
g <sub>2</sub>	[25, 35]	?	?	positif
e <sub>4</sub>	40	non	femme	positif
g <sub>3</sub>	[25, 40]	?	?	positif

Du point de vue de l'apprentissage supervisé, il est intéressant de généraliser ensemble des exemples d'une même classe, en prenant garde à ne pas couvrir d'exemples d'autres classes. Ce calcul, qui fournit la caractérisation d'un sous-concept au sein d'une classe, est celui d'un *moindre généralisé maximale correct* et est décrit formellement par l'algorithme 2 : étant donné un exemple-graine, on essaye de généraliser, un à un, les exemples de la même classe avec comme critère d'acceptation la non couverture d'exemples d'autres classes (*correction*). La généralisation produite est *maximale correcte* dans le sens où plus aucun exemple d'apprentissage ne peut lui être ajoutée sans perdre la correction. Notons également que cette hypothèse est fonction de l'ordre dans lequel les exemples sont testés pour l'ajout.

---

**Algorithm 2** MG\_Correct\_v1
 

---

**Entrées :**  $s$  un exemple graine,  $P = \{p_1, \dots, p_n\}$  un ensemble *ordonné* de  $n$  exemples de la même classe que la graine,  $N$  un ensemble de contre exemples. L'algorithme utilise l'opération élémentaire de **MoindreGénéralisé** entre une hypothèse et un exemple, ainsi qu'un test de subsomption noté  $\succeq$  permettant de décider si une hypothèse couvre ou non un exemple (ces deux opérations sont à définir en fonction du langage de description choisi).

**Sortie :**  $g$  une généralisation de  $P$ , maximale correcte par rapport à  $N$ .

```

 $g = s$ 
for  $i = 1$  to  $n$  do
     $g' = \text{MoindreGénéralisé}(g, p_i)$  {Généralisation entre deux hypothèses.}
    if  $(\forall e \in N : g' \not\succeq e)$  then
         $g = g'$  {Si  $g'$  est correcte, elle devient la généralisation courante.}
    end if
end for
return  $g$ 
    
```

---

L'exemple suivant montre trois exécutions possibles de l'algorithme 2 (les exemples qui provoquent une perte de correction pour la généralisation courante sont barrés) :

graine	autres positifs	→	paquet maximale correct	règle
$p_1$	$p_5$ $p_8$ <del><math>p_3</math></del> $p_{14}$ ...	→	$\{p_1, p_5, p_8, p_{14}, \dots\}$	$g_1$
$p_2$	<del><math>p_4</math></del> $p_3$ <del><math>p_1</math></del> $p_{12}$ ...	→	$\{p_2, p_3, p_{12}, \dots\}$	$g_2$
$p_3$	$p_7$ <del><math>p_4</math></del> <del><math>p_5</math></del> <del><math>p_8</math></del> <del><math>p_{12}</math></del> ...	→	$\{p_3, p_7, \dots\}$	$g_3$

On y voit l'importance de l'ordre de présentation des exemples et également que la contrainte de correction amène des généralisations qui prisent isolément ne couvrent pas tous les exemples positifs. Autrement dit, une telle généralisation ne peut conclure que sur une seule classe, la classe des exemples sur lesquels elle a été construite. Pour tous les autres exemples, ceux de la même classe qui ne sont pas couverts et ceux d'autres classes, la généralisation ne fournit pas de classe, elle s'abstient.

Cette brique de base a déjà été utilisée dans l'algorithme GLOBO [TOR 99] (cf. algorithme 3). Dans ce cas, chaque exemple joue successivement le rôle de graine et pour chacun, la liste des exemples de même classe est mélangée aléatoirement. Puis GLOBO opère une couverture minimale pour ne garder que le nombre minimum d'hypothèses permettant de couvrir tous les exemples. Le but de la couverture minimale dans GLOBO est de produire finalement un ensemble réduit de règles qui soit appréhendable par un expert humain. Cette approche s'est révélée pertinente puisque GLOBO a remporté le PTE Challenge [SRI 97, SRI 99] dans la catégorie des systèmes produisant une théorie compréhensible.

**Algorithm 3** GLOBO

---

**Entrées :**  $n$  exemples  $x_i$  et leurs classes  $y_i \in \{-1, +1\}$ .  
**Sortie :**  $H$  un ensemble de règles.

$H' = \emptyset$

**for**  $i = 1$  to  $n$  **do**

$P = \{x_j | y_j = y_i \wedge i \neq j\}$   
 $N = \{x_j | y_j \neq y_i\}$

**mélanger**  $P$  **aléatoirement**

$h_i = \text{MG\_Correct\_v1}(x_i, P, N)$  {Appel à l'algorithme 2}  
ajouter  $h_i$  à  $H'$

**end for**

$H = \emptyset$

**while**  $(\exists x_i, \forall h_j \in H, h_j \not\prec x_i)$  **do**

$h = \text{ArgMax}_{h_i \in H'} |\{x_j : h_i \succeq x_j \wedge \nexists h_k \in H : h_k \not\prec x_j\}|$   
ajouter  $h$  à  $H$

**end while**

return  $H$

---

Dans cet article, notre volonté est d'abandonner la compréhensibilité pour gagner en capacité prédictive et cela en utilisant le calcul de moindre généralisé correct comme apprenant faible dans des techniques de *boosting*.

## 2.2. Instanciation de l'algorithme ADABOOST

Nous proposons dans cette section une nouvelle version du calcul de moindres généralisés corrects destinée à devenir un apprenant faible pour ADABOOST. Cela implique de prendre en compte les poids des exemples, autrement dit de favoriser la couverture des exemples de poids élevé. Avec cet objectif, l'idée de notre apprenant faible est de réutiliser l'algorithme 2 en triant au préalable les exemples candidats à la généralisation par poids décroissants. Précisons immédiatement que, quels que soient les poids, nous ne relâchons pas notre critère de correction : la généralisation produite ne couvrira aucun contre-exemple, même de poids faible. Cette nouvelle version est présentée à l'algorithme 4.

Cet apprenant peut être fourni en l'état à ADABOOST (algorithme 1) dont l'utilisation ne nécessite pas d'autre information. Remarquons simplement que la formule de  $\alpha_t$  (poids des hypothèses) utilisée dans ADABOOST se simplifie en prenant en compte que les moindres généralisés ainsi calculés ne font pas d'erreur de classification (ils classifient bien ou s'abstiennent). En effet,  $W_- = 0$  et par suite  $W_0 = 1 - W_+$  et finalement

$$\alpha_t = \frac{1}{2} \log \left( \frac{1 + W_+}{1 - W_+} \right)$$



---

**Algorithm 4** MG\_Correct\_v2 (apprenant faible)

---

**Entrées :**  $n$  exemples  $(x_i, y_i, w_i)$  avec leurs étiquettes et poids.**Sortie :**  $g$  une généralisation maximale correcte d'exemples de poids élevés. $target$  = classe choisie au hasard $seed$  = l'exemple de la classe  $target$  ayant le poids le plus fort $P = \{x_i | y_i = target, x_i \neq seed\}$  $N = \{x_i | y_i \neq target\}$ **trier  $P$  par poids décroissants**return MG\_Correct\_v1( $seed, P, N$ ) {Appel à l'algorithme 2}

---

La section suivante présente les résultats de ADABOOST ainsi instancié, algorithme que nous appellerons ADABOOSTMG dans la suite (ou en plus concis ABMG).

### 3. Expérimentations avec ADABOOST

#### 3.1. Protocole expérimental

Notre version d'ADABOOST est maintenant comparée aux algorithmes suivants : le classique C4.5 ; GLOBO, système à base de moindres généralisés mais sans *boosting* présenté à l'algorithme 3 ; et ADABOOST (algorithme 1) utilisant comme apprenant faible des *decision stumps*, arbres de décision limités à un seul niveau, utilisés comme hypothèses faibles par exemple dans [FRE 96]. Ce dernier système a été simulé à l'aide du système ADTREE [De 01] en mode BOOSTEXTER et sera désigné dans la suite par ADABOOSTDS (ou en plus concis ABDS).

Ces systèmes sont évalués sur vingt problèmes classiques de l'UCI [BLA 98], représentant un éventail assez large des problèmes que l'on peut rencontrer : attributs tous continus, tous discrets, ou mélange des deux types, avec deux classes à prédire ou plus, avec ou sans valeurs manquantes, hautement disjonctif ou pas.

Pour chacun de ces problèmes, on utilise une validation croisée 10 fois, le découpage étant le même pour toutes les expérimentations<sup>1</sup>. Finalement, pour chaque problème et pour chaque système, on mesure le taux d'erreur moyen sur les 10 blocs de test. Précisons que les systèmes stochastiques (GLOBO et ADABOOSTMG) sont exécutés dix fois sur chaque bloc et c'est la moyenne sur ces dix exécutions qui est utilisée pour mesurer l'erreur du système.

---

1. Les fichiers de la validation croisée ainsi que des résultats complémentaires sont disponibles à l'adresse : <http://www.grappa.univ-lille3.fr/~torre/Recherche/Experiments/>.

<b>Problèmes</b>	<b>C4.5</b>	<b>GLOBO</b>	<b>ADABOOSTDS</b>	<b>ADABOOSTMG</b>
audiology	18.20	20.76	<b>13.43</b>	25.78
breast-cancer	4.87	3.89	4.16	<b>3.02</b>
car	<b>7.69</b>	10.17	13.42	17.92
cmc	48.07	50.60	<b>44.80</b>	56.32
crx	<b>14.79</b>	16.50	15.80	17.82
dermatology	6.23	8.51	<b>3.51</b>	6.58
ecoli	15.89	23.88	<b>15.59</b>	23.24
glass	28.72	5.57	23.62	<b>4.70</b>
hepatitis	20.70	20.09	<b>17.64</b>	17.88
horse-colic	<b>13.63</b>	22.29	18.44	20.90
house-votes-84	<b>3.22</b>	7.39	4.62	56.15
ionosphere	7.96	8.74	<b>6.78</b>	8.51
iris	<b>5.33</b>	7.47	<b>5.33</b>	5.93
pima	29.29	27.04	<b>25.13</b>	29.68
promoters	18.17	22.60	<b>6.83</b>	12.93
sonar	28.97	31.09	<b>16.86</b>	25.25
tic-tac-toe	14.40	1.11	13.78	<b>0.23</b>
vowel	21.21	23.99	<b>20.71</b>	21.30
wine	8.83	8.94	14.18	<b>5.41</b>
zoo	7.51	5.55	4.25	<b>3.67</b>
<i>Moyennes</i>	<i>16.18</i>	<i>16.31</i>	<b>14.45</b>	<i>18.16</i>

(a) Taux d'erreur.

	<b>C4.5</b>	<b>GLOBO</b>	<b>ABDS</b>	<b>ABMG</b>	<b>Bilan</b>
<b>C4.5</b>	-	14-6	5-14	12-8	31-28
<b>GLOBO</b>	6-14	-	5-15	6-14	17-43
<b>ABDS</b>	14-5	15-5	-	15-5	<b>44-15</b>
<b>ABMG</b>	8-12	14-6	5-15	-	27-33

(b) Confrontations.

**Tableau 1.** Résultats de C4.5, GLOBO, ADABOOSTDS et ADABOOSTMG sur vingt problèmes de l'UCI. Le nombre d'étapes de boosting vaut 100 pour ADABOOSTDS et ADABOOSTMG.

### 3.2. Résultats pour 100 étapes de boosting

Les taux d'erreur observés pour 100 étapes de *boosting* sont présentés dans le tableau 1(a) où l'on a en plus fait apparaître pour chaque système la moyenne des taux d'erreur sur l'ensemble des problèmes. Pour chacun des problèmes, l'erreur la plus faible est distinguée en gras.

Il y apparaît d'emblée que les systèmes à base d'arbres de décision (C4.5 et ADA-BOOSTDS) obtiennent les meilleurs résultats.

Le tableau 1(b) explicite cette impression en donnant les résultats des confrontations des méthodes prises deux à deux. Chaque ligne de cette table contient les résultats d'une méthode particulière : chaque case fournit le nombre de problèmes sur lesquels cette méthode a obtenu une erreur plus faible que la méthode en colonne, ainsi que le nombre de fois où la méthode considérée s'est montrée moins performante. Par souci de clarté, le nombre de nuls n'est pas rapporté mais peut être retrouvé aisément. Ces données apparaissent en italique si l'algorithme en ligne est meilleur que celui en colonne. Enfin, la dernière colonne totalise le nombre de victoires et de défaites obtenues par la méthode considérée.

À l'issue de ces expérimentations utilisant 100 tours de boosting, les méthodes à base de moindres généralisés sont battues. Notons cependant que ADABOOSTMG surpasse GLOBO sur 14 des 20 problèmes étudiés.

### 3.3. Résultats pour 1000 étapes de boosting

Le tableau 2 présente les résultats d'expérimentations identiques où le nombre d'étapes de *boosting* est passé de 100 à 1000 ; il amène les constatations suivantes :

- l'erreur de ADABOOSTMG diminue notablement sur quasiment tous les problèmes lorsque l'on passe de 100 étapes de *boosting* à 1000 (l'erreur globale chute de 18.16 % à 12.70 %) ;
- ADABOOSTMG surpasse ou égale C4.5 sur 14 problèmes sur 20 et cette supériorité est aussi apparente sur les taux d'erreur moyens (12.70 % contre 16.18 %) ;
- GLOBO est battu sur tous les problèmes par ADABOOSTMG et l'on mesure ainsi ce que nous avons gagné en abandonnant la compréhensibilité des théories produites par GLOBO (en moyenne, on passe de 16.31 % d'erreurs à 12.70 %) ;
- ADABOOSTDS supporte moins bien que ADABOOSTMG le passage à 1000 étapes de *boosting* et tous les indicateurs de performance se détériorent par rapport aux résultats précédents ; que ce soit avec 100 ou 1000 étapes de *boosting*, ADABOOSTDS est maintenant battu par ADABOOSTMG utilisant 1000 étapes de *boosting*.

Un dernier moyen de comparer deux méthodes d'apprentissage vis-à-vis d'un ensemble de problèmes consiste à construire un graphique où chaque problème est représenté par un point dont les coordonnées correspondent aux erreurs respectives obtenues par les deux compétiteurs. En plus de ces points, nous pouvons faire apparaître

Problèmes	ADABOOSTDS		ADABOOSTMG	
	100	1000	100	1000
audiology	<b>13.43</b>	15.79	25.78	19.22
breast-cancer	4.16	4.16	<b>3.02</b>	3.10
car	13.42	13.36	17.92	<b>9.74</b>
cmc	<b>44.80</b>	44.87	56.32	49.61
crx	15.80	16.53	17.82	<b>14.47</b>
dermatology	<b>3.51</b>	4.38	6.58	4.63
ecoli	<b>15.59</b>	<b>15.59</b>	23.24	19.25
glass	23.62	22.53	4.70	<b>4.59</b>
hepatitis	<b>17.64</b>	18.23	17.88	18.39
horse-colic	18.44	21.71	20.90	<b>18.42</b>
house-votes-84	<b>4.62</b>	5.31	56.15	6.12
ionosphere	<b>6.78</b>	12.19	8.51	7.44
iris	<b>5.33</b>	<b>5.33</b>	5.93	<b>5.33</b>
pima	25.13	25.39	29.68	<b>25.13</b>
promoters	<b>6.83</b>	18.67	12.93	7.00
sonar	16.86	<b>15.03</b>	25.25	23.34
tic-tac-toe	13.78	1.67	0.23	<b>0.00</b>
vowel	20.71	14.65	21.30	<b>9.30</b>
wine	14.18	18.07	<b>5.41</b>	5.73
zoo	4.25	5.25	3.67	<b>3.21</b>
<i>Moyennes</i>	<i>14.45</i>	<i>14.94</i>	<i>18.16</i>	<i>12.70</i>

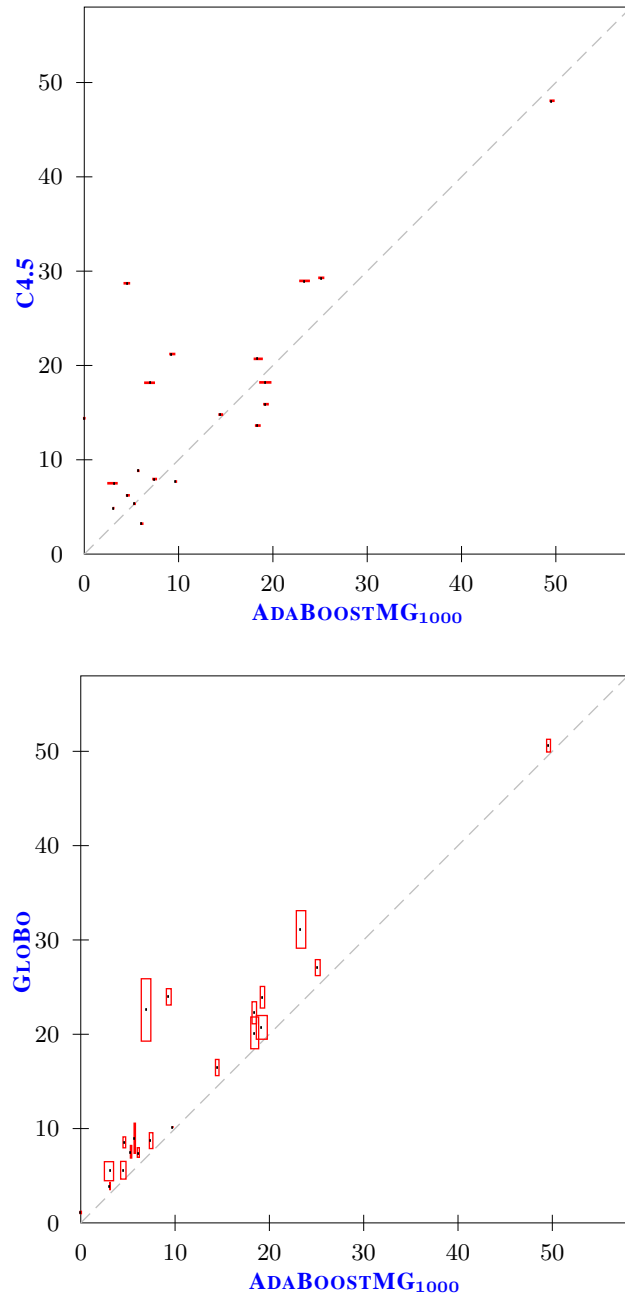
(a) Taux d'erreur.

	ABDS <sub>100</sub>	ABDS <sub>1000</sub>	ABMG <sub>100</sub>	ABMG <sub>1000</sub>	<i>Bilan</i>
ABDS <sub>100</sub>	-	12-5	15-5	9-10	36-20
ABDS <sub>1000</sub>	5-12	-	11-9	7-12	23-33
ABMG <sub>100</sub>	5-15	9-11	-	3-17	17-43
ABMG <sub>1000</sub>	10-9	12-7	17-3	-	<b>39-19</b>

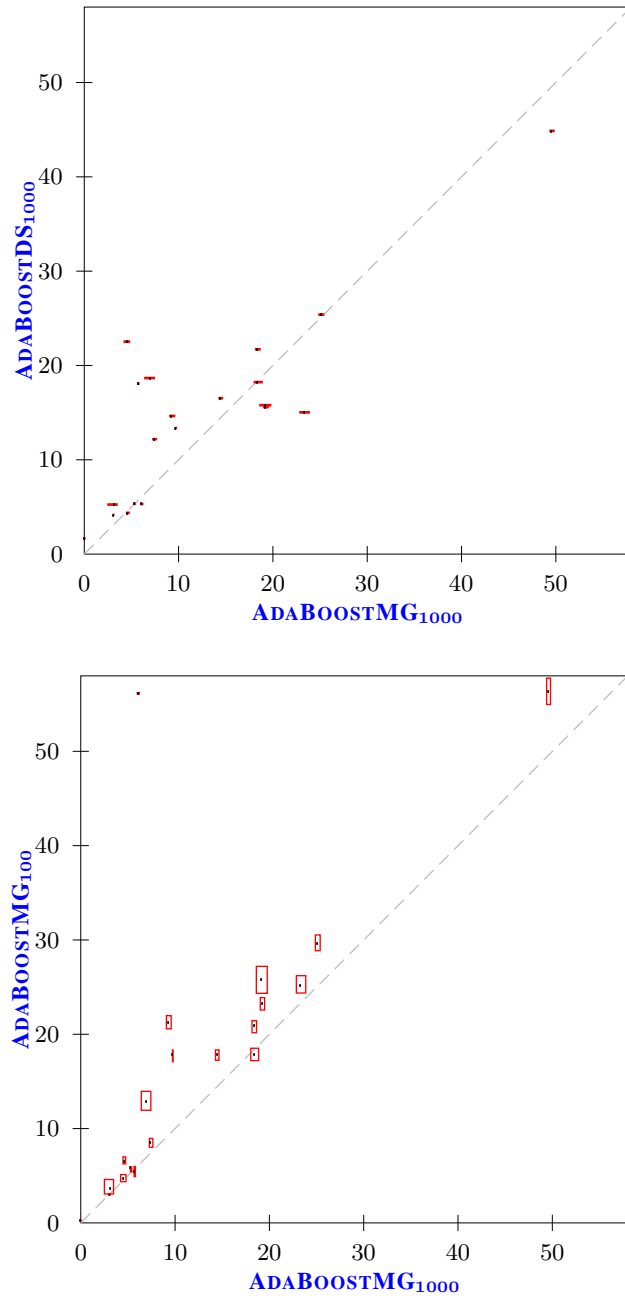
	C4.5	GLOBO	ABDS <sub>100</sub>	ABMG <sub>1000</sub>	<i>Bilan</i>
C4.5	-	14-6	5-14	6-13	25-33
GLOBO	6-14	-	5-15	0-20	11-49
ABDS <sub>100</sub>	14-5	15-5	-	9-10	38-20
ABMG <sub>1000</sub>	13-6	20-0	10-9	-	<b>43-15</b>

(c) Confrontations.

**Tableau 2.** ADABOOSTMG et ADABOOSTDS pour 100 et 1000 étapes de boosting.



**Figure 1.** ADABOOSTMG opposé à deux méthodes non boostées : C4.5 et GLOBO.



**Figure 2.**  $ADABOOSTMG_{1000}$  contre  $ADABOOSTDS_{1000}$  et  $ADABOOSTMG_{100}$ .

les écarts types observés pour les méthodes stochastiques. Dans la suite, la méthode testée est systématiquement placée en abscisse et sa supériorité par rapport à la seconde méthode (en ordonnée) est mesurée par le nombre de points au dessus de la droite  $y = x$ , ainsi que par les distances des points à cette droite.

Par exemple, les figures 1(a) et 1(b) montrent de tels graphiques pour ADABOOSTMG opposé à C4.5 puis à GLOBO. La première montre que la supériorité de ADABOOSTMG sur C4.5 vaut même en prenant en compte la variabilité du *boosting* de moindres généralisés. La seconde confirme que, sur tous nos problèmes, ADABOOSTMG surpasse GLOBO avec l'information supplémentaire que ADABOOSTMG présente des écarts types bien plus réduits que GLOBO.

La figure 2(a) illustre la confrontation ADABOOSTMG vs. ADABOOSTDS et, enfin, la figure 2(b) le passage de 100 à 1000 étapes de *boosting* : outre la réduction des taux d'erreur, la réduction des écarts types est également observée.

Ces résultats nous encouragent à augmenter encore le nombre d'étapes de *boosting* pour continuer à faire baisser l'erreur. Malheureusement, le calcul de moindres généralisés corrects est coûteux : pour chacun, il faut essayer d'ajouter les exemples de la classe cible un à un et chacun de ces ajouts doit être validé sur l'ensemble des contre-exemples.

Ainsi, le plus long des apprentissages avec 1000 étapes de *boosting* et l'utilisation de moindres généralisés prend une heure et vingt minutes sur une machine i686 à 2.5 Ghz pour une base qui contient moins de 1000 exemples. Ce temps n'est pas prohibitif en soi et est sans doute comparable au temps requis par bon nombre d'apprenants faibles. Cependant, notre protocole implique de réaliser 100 fois cet apprentissage (10 blocs de validation croisée et 10 exécutions de chaque bloc) ce qui cette fois est difficile à mener dans un temps raisonnable.

Il est donc impératif, pour que la méthode soit viable, de rendre l'apprentissage plus rapide. Nous avons choisi pour cela de paralléliser la génération des hypothèses et donc de différer l'affectation de poids à ces hypothèses. Cette étude est l'objet de la section suivante.

#### 4. Produire et combiner des moindres généralisés sans ADABOOST

Nous l'avons vu, la production d'un moindre généralisé correct est de complexité quadratique dans le nombre d'exemples disponibles et cela n'est pas acceptable pour un apprenant faible destiné à être appelé de nombreuses fois dans un algorithme de type ADABOOST.

Nous aurions pu, pour réduire le temps de calcul, échantillonner et distribuer les données à différentes unités de calcul. Nous avons préféré distribuer, non pas les données, mais le calcul des hypothèses, lequel se fera toujours sur l'ensemble complet des données disponibles. Cela nécessite de savoir :

- produire les hypothèses indépendamment les unes des autres ;
- affecter un poids à chacune de ces hypothèses *a posteriori*.

Nous commençons par le deuxième point en attribuant de nouveaux poids aux hypothèses disponibles, quel que soit leur mode de production.

#### 4.1. Affectation de poids aux hypothèses produites

Le premier de ces poids, dans la suite repéré par le nom *adalike*, cherche à imiter le mode de calcul des poids d'ADABOOST. Comme ADABOOST, nous commençons par définir le poids d'un exemple  $x_i$ , cela à partir de  $T_i$  le nombre de fois où cet exemple est couvert par les hypothèses produites :

$$w_i = \exp(-T_i)$$

Cette formule fait référence à ADABOOST (algorithme 1) qui multiplie le poids d'un exemple par  $\exp(-\alpha_t)$  pour chaque hypothèse  $h_t$  couvrant l'exemple. Une fois ces poids attribués aux exemples, on les normalise en divisant chacun par  $Z = \sum_i w_i$ . Puis nous définissons le poids d'une hypothèse  $h$ .  $W_+$  est défini comme la somme des poids des exemples couverts par  $h$  et  $n$  est le nombre de fois où cette hypothèse est apparue pendant la production. Finalement, le poids *adalike* de  $h$  est calculé comme suit :

$$\alpha(h) = \frac{n}{2} \log \left( \frac{1 + W_+}{1 - W_+} \right)$$

La valeur de  $n$  intervient pour simuler le cumul des poids effectué par ADABOOST lorsqu'une même hypothèse apparaît plusieurs fois.

Nous considérons également les poids suivants, plus simples :

- *covering* : chaque hypothèse distincte vote avec un poids qui correspond au nombre d'exemples qu'elle couvre dans l'ensemble d'apprentissage (dans notre cas, les hypothèses sont correctes par construction et les exemples couverts par une hypothèse sont donc tous de la même classe) ;
- *frequency* : chaque hypothèse distincte vote avec un poids qui correspond à son nombre d'apparitions pendant la phase de production ;
- *uniform* : toutes les hypothèses distinctes votent avec un poids égal.

#### 4.2. Production d'hypothèses sans ADABOOST

L'étape suivante est de trouver une alternative à ADABOOST pour produire des hypothèses. Nous proposons pour cela GLOBOOST (algorithme 5) : celui-ci produit des moindres généralisés corrects de façon purement stochastique et, par conséquent, sans aucun lien de dépendance entre eux. Cela nous permettra à terme de partager ce travail de génération entre plusieurs machines.



---

**Algorithm 5** GLOBOOST
 

---

**Entrées :**  $n$  exemples  $(x_i, y_i)$  avec étiquettes ;  $T$  un nombre d'itérations.

**Sortie :**  $H$  le classifieur final.

```

for  $t = 1$  to  $T$  do
     $target$  = classe choisie au hasard
     $seed$  = un exemple de la classe  $target$  choisi au hasard
     $P = \{x_i | y_i = target, x_i \neq seed\}$ 
     $N = \{x_i | y_i \neq target\}$ 
    mélanger  $P$  aléatoirement
     $h_t = \text{MG\_Correct\_v1}(seed, P, N)$  {Appel à l'algorithme 2, page 7}
end for
for all  $t$  tel que  $1 \leq t \leq T$  do
     $\alpha_t = \text{poids}(h_t)$  {fonction à instancier}
end for
return  $H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t \cdot h_t(x) \right)$ 
    
```

---

## 5. Évaluation expérimentale de GLOBOOST

Dans cette section, nous évaluons expérimentalement les idées de la section précédente : génération des hypothèses par GLOBOOST en remplacement d'ADABOOST et affectation de poids *a posteriori*. Le protocole et les données sont les mêmes qu'à la section 3.1.

### 5.1. Évaluation expérimentale des poids proposés

Dans un premier temps, nous conservons l'algorithme ADABOOST comme générateur d'hypothèses. Après la phase de *boosting*, nous obtenons un ensemble d'hypothèses avec leurs poids attribués par ADABOOST et nous leur affectons les poids supplémentaires définis précédemment. On obtient ainsi des systèmes de vote distincts même si tous sont fondés sur les mêmes votants.

Notre intuition, forgée par les différents travaux qui cherchent à redéfinir les poids pour maximiser la marge [SCH 97, Rät 02, Rät 03, GRO 98, HAR 99], est que ADABOOST a un comportement quasi optimal dans le réglage des poids des hypothèses. Il s'agit donc ici de quantifier ce que nous perdons en abandonnant les poids des hypothèses fixés par ADABOOST.

Les résultats pour 100 (respectivement 1000) étapes de *boosting* sont présentés au tableau 3 (respectivement tableau 4). La première constatation est que ADABOOST reste le meilleur pour fixer les poids des hypothèses qu'il produit, en particulier lorsque le nombre d'étapes de *boosting* augmente.

*A contrario*, notre imitation *adalike* a un meilleur comportement pour un nombre d'étapes de *boosting* de 100 : ce poids obtient les taux d'erreur les plus faibles sur 10

<b>Problèmes</b>	<b>adaboost</b>	<b>adalike</b>	<b>covering</b>	<b>frequency</b>	<b>uniform</b>
audiology	25.78	<b>24.14</b>	28.32	25.97	26.25
breast-cancer	3.02	4.17	<b>2.96</b>	4.63	4.65
car	17.92	<b>17.89</b>	18.14	18.39	18.39
cmc	56.32	56.31	<b>56.27</b>	56.62	56.62
crx	17.82	<b>17.12</b>	17.77	20.45	20.43
dermatology	6.58	<b>6.55</b>	7.20	7.18	7.32
ecoli	<b>23.24</b>	24.21	25.25	25.37	25.32
glass	<b>4.70</b>	<b>4.70</b>	4.83	<b>4.70</b>	4.83
hepatitis	17.88	<b>15.65</b>	18.98	20.67	20.72
horse-colic	<b>20.90</b>	21.06	21.24	23.16	23.40
house-votes-84	56.15	37.59	14.94	56.38	<b>6.20</b>
ionosphere	8.51	27.95	9.61	<b>8.14</b>	<b>8.14</b>
iris	5.93	<b>5.07</b>	6.00	6.27	6.27
pima	29.68	<b>29.61</b>	30.23	32.46	32.46
promoters	<b>12.93</b>	16.12	16.45	15.50	15.33
sonar	25.25	40.70	25.28	<b>24.26</b>	<b>24.26</b>
tic-tac-toe	<b>0.23</b>	<b>0.23</b>	<b>0.23</b>	0.71	1.13
vowel	21.30	22.60	21.38	<b>20.14</b>	<b>20.14</b>
wine	5.41	16.40	<b>5.35</b>	5.46	<b>5.35</b>
zoo	3.67	<b>3.52</b>	5.07	4.20	5.40
<i>Moyennes</i>	<i>18.16</i>	<i>19.58</i>	<i>16.78</i>	<i>19.03</i>	<b>16.63</b>

(a) Taux d'erreur.

	<b>adaboost</b>	<b>adalike</b>	<b>covering</b>	<b>frequency</b>	<b>uniform</b>	<b>Bilan</b>
<b>adaboost</b>	-	8-10	14-5	16-3	15-5	<b>53-23</b>
<b>adalike</b>	10-8	-	12-7	14-5	14-6	50-26
<b>covering</b>	5-14	7-12	-	12-8	12-6	36-40
<b>frequency</b>	3-16	5-14	8-12	-	8-5	24-47
<b>uniform</b>	5-15	6-14	6-12	5-8	-	22-49

(b) Confrontations.

**Tableau 3.** Résultats selon les poids attribués aux moindres généralisés corrects produits par ADABOOSTMG pour 100 étapes de boosting.

Problèmes	adaboost	adalike	covering	frequency	uniform
audiology	19.22	<b>19.17</b>	26.89	20.66	21.72
breast-cancer	3.10	10.36	<b>3.00</b>	3.96	3.94
car	9.74	<b>9.56</b>	9.74	15.83	10.20
cmc	49.61	49.24	<b>48.98</b>	51.25	51.03
crx	<b>14.47</b>	26.50	14.60	21.22	20.71
dermatology	<b>4.63</b>	12.93	5.92	5.51	7.15
ecoli	<b>19.25</b>	23.92	21.02	22.46	21.62
glass	<b>4.59</b>	8.65	6.33	<b>4.59</b>	6.33
hepatitis	<b>18.39</b>	22.74	19.28	21.21	23.06
horse-colic	18.42	<b>18.17</b>	19.54	19.80	18.86
house-votes-84	6.12	34.84	<b>5.07</b>	30.06	5.80
ionosphere	<b>7.44</b>	57.79	9.23	7.64	7.56
iris	<b>5.33</b>	6.67	6.00	6.00	6.80
pima	25.13	<b>24.75</b>	25.55	29.37	29.39
promoters	<b>7.00</b>	29.38	12.83	8.02	13.63
sonar	23.34	49.52	22.53	<b>21.83</b>	<b>21.83</b>
tic-tac-toe	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	0.72	2.71
vowel	9.30	34.43	9.85	9.17	<b>9.16</b>
wine	5.73	45.90	<b>4.62</b>	5.23	4.67
zoo	<b>3.21</b>	3.55	6.78	4.59	5.58
<i>Moyennes</i>	<b>12.70</b>	<i>24.40</i>	<i>13.89</i>	<i>15.46</i>	<i>14.59</i>

(a) Taux d'erreur.

	adaboost	adalike	covering	frequency	uniform	<i>Bilan</i>
<b>adaboost</b>	-	<i>14-5</i>	<i>13-6</i>	<i>16-3</i>	<i>16-4</i>	<b>59-18</b>
<b>adalike</b>	5-14	-	5-14	7-13	9-11	26-52
<b>covering</b>	6-13	<i>14-5</i>	-	<i>11-8</i>	<i>13-6</i>	44-32
<b>frequency</b>	3-16	<i>13-7</i>	8-11	-	9-10	33-44
<b>uniform</b>	4-16	<i>11-9</i>	6-13	<i>10-9</i>	-	31-47

(b) Confrontations.

**Tableau 4.** Résultats selon les poids attribués aux moindres généralisés corrects produits par ADABOOSTMG pour 1000 étapes de boosting.

des vingt problèmes considérés et est le meilleur de tous les poids dans les confrontations deux à deux. Malheureusement les performances de *adalike* se dégradent avec l'augmentation du nombre d'étapes de *boosting* au point de devenir le moins bon de tous les poids étudiés.

Les poids *covering*, *frequency* et même *uniform* présentent également des résultats intéressants et qui en plus résistent au passage à 1000 étapes de *boosting*; *covering* en particulier semble très proche du poids d'ADABOOST pour 1000 étapes de *boosting*. Il est important de noter que, même si ces poids sont finalement moins bons que le poids *adaboost*, ils maintiennent de meilleures performances que celles des systèmes C4.5 et GLOBO évalués à la section 3, tableau 1(a).

En conclusion de ces expérimentations, nous disposons de poids qui peuvent raisonnablement se substituer aux poids traditionnels d'ADABOOST.

## 5.2. Évaluation expérimentale de GLOBOOST

Toujours avec le protocole et les données de la section 3.1, nous évaluons l'algorithme GLOBOOST muni des différents poids dont nous disposons (*adalike*, *covering*, *frequency* et *uniform*). Les résultats sont présentés aux tableaux 5 et 6.

À nouveau, notre poids *adalike* fait apparaître des valeurs aberrantes et supporte mal l'augmentation du nombre d'étapes de *boosting*. Par contre, les poids *covering*, *frequency* et *uniform* confirment leurs bons résultats au point de concurrencer ADABOOSTMG. De plus, cette tendance se confirme lorsque l'on augmente le nombre d'étapes de 100 à 1000. Même le poids *uniform* obtient des résultats très honorables en regard de sa simplicité. Ce résultat est surprenant pour nous : ces poids faisaient baisser les performances d'ADABOOST lorsque celui-ci était utilisé comme générateur et on s'attendait à ce que les taux d'erreur augmentent un peu plus avec un générateur purement aléatoire. Au contraire, ces poids se révèlent mieux adaptés à ce mode de production des hypothèses.

Précisons que, dans le cas de la génération purement aléatoire de GLOBOOST, la fréquence et la couverture sont fortement liées : une hypothèse ayant un fort support sortira souvent, et inversement, une hypothèse fréquente couvre nécessairement beaucoup d'exemples. Ceci explique les résultats proches des poids basés sur ces quantités.

Finalement, c'est le poids *frequency* qui apparaît comme le meilleur de tous : que ce soit à 100 ou à 1000 étapes de *boosting*, ce poids gagne toutes ces confrontations et obtient le taux d'erreur moyen le plus faible.

<b>Problèmes</b>	<b>adalike</b>	<b>covering</b>	<b>frequency</b>	<b>uniform</b>
audiology	<b>26.27</b>	28.82	26.73	27.14
breast-cancer	3.86	4.34	<b>3.81</b>	3.82
car	<b>22.12</b>	22.22	22.25	22.27
cmc	54.19	<b>54.16</b>	54.25	54.24
crx	15.12	<b>14.17</b>	14.45	14.45
dermatology	7.42	5.84	5.49	<b>5.13</b>
ecoli	23.49	23.94	23.49	<b>23.46</b>
glass	<b>4.59</b>	4.73	<b>4.59</b>	4.73
hepatitis	<b>16.99</b>	17.33	18.04	18.18
horse-colic	19.74	19.52	<b>19.43</b>	<b>19.43</b>
house-votes-84	9.98	6.25	5.26	<b>5.00</b>
ionosphere	18.28	7.47	<b>7.22</b>	<b>7.22</b>
iris	<b>6.33</b>	6.73	6.67	6.67
pima	27.33	<b>26.65</b>	27.02	27.02
promoters	15.90	16.47	<b>14.38</b>	14.98
sonar	36.88	26.06	<b>25.85</b>	<b>25.85</b>
tic-tac-toe	<b>0.54</b>	<b>0.54</b>	0.58	2.00
vowel	28.81	29.07	<b>28.22</b>	<b>28.22</b>
wine	14.17	4.39	4.32	<b>4.27</b>
zoo	<b>3.92</b>	5.39	4.12	5.60
<i>Moyennes</i>	<i>17.80</i>	<i>16.20</i>	<b>15.81</b>	<i>15.98</i>

(a) Taux d'erreur.

	<b>adalike</b>	<b>covering</b>	<b>frequency</b>	<b>uniform</b>	<b>Bilan</b>
<b>adalike</b>	-	10-9	7-12	8-12	25-33
<b>covering</b>	9-10	-	6-14	7-12	22-36
<b>frequency</b>	12-7	14-6	-	8-5	<b>34-18</b>
<b>uniform</b>	12-8	12-7	5-8	-	29-23

(b) Confrontations.

**Tableau 5.** Résultats de GLOBOOST et des poids adalike, covering, frequency, uniform pour 100 étapes.

<b>Problèmes</b>	<b>adalike</b>	<b>covering</b>	<b>frequency</b>	<b>uniform</b>
audiology	18.10	24.35	<b>17.07</b>	20.67
breast-cancer	5.91	4.39	<b>3.65</b>	3.65
car	<b>10.45</b>	11.37	11.52	11.84
cmc	50.05	<b>47.02</b>	47.67	47.63
crx	24.15	<b>13.90</b>	13.93	13.91
dermatology	13.14	4.14	3.94	<b>3.08</b>
ecoli	23.72	19.48	<b>19.19</b>	19.28
glass	8.69	6.33	<b>4.55</b>	6.33
hepatitis	18.70	<b>17.16</b>	17.62	17.84
horse-colic	21.81	16.98	<b>16.54</b>	16.65
house-votes-84	29.72	4.38	5.00	<b>4.20</b>
ionosphere	57.47	<b>6.36</b>	7.05	7.07
iris	12.33	5.87	<b>5.67</b>	5.73
pima	30.15	<b>24.59</b>	24.88	24.88
promoters	27.60	10.88	<b>6.03</b>	9.07
sonar	49.21	<b>23.90</b>	23.96	23.96
tic-tac-toe	0.24	0.39	<b>0.14</b>	7.08
vowel	25.95	14.51	13.44	<b>13.42</b>
wine	47.49	<b>3.98</b>	4.37	4.03
zoo	<b>3.71</b>	6.44	4.05	4.74
<i>Moyennes</i>	<i>23.93</i>	<i>13.32</i>	<b>12.51</b>	<i>13.25</i>

(a) Taux d'erreur.

	<b>adalike</b>	<b>covering</b>	<b>frequency</b>	<b>uniform</b>	<b>Bilan</b>
<b>adalike</b>	-	4-16	2-18	4-16	10-50
<b>covering</b>	16-4	-	9-11	9-10	34-25
<b>frequency</b>	18-2	11-9	-	12-6	<b>41-17</b>
<b>uniform</b>	16-4	10-9	6-12	-	32-25

(b) Confrontations.

**Tableau 6.** Résultats de GLOBOOST et des poids adalike, covering, frequency, uniform pour 1000 étapes.

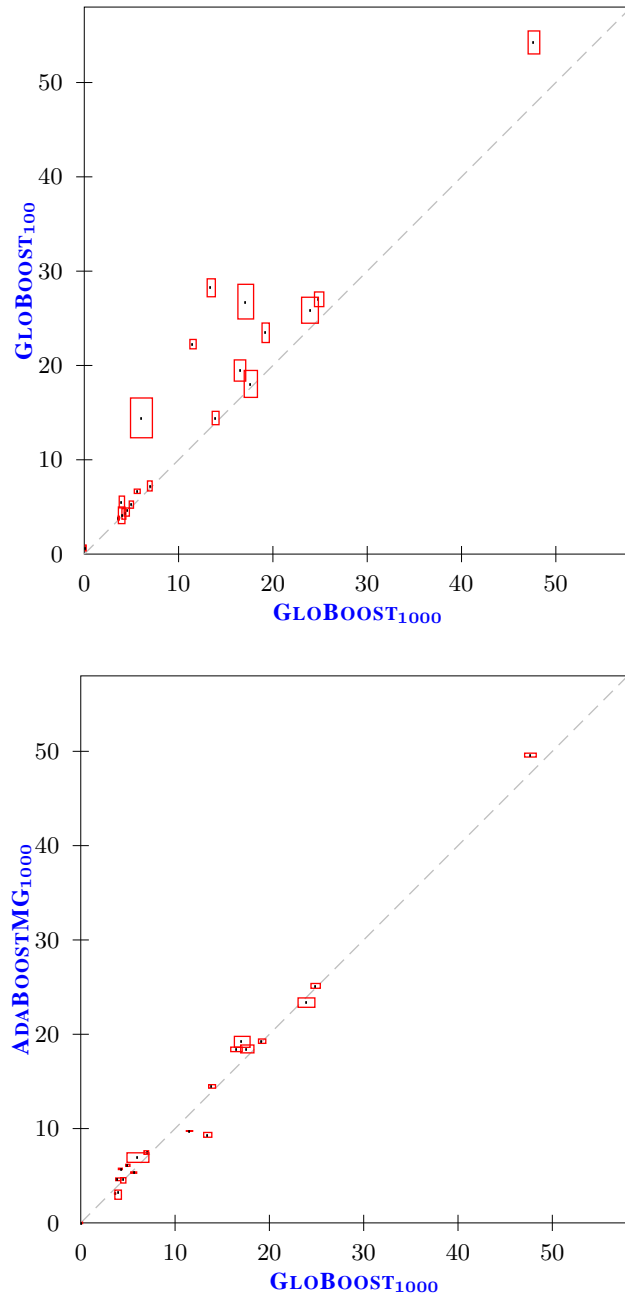
Problèmes	C4.5	GLOBo	ABDS	ABMG	GLOBOOST
audiology	18.20	20.76	<b>15.79</b>	19.22	17.07
breast-cancer	4.87	3.89	4.16	<b>3.10</b>	3.65
car	<b>7.69</b>	10.17	13.36	9.74	11.52
cmc	48.07	50.60	<b>44.87</b>	49.61	47.67
crx	14.79	16.50	16.53	14.47	<b>13.93</b>
dermatology	6.23	8.51	4.38	4.63	<b>3.94</b>
ecoli	15.89	23.88	<b>15.59</b>	19.25	19.19
glass	28.72	5.57	22.53	4.59	<b>4.55</b>
hepatitis	20.70	20.09	18.23	18.39	<b>17.62</b>
horse-colic	<b>13.63</b>	22.29	21.71	18.42	16.54
house-votes-84	<b>3.22</b>	7.39	5.31	6.12	5.00
ionosphere	7.96	8.74	12.19	7.44	<b>7.05</b>
iris	<b>5.33</b>	7.47	<b>5.33</b>	<b>5.33</b>	5.67
pima	29.29	27.04	25.39	25.13	<b>24.88</b>
promoters	18.17	22.60	18.67	7.00	<b>6.03</b>
sonar	28.97	31.09	<b>15.03</b>	23.34	23.96
tic-tac-toe	14.40	1.11	1.67	<b>0.00</b>	0.14
vowel	21.21	23.99	14.65	<b>9.30</b>	13.44
wine	8.83	8.94	18.07	5.73	<b>4.37</b>
zoo	7.51	5.55	5.25	<b>3.21</b>	4.05
<i>Moyennes</i>	<i>16.18</i>	<i>16.31</i>	<i>14.94</i>	<i>12.70</i>	<b>12.51</b>

(a) Taux d'erreur.

	C4.5	GLOBo	ABDS	ABMG	GLOBOOST	<i>Bilan</i>
<b>C4.5</b>	-	14-6	7-12	6-13	5-15	32-46
<b>GLOBo</b>	6-14	-	7-13	0-20	1-19	14-66
<b>ABDS</b>	12-7	13-7	-	7-12	5-15	37-41
<b>ABMG</b>	13-6	20-0	12-7	-	7-13	52-26
<b>GLOBOOST</b>	15-5	19-1	15-5	13-7	-	<b>62-18</b>

(b) Confrontations.

**Tableau 7.** Comparaison de GLOBOOST (frequency) avec C4.5, GLOBo et ADA-BOOSTDS et ADABOOSTMG. GLOBOOST et ADABOOSTMG utilisent ici 1000 étapes de génération.



**Figure 3.**  $GLOBBOOST_{1000}$  contre  $ADABOOSTMG_{1000}$  et  $GLOBBOOST_{100}$ .



## 6. Discussion

### 6.1. Bilan

Nous avons montré que les résultats sur le *boosting* laissaient une large place à l'investigation, en particulier sur les apprenants faibles qui peuvent être utilisés. Nous avons proposé d'utiliser dans ce rôle la production de *moindres généralisés corrects*. Devant la complexité de ce calcul et l'usage intensif que doit en faire ADABOOST, nous avons proposé une génération des hypothèses parallélisable et avons démontré que des poids simples pouvaient ensuite être affectés aux hypothèses. Toutes ces idées ont été validées expérimentalement.

Aujourd'hui, le système GLOBOOST implémente la production aléatoire vue à l'algorithme 5 avec le poids *frequency*. Ce poids offre de bonnes performances et n'occasionne pas de calculs importants après la génération.

Le tableau 7 et la figure 3 dressent le bilan de nos expérimentations pour cet algorithme. En résumé GLOBOOST présente des taux d'erreur et des écarts types faibles ; de plus, il supporte favorablement l'augmentation du nombre d'étapes de *boosting*. Ces résultats sont tout à fait comparables à ceux obtenus avec ADABOOSTMG, avec l'avantage que le calcul peut être réparti sur plusieurs machines.

### 6.2. Travaux apparentés

Les deux systèmes les plus proches de GLOBOOST sont de notre point de vue RISE [DOM 96] et SLIPPER [COH 99]. Le premier utilise également la construction d'un moindre généralisé comme brique de base (dans ce cas, il est permis au moindre généralisé d'être incorrect) et son originalité repose sur le mode de classification d'un exemple : RISE considère la règle la plus proche de l'exemple à classer, au sens d'une distance.

SLIPPER n'utilise pas les moindres généralisés mais des règles assez proches dans la forme et qui de plus sont boostées avec ADABOOST : dans SLIPPER, l'apprenant faible choisit la règle qui minimise la valeur de  $Z_t$  apparaissant dans l'algorithme 1, comme préconisé dans [SCH 99]. Cette stratégie permet de faire chuter très rapidement l'erreur observée (mais pas nécessairement l'erreur en généralisation). Notons simplement que, dans notre cadre, la correction des règles produites fait que l'erreur observée est nulle dès que suffisamment de règles ont été produites pour couvrir tous les exemples disponibles.

Au delà des méthodes elles-mêmes, notre démarche s'inscrit dans une tendance qui consiste à englober toutes les techniques de production et combinaison d'hypothèses dans un cadre plus large que celui du *boosting*.

Une première tentative dans ce sens a consisté à rassembler les méthodes de type *arcing* (pour *adaptively resample and combine*, [BRE 96a]), c'est-à-dire essentiellement le *bagging* et le *boosting*.

Plus proche de nous, [DIE 00b] décrit des expérimentations semblables aux nôtres mais dédiées aux arbres de décisions (comparaison de *boosting*, *bagging* et de production aléatoire), travail qui conduit finalement à la notion de *méthodes d'ensemble* [DIE 00a] ([MEI 03] apporte une nuance supplémentaire pour définir les techniques de type *Leveraging*).

Enfin, et d'un point de vue plus théorique, des travaux sont menés pour voir toutes ces méthodes comme des algorithmes de type *Monte Carlo* et par là expliquer leurs performances [ESP 03, ESP 04].

### 6.3. Perspectives

La question se pose maintenant de savoir si les résultats présentés dans cet article sont propres aux moindres généralisés corrects. La seule particularité de notre apprenant faible est de fournir des hypothèses qui ne se trompent pas ; on a vu que cela simplifiait le calcul des poids dans ADABOOST. Peut-être que cela rend les poids et les performances d'ADABOOST plus faciles à approcher. Il faudra également déterminer si autoriser un moindre généralisé à être incorrect permet à ADABOOST d'atteindre de meilleurs performances. Précisons que cela va dans le sens de notre souci d'efficacité : un apprenant faible incorrect n'a plus besoin de considérer tous les exemples mais seulement ceux de poids forts. Enfin, il faudra déterminer si, dans ce cadre, nous sommes encore capables de proposer des poids pertinents.

Dans l'idée d'améliorer l'apprenant faible fourni à ADABOOST, il semble aisé de contraindre un peu plus notre construction de moindres généralisés corrects pour focaliser sur les exemples de poids forts (même si c'est le principe de l'algorithme 4, cette mise en œuvre reste assez lâche).

Nous avons vu que ADABOOST et GLOBOOST, boostant des moindres généralisés corrects, s'amélioreraient avec l'augmentation du nombre d'hypothèses produites. Il faudra donc poursuivre dans cette voie, cela étant plus facile en pratique pour la version parallèle de GLOBOOST.

La recherche de nouveaux poids reste elle aussi ouverte. En particulier, il faut déterminer pourquoi notre poids *adalike* se comporte bien pour des étapes de *boosting* en nombre faible mais mal sur certains problèmes lorsque ce nombre augmente.

Une perspective qui nous semble importante est celle de la compréhensibilité du classifieur produit. Il s'agit maintenant de savoir si, à partir des hypothèses produites par ADABOOST ou GLOBOOST, on peut revenir à des théories compréhensibles comme celles de GLOBO. La solution la plus simple est d'appliquer la couverture minimale de GLOBO sur les hypothèses produites mais il est peu probable que les performances de la théorie obtenue restent proches de celles du système de votants. Il serait donc préférable de chercher un ensemble d'hypothèses d'assez petite taille pour être appréhendable par un humain mais en conservant au maximum le pouvoir prédictif du classifieur complet. Précisons que, tout proche de cette problématique, le

problème de fournir un nombre minimal de votants a été suggéré dans [QUI 99] et une solution proposée dans [LON 02], solution qui implique de calculer d’abord toutes les hypothèses faibles possibles (dans ce cas, des *decision stumps*).

Enfin, un dernier objectif est de poursuivre la parallélisation en distribuant également les données : la génération des hypothèses serait distribuée comme nous l’avons décrit, mais en plus elle se ferait sur des échantillons de l’ensemble de données initial. Cela conduira à des hypothèses incorrectes mais nous avons vu que cela pouvait être un avantage. À noter qu’une fois cette distribution mise en œuvre, le poids *frequency* sera un candidat privilégié puisqu’il ne nécessite pas de calcul sur l’ensemble de données, comme ce serait le cas par exemple pour le poids *covering*.

## 7. Bibliographie

- [BLA 98] BLAKE C., MERZ C., « UCI Repository of machine learning databases [<http://www.ics.uci.edu/~mllearn/MLRepository.html>] », 1998.
- [BRE 96a] BREIMAN L., « Bias, Variance, and Arcing Classifiers », 1996, Technical Report 460, Statistics Department, University of California.
- [BRE 96b] BREIMAN L., « Bagging Predictors », *Machine Learning*, vol. 24, n° 2, 1996, p. 123-140.
- [COH 99] COHEN W. W., SINGER Y., « A simple, fast, and effective rule learner », *Proceedings of the 16th National Conference on Artificial Intelligence (AAAI-99)*, AAAI/MIT Press, 1999, p. 335–342.
- [De 01] DE COMITÉ F., GILLERON R., TOMMASI M., « Learning multi-label alternating decision trees and applications », BISSON G., Ed., *Actes de la Conférence d’Apprentissage Automatique*, 2001, p. 195–210.
- [DIE 00a] DIETTERICH T. G., « Ensemble Methods in Machine Learning », KITTLER J., ROLI F., Eds., *First International Workshop on Multiple Classifier Systems*, vol. 1857 de *Lecture Notes in Computer Science*, Springer Verlag, 2000, p. 1-15.
- [DIE 00b] DIETTERICH T. G., « An Experimental Comparison of Three Methods for Constructing Ensembles of Decision Trees : Bagging, Boosting, and Randomization », *Machine Learning*, vol. 40, n° 2, 2000, p. 139–158.
- [DOM 96] DOMINGOS P., « Unifying Instance-Based and Rule-Based Induction », *Machine Learning*, vol. 24, n° 2, 1996, p. 141-168.
- [ESP 03] ESPOSITO R., SAITTA L., « Monte Carlo Theory as an Explanation of Bagging and Boosting », GOTTLÖB G., WALSH T., Eds., *Proceeding of the Eighteenth International Joint Conference on Artificial Intelligence*, Morgan Kaufman, 2003, p. 499–504.
- [ESP 04] ESPOSITO R., SAITTA L., « A Monte Carlo Analysis of Ensemble Classification », GREINER R., SCHUURMANS D., Eds., *Proceedings of the twenty-first International Conference on Machine Learning*, Banff, Canada, July 2004, ACM Press, New York, NY, p. 265-272.
- [FRE 95a] FREUND Y., « Boosting a Weak Learning Algorithm by Majority », *Information and Computation*, vol. 121, n° 2, 1995, p. 256–285.

- [FRE 95b] FREUND Y., SCHAPIRE R. E., « A decision-theoretic generalization of on-line learning and an application to boosting », *European Conference on Computational Learning Theory*, 1995, p. 23-37.
- [FRE 96] FREUND Y., SCHAPIRE R. E., « Experiments with a New Boosting Algorithm », *International Conference on Machine Learning*, 1996, p. 148–156.
- [FRE 97] FREUND Y., SCHAPIRE R. E., « A decision-theoretic generalization of on-line learning and an application to boosting », *Journal of Computer and System Sciences*, vol. 55, n° 1, 1997, p. 119-139.
- [FRE 98] FREUND Y., SCHAPIRE R. E., « Discussion of the paper Arcing Classifiers by Leo Breiman », *The Annals of Statistics*, vol. 26, 1998, p. 824–832.
- [GRO 98] GROVE A. J., SCHUURMANS D., « Boosting in the Limit : Maximizing the Margin of Learned Ensembles », *AAAI/IAAI*, 1998, p. 692–699.
- [HAR 99] HARRIES M., « Boosting a strong learner : evidence against the minimum margin », *Proc. 16th International Conf. on Machine Learning*, Morgan Kaufmann, San Francisco, CA, 1999, p. 171–180.
- [KEA 89] KEARNS M., VALIANT L. G., « Cryptographic limitations on learning Boolean formulae and finite automata », *Proceedings of the 21st Annual ACM Symposium on Theory of Computing*, 1989, p. 433–444.
- [KEA 94] KEARNS M. J., VAZIRANI U. V., *An Introduction to Computational Learning Theory*, MIT Press, 1994.
- [LON 02] LONG P. M., « Minimum Majority Classification and Boosting », *Proceedings of the The Eighteenth National Conference on Artificial Intelligence*, 2002.
- [MEI 03] MEIR R., RÄTSCH G., « An Introduction to Boosting and Leveraging », MENDELSON S., SMOLA A., Eds., *Advanced Lectures on Machine Learning*, n° 2600 LNAI, p. 119-184, Springer, January 2003.
- [QUI 99] QUINLAN J. R., « Some Elements of Machine Learning », 1999, Invited talk (ILP / ICML).
- [Rät 02] RÄTSCH G., WARMUTH M., « Maximizing the Margin with Boosting », KIVINEN J., SLOAN R. H., Eds., *Proceedings of the Annual Conference on Computational Learning Theory (COLT)*, vol. 2375 de *Lecture Notes in Computer Science*, Springer, 2002, p. 334–350.
- [Rät 03] RÄTSCH G., WARMUTH M. K., « Efficient Margin Maximizing with Boosting », 2003, submitted to *Journal of Machine Learning Research (JMLR)*.
- [SCH 90] SCHAPIRE R. E., « The Strength of Weak Learnability », *Machine Learning*, vol. 5, 1990, p. 197-227.
- [SCH 97] SCHAPIRE R. E., FREUND Y., BARTLETT P., LEE W. S., « Boosting the margin : a new explanation for the effectiveness of voting methods », *Proc. 14th International Conference on Machine Learning (ICML)*, Morgan Kaufmann, 1997, p. 322–330.
- [SCH 99] SCHAPIRE R. E., SINGER Y., « Improved Boosting Algorithms using Confidence-Rated Predictions », *Machine Learning*, vol. 37, n° 3, 1999, p. 297–336.
- [SCH 00] SCHAPIRE R. E., SINGER Y., « BoosTexter : A Boosting-based System for Text Categorization », *Machine Learning*, vol. 39, n° 2/3, 2000, p. 135-168.
- [SRI 97] SRINIVASAN A., KING R. D., MUGGLETON S. H., STERNBERG M. J. E., « The Predictive Toxicology Evaluation Challenge », *Proceedings of the 15th International Joint*

*Conference on Artificial Intelligence*, Morgan Kaufmann, 1997, p. 4–9.

- [SRI 99] SRINIVASAN A., KING R., BRISTOL D., « An assessment of submissions made to the Predictive Toxicology Evaluation Challenge », *Proceedings of the 16th International Joint Conference on Artificial Intelligence*, Morgan Kaufmann, 1999, p. 270–276.
- [TOR 99] TORRE F., « GloBo : un algorithme stochastique pour l'apprentissage supervisé et non-supervisé », SEBAG M., Ed., *Actes de la Première Conférence d'Apprentissage*, 1999, p. 161–168.
- [VAL 84] VALIANT L. G., « A theory of the Learnable », *Communications of the ACM*, vol. 27, 1984, p. 1134-1142.