

DATA MINING 2

TP Ensembles de classifieurs

Thomas ROBERT

1 Le *bootstrap* / *bagging*

1.1 Principe

La première partie du TP consiste à appliquer la méthode de *bootstrap* / *bagging*. De manière générale, cette méthode consiste à tirer B jeux d'apprentissage par tirages aléatoires avec remise sur le jeu de données initial. Le résultat final consiste à « mélanger » les B résultats unitaires obtenus pour chaque jeu d'apprentissage afin d'obtenir un résultat plus stable et plus performant en généralisation.

Le « mélange » dépend de ce qu'on évalue. Par *bootstrap* on évalue une valeur (par exemple l'écart-type), le « mélange » consiste donc dans ce cas à faire la moyenne des écarts-type. Par *bagging*, on évalue un classifieur cette fois, le « mélange » pourra par exemple être un vote majoritaire (mieux vaut dans ce cas choisir B astucieusement pour éviter des égalités...).

Notons également que le *bootstrap* permet également de calculer une sorte d'intervalle de confiance sur le résultat que l'on obtient. Il suffit de calculer l'écart-type de l'ensemble des B valeurs trouvées. Pour le *bagging*, un principe similaire permet de calculer l'erreur en test (données non tirées dans l'apprentissage) pour chaque tirage, et l'écart-type de cette erreur.

1.2 Application du *bootstrap* pour le calcul de *skewness*

On commence par estimer la *skewness* d'un jeu de données par *bootstrap*. (Voir annexe 4.1.2). On choisit $B = 50$ tirages comme recommandé tirés par la fonction `tireBootstrap` (voir annexe 4.1.1).

On obtient le résultat suivant : $\gamma = 0.52 \pm 0.072$

1.3 Application du *bagging* pour calcul d'un mélange de kppv et arbre de décision

On applique cette fois la méthode du *bagging* pour déterminer un mélange de classifieur. On code donc une fonction `[tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = baggingEval(X, Y, Xtest, Ytest, B, classifieur, estimateur)` (voir annexe 4.2.1) dont les paramètres sont :

- Entrée :
 - `X, Y, Xtest, Ytest` : jeu d'apprentissage et de test
 - `B` : nombre d'itération pour le *bagging*
 - `classifieur` : *function handler* prtools vers un classifieur (exemple : `@knn` ou `@treec`)
 - `estimateur` : *function handler* prtools vers un estimateur de classifieur (exemple : `@knn_map` ou `@tree_map`)
- Sortie :
 - `tauxErreurOOBMoy` : taux d'erreur *out-of-bag* moyen
 - `ecartErreurOOB` : écart type du taux d'erreur *out-of-bag* moyen
 - `tauxErreurTestBagging` : taux d'erreur du mélange de classifieur produit par *bagging* sur l'ensemble de test
 - `tauxErreurTestSansBagging` : taux d'erreur en test du classifieur calculé seul sur l'ensemble d'apprentissage

Cette méthode a été testée sur divers jeux de données avec divers paramètres (voir annexe 4.2.2), pour les classifieurs *k-plus proches voisins* et *arbre de décision*. On constate qu'en moyenne, le *bagging* permet d'obtenir un résultat légèrement voire bien meilleur que le classifieur seul.

Le résultat le plus impressionnant est sur le jeu de données *pima* avec la méthode *knn*, où le taux d'erreur en *bagging* est de 0,3% contre 21,6% sans *bagging*.

Jeu	Taux erreur OOB	Taux erreur test bagging	Taux erreur test sans bagging
Clown - knn	$8.2 \pm 6.8 \%$	9,3 %	10,7 %
Clown - tree	$8,2 \pm 6,4 \%$	8,6 %	9,3 %
Pima - knn	$32,7 \pm 2,2 \%$	0,37 %	21,6 %
Pima - tree	$34,4 \pm 3,0 \%$	0 %	0 %
Satimage - knn	$10,5 \pm 0,64 \%$	9,7 %	9,8 %
Satimage - tree	$18,0 \pm 0,9 \%$	10,5 %	17,0 %

2 Le Boosting

2.1 Principe

La seconde méthode de mélange de classifieurs testée est la méthode de *boosting*, qui consiste à pondérer un ensemble de classifieurs et à pondérer les données au fur et à mesure du calcul des classifieurs afin d'avoir en moyenne la bonne estimation.

2.2 Application

On veut tester la *boosting* avec un ensemble de souches binaire.

Afin de tester la méthode d'*AdaBoost*, on code une fonction (voir annexe refboostingEval)

```
[YappClassCumule, tauxErreurApp, YtestClassCumule, tauxErreurTest] = AdaBoostEval(Xapp, Yapp, Xtest, Ytest, T)
```

dont les paramètres sont les suivants :

- Entrée :
 - **Xapp, Yapp, Xtest, Ytest** : jeu d'apprentissage et de test
 - **T** : nombre d'itérations
- Sortie :
 - **YappClassCumule** : Résultat de l'évaluation sur l'apprentissage au fur et à mesure des itérations
 - **tauxErreurApp** : Taux d'erreur sur le test au fur et à mesure des itérations
 - **YtestClassCumule** : Résultat de l'évaluation sur l'apprentissage au fur et à mesure des itérations
 - **tauxErreurTest** : Taux d'erreur sur le test au fur et à mesure des itérations

On lance cette fonction sur plusieurs jeux de données (voir figure 1). Le fait d'avoir des résultats au fur et à mesure des itérations permet de voir à partir de quand on commence à faire du sur-apprentissage sur le jeu de données.

On constate que sur le jeu de données clown, le boosting est efficace et donne son meilleur résultat en test pour $T = 31$. Par contre, sur le jeu *pima*, le boosting fonctionne très mal et se stabilise à un taux d'erreur plutôt élevé.

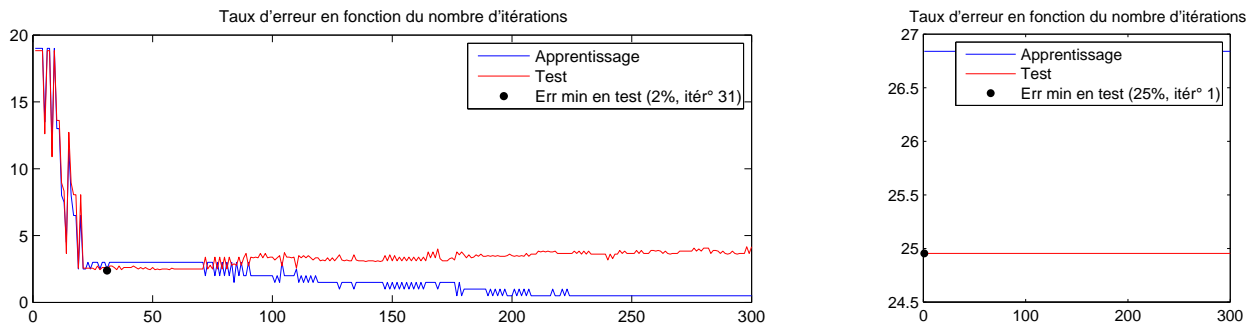


FIGURE 1 – Résultats du *boosting* sur les jeux *clown* (gauche) et *pima* (droite)

3 Les *Random Forests*

3.1 Principe

Le but des forêts aléatoires est de stabiliser les arbres de décision. Pour cela, on ajoute de l'aléatoire aux données d'apprentissage et on construit donc plusieurs arbres différents afin de construire une « forêt » d'arbres de décision plus robuste en généralisation qu'un arbre de décision seul, qui est lui instable.

Il y a de nombreuses façons d'ajouter de l'aléatoire lors de la construction d'une forêt aléatoire. Dans notre cas, pour chaque arbre nous allons choisir un sous-ensemble des composantes initiales et faire un tirage *bootstrap* des données qui serviront à faire l'apprentissage.

3.2 Mise en oeuvre

Le codage a été fait en plusieurs étapes : codage d'une fonction de calcul et d'évaluation d'arbre de décision (voir annexe 4.4), codage d'une fonction de calcul et d'évaluation de forêt aléatoire, et enfin test sur 2 jeux de test (voir annexe 4.5).

Il est important de noter que la construction d'une forêt aléatoire se fait à partir de 6 paramètres :

- **X, Y** : Données d'apprentissage
- **hauteurMax** : La hauteur maximale d'un arbre
- **pureteSeuil** : La pureté minimale à atteindre avant de s'arrêter (sauf si **hauteurMax** dépassée)
- **nbArbres** : Le nombre d'arbres dans la forêt
- **nbFeaturesParArbre** : Le nombre de composantes utilisées par arbre

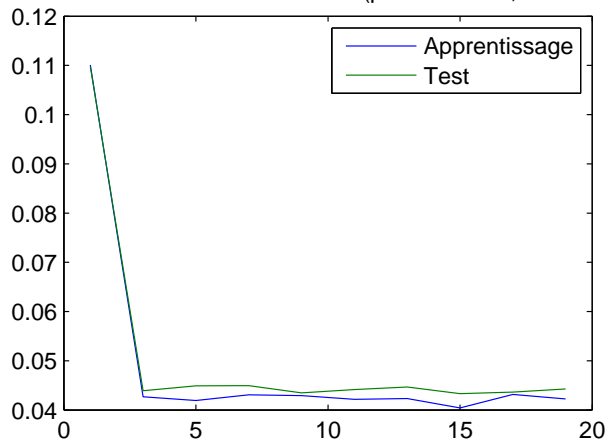
Ces divers paramètres ont donc dû être testés. Pour cela, j'ai choisi de faire varier ces paramètres un par un et de tracer l'évolution de l'erreur en fonction du paramètre testé. (Pour stabiliser la mesure, on fait 20 fois la génération de la forêt avec les mêmes paramètres et on moyenne les 20 mesures d'erreur.)

3.3 Résultats

3.3.1 Jeu de données *clown*

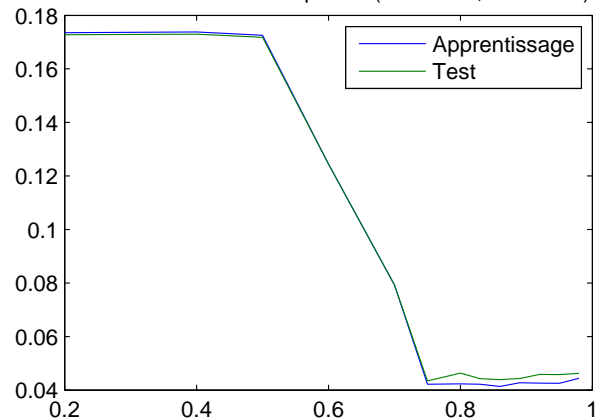
On teste la méthode avec divers paramètres et on obtient les résultats suivants :

Erreur en fonction de la hauteur (pureté = 95%, 49 arbres)

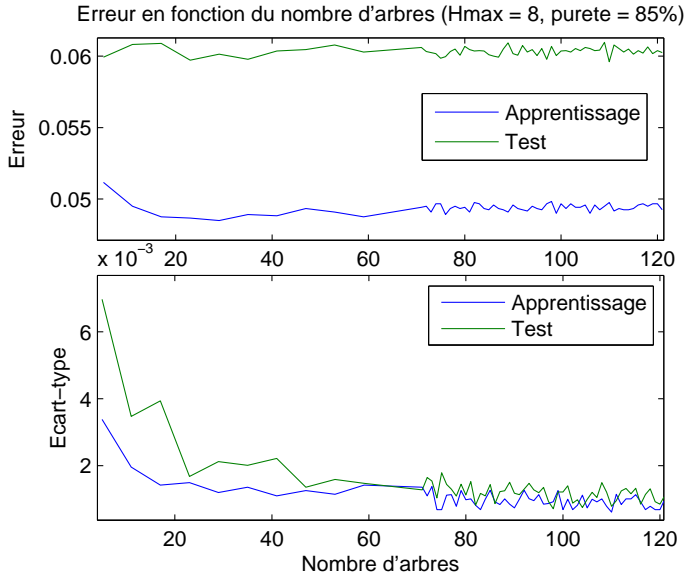


Hauteur Une hauteur supérieure à 3 semble suffisante pour avoir de bons résultats. On choisira 8 pour la suite.

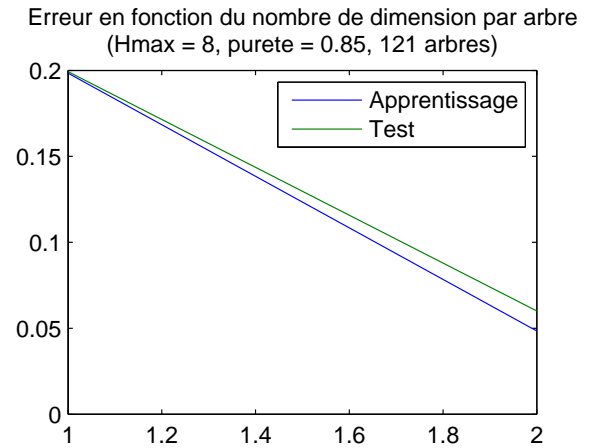
Erreur en fonction de la pureté (Hmax = 8, 49 arbres)



Pureté On voit apparaître un palier de 0,2 à 0,5 puis une chute relativement linéaire et un nouveau plateau de 0,75 à 1. On choisira donc pour la suite une pureté minimale de 0,85.



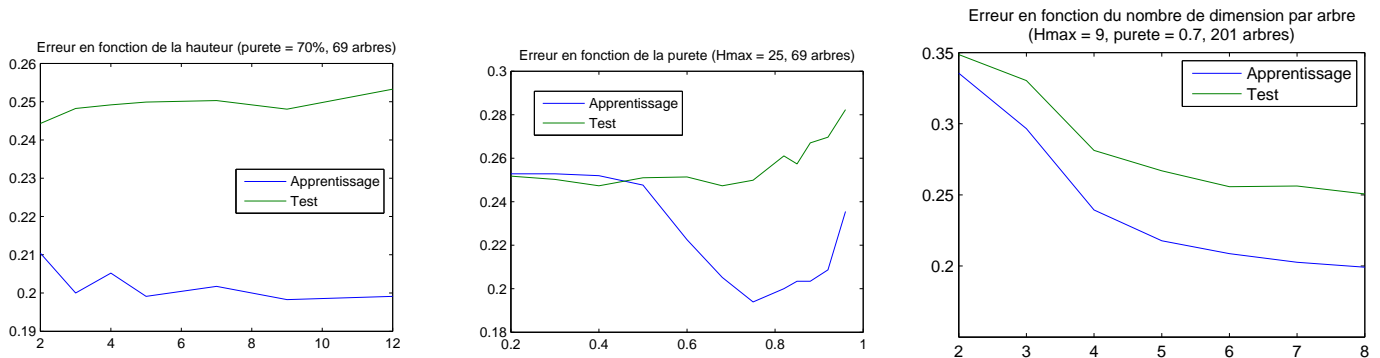
Nombre d'arbres Le nombre d'arbre a peu d'impact sur les résultats moyens. Par contre, on remarque que l'écart-type indique que l'erreur varie beaucoup d'une forêt à l'autre s'il n'y a pas assez d'arbres dans la forêt pour stabiliser l'aléatoire introduit. A partir de 50 arbres, le résultat commence à être relativement stable.



Nombre de dimensions Nous n'avons ici que 2 dimensions, il semble donc étrange de n'en garder qu'une seule par arbre. Cela a quand même été essayé avec 121 arbres dans la forêt (moins de dimensions à la fois nécessite logiquement plus d'arbres). Le résultat indique qu'il est clairement mieux de choisir de garder les 2 dimensions.

Conclusion Finalement, on choisi dont une hauteur maximale de **8 niveaux**, une pureté minimale de **85%**, **49 arbres** et **2 dimensions** par arbre pour ce jeu de données. On obtient une erreur de **4,9% en apprentissage** et de **5,9% en test**. La proximité entre ces deux valeurs montre une bonne robustesse à la généralisation.

3.3.2 Jeu de données *pima*



Hauteur L'erreur évolue peu en fonction du nombre de niveaux. On choisi 9 maximum dans l'arbre pour la suite.

Pureté L'erreur en test dépend étrangement peu de la pureté. On remarque du sur-apprentissage pour une pureté supérieur à 75%. On choisira 70% de pureté pour la suite.

Nombre de dimensions Il y a 8 dimensions, il pourrait donc être intéressant de ne pas toutes les sélectionner. On voit cependant que les résultats restent meilleurs avec les 8 dimensions. Cependant, l'erreur varie relativement peu entre 5 et 8 dimensions.

Conclusion Finalement, on choisi dont une hauteur maximale de **9 niveaux**, une pureté minimale de **65%**, **69 arbres** et **6 dimensions** par arbre pour ce jeu de données. On obtient une erreur de **20% en apprentissage** et de **25% en test**.

4 Annexes (code Matlab)

Contents

4.1 Bootstrap	5
4.1.1 Tirage de bootstrap	5
4.1.2 Bootstrap sur Skewness	5
4.2 Bagging	5
4.2.1 Evaluation du bagging	5
4.2.2 Lancement d'évaluations de bagging	6
4.3 Boosting	7
4.3.1 Evaluation du boosting	7
4.3.2 Lancement d'évaluations de boosting	8
4.4 Arbres de décision	9
4.4.1 Construction d'un arbre de décision	9
4.4.2 Évaluation d'un arbre de décision	9
4.5 Forêts aléatoires	9
4.5.1 Construction d'une forêt aléatoire	9
4.5.2 Évaluation d'une forêt aléatoire	10
4.5.3 Test d'une forêt aléatoire	10
4.5.4 Test des forêts aléatoires	10

4.1 Bootstrap

4.1.1 Tirage de bootstrap

```

1 % Tire des indices de bootstrap
2 %
3 % n : nombre d'éléments dans l'ensemble
4 % m : nombre d'éléments à tirer
5 function [bag, obag] = tireBootstrap(n, m)
6
7     bag = randi(n, m, 1);
8     obag = setdiff(1:n, bag);
9
10 end

```

4.1.2 Bootstrap sur Skewness

```

1 % Chargement CSV
2 X = csvread('skewnormal.csv');
3 n = length(X);
4 B = 50;
5
6 % On calcule B skewness
7 skewVals = zeros(B,1);
8 for i = 1:B
9     skewVals(i) = skewness(X(tireBootstrap(n,n)));
10 end
11 skew = mean(skewVals)
12 errSkew = std(skewVals)

```

4.2 Bagging

4.2.1 Evaluation du bagging

```

1 function [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
2     baggingEval(X, Y, Xtest, Ytest, B, classifieur, estimateur)
3
4 nX = length(X);
5 nXtest = length(Xtest);
6

```

```

7 % BAGGING
8 tauxErreurs = zeros(B,1);
9 classifieurs = cell(1,B);
10 for i = 1:B
11     % Tirage des jeux
12     [bag, obag] = tireBootstrap(nX, nX);
13     Xapp = X(bag, :);
14     Yapp = Y(bag);
15     Xval = X(obag, :);
16     Yval = Y(obag);
17
18     % Création jeux
19     dataApp = prdataset(Xapp, Yapp);
20     dataVal = prdataset(Xval, Yval);
21
22     % Classification
23     classifieurs{i} = classifieur(dataApp);
24     probas = estimateur(dataVal, classifieurs{i});
25     [~, classesVal] = max(probas.data, [], 2);
26     tauxErreurs(i) = sum(dataVal.nlab ~= classesVal)/length(classesVal)*100;
27 end
28 tauxErreurOOBMoy = mean(tauxErreurs);
29 ecartErreurOOB = std(tauxErreurs);
30
31 % ENSEMBLE DE CLASSIFIEURS
32 dataTest = prdataset(Xtest, Ytest);
33 target = dataTest.nlab;
34
35 classes = zeros(nXtest, B);
36 for i = 1:B
37     probas = estimateur(dataTest, classifieurs{i});
38     [~, classesTest] = max(probas.data, [], 2);
39     classes(:,i) = classesTest;
40 end
41
42 Ytarget = dataTest.nlab;
43 Yclassif = mode(classes,2);
44
45 tauxErreurTestBagging = sum(Ytarget ~= Yclassif)/nXtest*100;
46
47 % SANS BAGGING
48
49 % Création du jeu
50 dataApp = prdataset(X, Y);
51 dataTest = prdataset(Xtest, Ytest);
52
53 % Classification
54 classifieurSeul = classifieur(dataApp);
55 probas = estimateur(dataTest, classifieurSeul);
56 [~, classesSeul] = max(probas.data, [], 2);
57 Ytarget = dataTest.nlab;
58 tauxErreurTestSansBagging = sum(Ytarget ~= classesSeul)/length(classesSeul)*100;

```

4.2.2 Lancement d'évaluations de bagging

```

1 %% Clown data
2
3 clear all
4
5 nX=200;
6 X = zeros(nX,2);
7
8 X(1:nX/2, :) = randn(nX/2,2) + repmat([0 6], nX/2, 1);
9 X(nX/2+1:nX,1) = (rand(nX/2,1) - 0.5) * 6;
10 X(nX/2+1:nX,2) = X(nX/2+1:nX,1).^2 + 0.7*randn(nX/2,1);
11
12 Y = [ones(nX/2,1); zeros(nX/2,1)];
13
14 % découpage en apprentissage et test
15 [X, Y, Xtest, Ytest] = splitdata(X, Y, 0.3);
16 nX=length(X);
17 nXtest=length(Xtest);
18 plot(X(Y==1,1), X(Y==1,2), 'r'); hold on;
19 plot(X(Y==0,1), X(Y==0,2), 'b');
20 plot(Xtest(Ytest==1,1), Xtest(Ytest==1,2), 'or'); hold on;

```

```
21 plot(Xtest(Ytest==0,1), Xtest(Ytest==0,2), 'ob');
22 title('Jeu de test');
23
24 % Bagging
25
26 disp('K plus proches voisins');
27
28 [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
29     baggingEval(X, Y, Xtest, Ytest, 49, @knn, @knn_map)
30
31 disp('Arbre de decision');
32
33 [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
34     baggingEval(X, Y, Xtest, Ytest, 49, @treec, @tree_map)
35
36 %% Jeu de données pima
37
38 data = csvread('pima-indians-diabetes.data');
39
40 X = data(:,1:end-1);
41 Y = data(:,end);
42
43 % découpage en apprentissage et test
44 [Xapp, Yapp, Xtest, Ytest] = splitdata(X, Y, 0.3);
45
46 %%
47 disp('K plus proches voisins');
48
49 [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
50     baggingEval(X, Y, Xtest, Ytest, 49, @knn, @knn_map)
51 %%
52 disp('Arbre de decision');
53
54 [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
55     baggingEval(X, Y, Xtest, Ytest, 49, @treec, @tree_map)
56
57
58 %%
59
60 dataApp = tblread('sat.trn');
61 dataTest = tblread('sat.tst');
62
63 X = dataApp(:,1:end-1);
64 Y = dataApp(:,end);
65 Xtest = dataTest(:,1:end-1);
66 Ytest = dataTest(:,end);
67
68 %%
69 disp('K plus proches voisins');
70
71 [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
72     baggingEval(X, Y, Xtest, Ytest, 49, @knn, @knn_map)
73 %%
74 disp('Arbre de decision');
75
76 [tauxErreurOOBMoy, ecartErreurOOB, tauxErreurTestBagging, tauxErreurTestSansBagging] = ...
77     baggingEval(X, Y, Xtest, Ytest, 49, @treec, @tree_map)
78
79
80 %%
```

4.3 Boosting

4.3.1 Evaluation du boosting

```
1 function [YappClassCumule, tauxErreurApp, YtestClassCumule, tauxErreurTest] = ...
2     AdaBoostEval(Xapp, Yapp, Xtest, Ytest, T)
3
4 % tailles
5 nXapp = size(Xapp,1);
6 nXtest = size(Xtest,1);
7
8 % initialisations
9 D = zeros(nXapp, T+1); % poids
```

```

10 h = cell(T,1); % classifieurs
11 eps = zeros(T,1); % erreurs pondérées
12 alpha = zeros(T,1); % pondérations de classifieurs
13 YappClass = zeros(nXapp, T); % classifications
14 YtestClass = zeros(nXtest, T); % classifications
15
16 % initialisation des poids pour la première itération
17 D(:,1) = 1/nXapp;
18
19 % pour chaque itération
20 for t=1:T
21
22     % apprentissage du classifieur
23     h{t} = souchebinairetrain(Xapp, Yapp, D(:,t));
24
25     % evaluation avec le classifieur
26     YappClass(:,t) = souchebinaireval(h{t}, Xapp);
27
28     % calcul de l'erreur pondérée
29     eps(t) = sum(D(YappClass(:,t) ~= Yapp, t));
30
31     % calcul du coef de pondération du classifieur
32     alpha(t) = 1/2 * log( (1 - eps(t)) / eps(t) );
33
34     % calcul des pondérations pour l'étape suivante
35     D(:, t+1) = D(:, t) .* exp( -alpha(t) * YappClass(:,t) .* Yapp );
36     D(:, t+1) = D(:, t+1) / sum(D(:, t+1));
37
38     % evaluation avec le classifieur sur l'ensemble de test
39     YtestClass(:,t) = souchebinaireval(h{t}, Xtest);
40 end
41
42 % évaluation des performances au fur et à mesure des itérations
43 % en apprentissage
44 YappClassCumule = sign(cumsum( (YappClass .* repmat(alpha', nXapp, 1)) , 2));
45 tauxErreurApp = sum(YappClassCumule ~= repmat(Yapp, 1, T))/nXapp * 100;
46 % en test
47 YtestClassCumule = sign(cumsum( (YtestClass .* repmat(alpha', nXtest, 1)) , 2));
48 tauxErreurTest = sum(YtestClassCumule ~= repmat(Ytest, 1, T))/nXtest * 100;
49 [tauxErreurTestMin, tauxErreurTestMinInd] = min(tauxErreurTest);
50
51 % affichage des résultats
52 figure();
53 plot(tauxErreurApp, '-b'); hold on;
54 plot(tauxErreurTest, '-r');
55 plot(tauxErreurTestMinInd, tauxErreurTestMin, '.k', 'MarkerSize', 15);
56
57 title('Taux d''erreur en fonction du nombre d''itérations');
58 legend('Apprentissage', 'Test', ['Err min en test (' int2str(round(tauxErreurTestMin))...
59     '%, itér° ' int2str(round(tauxErreurTestMinInd)) ')'] );

```

4.3.2 Lancement d'évaluations de boosting

```

1 %% Jeu clown
2
3 clear all
4
5 nX=2000;
6 X = zeros(nX,2);
7
8 X(1:nX/2, :) = randn(nX/2,2) + repmat([0 6], nX/2, 1);
9 X(nX/2+1:nX,1) = (rand(nX/2,1) - 0.5) * 6;
10 X(nX/2+1:nX,2) = X(nX/2+1:nX,1).^2 + 0.7*randn(nX/2,1);
11
12 Y = [-ones(nX/2,1) ; ones(nX/2,1)];
13
14 % découpage en apprentissage et test
15 [Xapp, Yapp, Xtest, Ytest] = splitdata(X, Y, 0.1);
16 plot(X(Y==1,1), X(Y==1,2), '.r'); hold on;
17 plot(X(Y==-1,1), X(Y==-1,2), '.b');
18 plot(Xtest(Ytest==1,1), Xtest(Ytest==1,2), 'or'); hold on;
19 plot(Xtest(Ytest==-1,1), Xtest(Ytest==-1,2), 'ob');
20 title('Jeu de test');
21
22 AdaBoostEval(Xapp, Yapp, Xtest, Ytest, 300);

```



```
23
24 %% Jeu pima
25
26 clear all
27
28 data = csvread('pima-indians-diabetes.data');
29
30 X = data(:,1:end-1);
31 Y = data(:,end);
32 Y(Y==0) = -1;
33
34 % découpage en apprentissage et test
35 [Xapp, Yapp, Xtest, Ytest] = splitdata(X, Y, 0.3);
36
37 AdaBoostEval(Xapp, Yapp, Xtest, Ytest, 300);
```

4.4 Arbres de décision

4.4.1 Construction d'un arbre de décision

```
1 function tree = decisionTreeTrain(X, Y, hauteurRestante, pureteSeuil)
2
3 % taille
4 [n, p] = size(X);
5
6 % création d'un noeud
7 tree = struct();
8
9 % décision
10 tree.decision = souchebinairetrain(X, Y, ones(n,1));
11
12 % calcul des résultats dans les fils
13 Yhat = souchebinaireval(tree.decision, X);
14
15 % si des exemples dans le fils gauche ("-1") sont mal classés
16 Xgauche = X(Yhat == -1, :);
17 Ygauche = Y(Yhat == -1);
18 pureteG = abs(sum(Ygauche == 1) - sum(Ygauche == -1))/length(Ygauche);
19
20 if (pureteG < pureteSeuil && hauteurRestante > 0)
21     tree.filsGauche = decisionTreeTrain(Xgauche, Ygauche, hauteurRestante-1, pureteSeuil);
22 end
23
24 % si des exemples dans le fils droit ("1") sont mal classés
25 Xdroit = X(Yhat == 1, :);
26 Ydroit = Y(Yhat == 1);
27 pureteD = abs(sum(Ydroit == 1) - sum(Ydroit == -1))/length(Ydroit);
28
29 if (pureteD < pureteSeuil && hauteurRestante > 0)
30     tree.filsDroit = decisionTreeTrain(Xdroit, Ydroit, hauteurRestante-1, pureteSeuil);
31 end
```

4.4.2 Évaluation d'un arbre de décision

```
1 function Yhat = decisionTreeVal(tree, X)
2
3 Yhat = souchebinaireval(tree.decision, X);
4
5 % si fils gauche existant, recalcul
6 if (isfield(tree, 'filsGauche'))
7     Yhat(Yhat == -1) = decisionTreeVal(tree.filsGauche, X(Yhat == -1, :));
8 end
9
10 % si fils droit existant, recalcul
11 if (isfield(tree, 'filsDroit'))
12     Yhat(Yhat == 1) = decisionTreeVal(tree.filsDroit, X(Yhat == 1, :));
13 end
```

4.5 Forêts aléatoires

4.5.1 Construction d'une forêt aléatoire

```

1 function forest = randomForestTrain(X, Y, ...
2     hauteurMax, pureteSeuil, ...
3     nbArbres, nbFeaturesParArbre)
4
5 % taille donnees
6 [n,p] = size(X);
7
8 forest = cell(nbArbres, 1);
9
10 for i = 1:nbArbres
11
12     % features de l'arbre i
13     forest{i}.features = randperm(p, nbFeaturesParArbre);
14
15     % tirage du bagging
16     [bag, obag] = tireBootstrap(n, n);
17
18     % arbre i
19     forest{i}.tree = decisionTreeTrain(...
20         X(bag, forest{i}.features), Y(bag), ...
21         hauteurMax, pureteSeuil);
22
23 end

```

4.5.2 Évaluation d'une forêt aléatoire

```

1 function Yhat = randomForestVal(forest, X)
2
3 [n,p] = size(X);
4 L = length(forest);
5
6 Yhats = zeros(n,L);
7
8 for i = 1:L
9     Yhats(:,i) = decisionTreeVal(forest{i}.tree, X(:,forest{i}.features));
10 end
11
12 Yhat = mode(Yhats,2);

```

4.5.3 Test d'une forêt aléatoire

```

1 function [tauxErreurTrain, tauxErreurTest, stdTrain, stdTest] = randomForestEval(X, Y, Xtest, Ytest,
2     hmax, purete, nbArbres, nbDim, nbIte)
3
4 if (nargin < 9)
5     nbIte = 20;
6 end
7
8 tauxErreurTrainMat = zeros(1,nbIte);
9 tauxErreurTestMat = zeros(1,nbIte);
10
11 for i=1:nbIte
12     forest = randomForestTrain(X, Y, hmax, purete, nbArbres, nbDim);
13
14     Yhat = randomForestVal(forest, X);
15     tauxErreurTrainMat(i) = sum(Yhat ~= Y) / length(Y);
16
17     Ytestthat = randomForestVal(forest, Xtest);
18     tauxErreurTestMat(i) = sum(Ytestthat ~= Ytest) / length(Ytest);
19 end
20
21 tauxErreurTrain = mean(tauxErreurTrainMat);
22 tauxErreurTest = mean(tauxErreurTestMat);
23 stdTrain = std(tauxErreurTrainMat);
24 stdTest = std(tauxErreurTestMat);

```

4.5.4 Test des forêts aléatoires

```

1 %% Clown data
2
3 clear all
4
5 nX=2000;

```

```

6 X = zeros(nX,2);
7
8 X(1:nX/2, :) = randn(nX/2,2) + repmat([0 6], nX/2, 1);
9 X(nX/2+1:nX,1) = (rand(nX/2,1) - 0.5) * 6;
10 X(nX/2+1:nX,2) = X(nX/2+1:nX,1).^2 + 0.7*randn(nX/2,1);
11
12 Y = [ones(nX/2,1) ; -ones(nX/2,1)];
13
14 % découpage en apprentissage et test
15 [X, Y, Xtest, Ytest] = splitdata(X, Y, 0.3);
16 nX=length(X);
17 nXtest=length(Xtest);
18 plot(X(Y==1,1), X(Y==1,2), 'r'); hold on;
19 plot(X(Y==-1,1), X(Y==-1,2), 'b');
20 plot(Xtest(Ytest==1,1), Xtest(Ytest==1,2), 'or'); hold on;
21 plot(Xtest(Ytest==-1,1), Xtest(Ytest==-1,2), 'ob');
22 title('Jeu de test');
23
24 %%
25
26 hauteurVals = 1:2:20;
27 erreurs = zeros(length(hauteurVals), 2);
28
29 for i = 1:length(hauteurVals)
30
31     hauteur = hauteurVals(i);
32     [a,b] = randomForestEval(X, Y, Xtest, Ytest, hauteur, 0.95, 49, 2);
33     erreurs(i,:) = [a b];
34
35 end
36
37 figure
38 plot(hauteurVals, erreurs);
39 legend('Apprentissage', 'Test');
40 title('Erreur en fonction de la hauteur (purete = 95%, 49 arbres)');
41
42 % choix = 8
43
44 %%
45
46 pureteVals = [0.2 0.4 0.5 0.6 0.7 0.75 0.8:0.03:0.98];
47 erreurs = zeros(length(pureteVals), 2);
48
49 for i = 1:length(pureteVals)
50
51     purete = pureteVals(i);
52     [a,b] = randomForestEval(X, Y, Xtest, Ytest, 8, purete, 49, 2);
53     erreurs(i,:) = [a b];
54
55 end
56
57 figure
58 plot(pureteVals, erreurs);
59 legend('Apprentissage', 'Test');
60 title('Erreur en fonction de la purete (Hmax = 8, 49 arbres)');
61
62 % choix = 0.85
63
64 %%
65
66 nbArbresVals = [5:6:60 71:121];
67 erreurs = zeros(length(nbArbresVals), 4);
68
69 for i = 1:length(nbArbresVals)
70
71     nbArbres = nbArbresVals(i);
72     [a,b,c,d] = randomForestEval(X, Y, Xtest, Ytest, 8, 0.85, nbArbres, 2);
73     erreurs(i,:) = [a b c d];
74
75 end
76
77 %%
78
79 figure
80 subplot(2,1,1);

```

```
81 plot(nbArbresVals, erreurs(:,1:2));
82 legend('Apprentissage', 'Test');
83 title('Erreur en fonction du nombre d\'arbres (Hmax = 8, purete = 85%)');
84 ylabel('Erreur');
85
86 subplot(2,1,2);
87 plot(nbArbresVals, erreurs(:,3:4));
88 ylabel('Ecart-type');
89 xlabel('Nombre d\'arbres');
90 legend('Apprentissage', 'Test');
91
92
93 % choix 35
94
95 %%
96
97 erreurs = zeros(2, 2) ;
98
99 for i = 1:2
100
101     nbArbres = nbArbresVals(i);
102     [a,b] = randomForestEval(X, Y, Xtest, Ytest, 8, 0.85, 121, i);
103     erreurs(i,:) = [a b];
104
105 end
106
107 figure
108 plot(1:2, erreurs);
109 legend('Apprentissage', 'Test');
110 title('Erreur en fonction du nombre de dimension par arbre (Hmax = 8, purete = 0.85, 121 arbres)');
111
112 % 2 dimensions
113
114 %%
115
116 [erreurTrain, erreurTest] = randomForestEval(X, Y, Xtest, Ytest, 8, 0.85, 35, 2)
117
118 %% Jeu de données pima
119
120 data = csvread('pima-indians-diabetes.data');
121
122 X = data(:,1:end-1);
123 Y = data(:,end);
124 Y(Y==0) = -1;
125
126 % découpage en apprentissage et test
127 [X, Y, Xtest, Ytest] = splitdata(X, Y, 0.3);
128 p = size(X,2);
129
130 %%
131
132 hauteurVals = [2 3 4 5 7 9 12];
133 erreurs = zeros(length(hauteurVals), 2) ;
134
135 for i = 1:length(hauteurVals)
136
137     hauteur = hauteurVals(i);
138     [a,b] = randomForestEval(X, Y, Xtest, Ytest, hauteur, 0.7, 69, p, 10);
139     erreurs(i,:) = [a b];
140
141 end
142
143 figure
144 plot(hauteurVals, erreurs);
145 legend('Apprentissage', 'Test');
146 title('Erreur en fonction de la hauteur (purete = 70%, 69 arbres)');
147
148 % choix = 8
149
150 %%
151
152 pureteVals = [0.2 0.3 0.4 0.5 0.6 0.68 0.75 0.82 0.85 0.88 0.92 0.96];
153 erreurs = zeros(length(pureteVals), 2) ;
154
155 for i = 1:length(pureteVals)
```

```
156
157     purete = pureteVals(i);
158     [a,b] = randomForestEval(X, Y, Xtest, Ytest, 25, purete, 69, p, 5);
159     erreurs(i,:) = [a b];
160
161 end
162
163 figure
164 plot(pureteVals, erreurs);
165 legend('Apprentissage', 'Test');
166 title('Erreur en fonction de la purete (Hmax = 25, 69 arbres)');
167
168 %%
169
170 nbDimVals = [2 3 4 5 6 7 8];
171 erreurs = zeros(nbDimVals, 2) ;
172
173 for i = 1:length(nbDimVals)
174
175     nbDim = nbDimVals(i);
176     [a,b] = randomForestEval(X, Y, Xtest, Ytest, 9, 0.7, 201, nbDim, 10);
177     erreurs(i,:) = [a b];
178
179 end
180
181 figure
182 plot(nbDimVals, erreurs);
183 legend('Apprentissage', 'Test');
184 title('Erreur en fonction du nombre de dimension par arbre (Hmax = 9, purete = 0.7, 201 arbres)');
185
186 %%
187 [erreurTrain, erreurTest] = randomForestEval(X, Y, Xtest, Ytest, 9, 0.70, 69, 7)
```