

# DATA MINING 2

## TP1 & 2 - Toolbox réseaux de neurones

Thomas ROBERT

### 1 Codage et test des fonctions

Afin de vérifier le fonctionnement de la toolbox, j'ai testé les diverses fonctions avec deux problèmes très simple : un problème de régression de la fonction  $y = 2 \times x$ , et un problème de classification de dimension 1.

Le réseau testé est un réseau MLP avec une couche cachée, 1 entrée et 1 ou 2 sorties selon les cas.

En utilisant la toolbox, on se rend compte que le choix des fonctions d'activation à un impact fort sur les performances du modèle. Par ailleurs, le fait que la méthode utilise un pas fixe et un nombre d'itération fixé fait que le choix de ces paramètres est très important : un pas trop faible rend la convergence très lente, un pas trop grand fait diverger au lieu de converger.

Une bonne amélioration serait de fixer un critère de fin en plus d'un nombre d'itération maximum, et d'utiliser une méthode à pas variable pour accélérer la convergence et s'assurer de ne pas faire augmenter le coût, comme ça peut être le cas actuellement si le pas est choisi trop grand.

Notons également que pour augmenter légèrement la rapidité des calculs et surtout le confort d'utilisation, j'ai ajouté un paramètre à la fonction `onlinegrad` qui possède désormais une option `verbose` indiquant si on doit ou non afficher toutes les itérations.

### 2 Générer le jeu de test

On génère un jeu de données que l'on découpe en apprentissage (10%) et en validation (90%) (voir figure 1).

### 3 Création d'un réseau

On se propose de tester un réseau avec 2 fonctions d'activation `tanh`. Nous avons 2 entrées qui sont les 2 dimensions de chaque donnée  $x$ , et une sortie (la classe).

On teste donc ce réseau avec entre 1 et 10 neurones dans la couche cachée.

Notons qu'à partir des données initiales, il est important de construire le bon vecteur cible à passer à la toolbox. Dans le cas d'une sortie de `tanh`, la cible contient des 0 et des 1.

On notera que le choix d'un critère type MSE se justifie par le fait qu'il faut pouvoir dériver la fonction, mais la vraie mesure de qualité du réseau est le taux de bonne classification, que l'on obtient en arrondissant la sortie afin d'obtenir des 0 et des 1, et de comparer à la cible.

En effet, il faut distinguer le fait que la sortie soit réelle alors que l'on réalise en réalité de la classification et donc que l'ensemble cible est un ensemble de valeurs de cardinalité finie. Il n'est donc pas "grave" d'être un peu éloigné de la valeur cible tant que la valeur la plus proche est la bonne.

Les résultats sont visibles sur le graphe 2.

On constate qu'on obtient des résultats satisfaisants et environ identiques dès que l'on a au moins 3 neurones dans la couche cachée. Ceci est peut-être dû au fait que le problème n'est pas linéaire, ce qui augmente la complexité du problème et demande sans doute plus de neurones qu'un problème de séparation linéairement séparable. Un problème plus simple aurait sans doute nécessité moins de neurones pour avoir des résultats satisfaisants. Il est tout de même assez impressionnant de remarquer que 3 neurones dans la couche cachée suffisent à résoudre un problème de classification quadratique.

On constate tout de même que l'ajout de neurones ne semble pas améliorer la résolution du problème, comme on pourrait le penser naïvement.

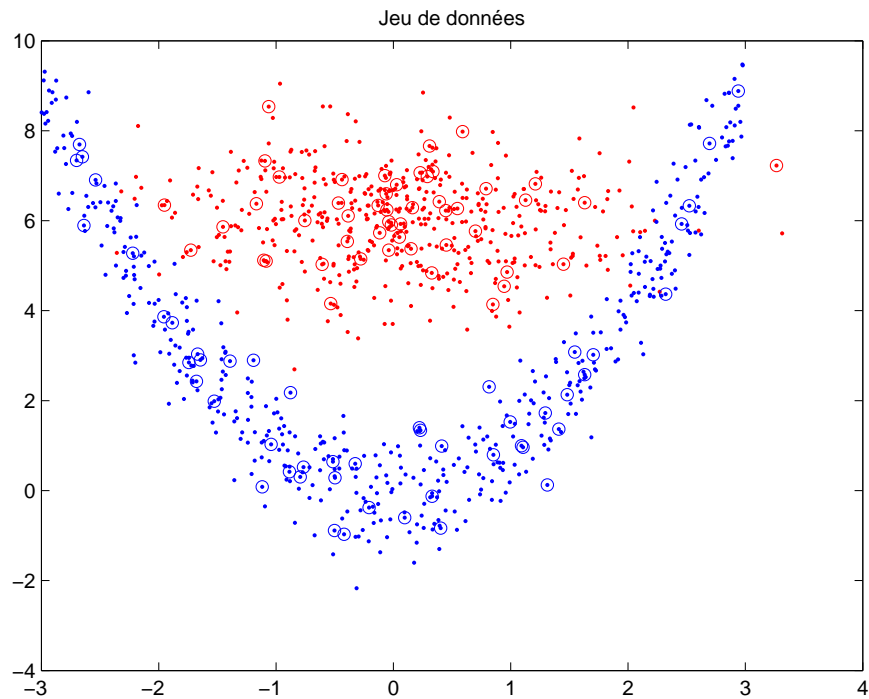


FIGURE 1 – Encerclés, les points du jeu d'apprentissage. Les autres sont le jeu de test.

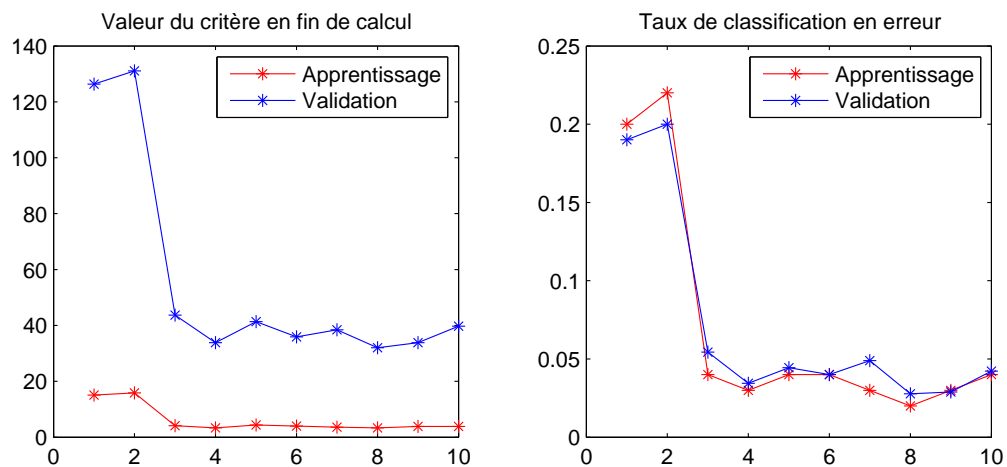


FIGURE 2 – Synthèse des résultats (en abscisse le nombre de neurones dans la couche cachée)

## 4 Code Matlab

### 4.1 Test des fonctions

```

1 % Exemple de problème de régression
2
3 clear all
4
5 net=ASINETfactory(1,[3 1],{'linear','linear'});
6
7 X = [1 2 3 4 5 6]';
8 Y = [2 4 6 8 10 12]';
9
10 [netout, learningErr, valError]=ASINETonlinegrad(
11     net,X,Y,0.01,100,'mse', false);
12 YE=ASINETforward(netout,X);
13 %YE =
14 %     2.0510
15 %     4.0419
16 %     6.0328
17 %     8.0237
18 %    10.0146
19 %    12.0055
20
21
22 % Exemple de problème de classification
23
24 clear all
25
26 net=ASINETfactory(1,[5 2],{'tanh','softmax'});
27
28 X = [1 2 3 7 8 9]';
29 Y = [0 0 0 1 1 1; 1 1 1 0 0 0]';
30
31 [netout, learningErr, valError]=ASINETonlinegrad(
32     net,X,Y,0.001,1000,'nll', false);
33 YE=ASINETforward(netout,X);
34 %YE =
35 %     0.0095     0.9905
36 %     0.0185     0.9815
37 %     0.0518     0.9482
38 %     0.9666     0.0334
39 %     0.9761     0.0239
40 %     0.9791     0.0209

```

### 4.2 Génération du jeu de test

```

1 clear all
2
3 nX=1000;
4 X = zeros(nX,2);
5
6 X(1:nX/2,:) = randn(nX/2,2) + repmat([0 6], nX/2, 1);
7 X(nX/2+1:nX,1) = (rand(nX/2,1) - 0.5) * 6;
8 X(nX/2+1:nX,2) = X(nX/2+1:nX,1).^2 + 0.7*randn(nX/2,1);
9
10 Y = [ones(nX/2,1) ; zeros(nX/2,1)];
11
12 plot(X(Y==1,1), X(Y==1,2), 'r'); hold on;
13 plot(X(Y==0,1), X(Y==0,2), 'b');
14
15 % découpage en apprentissage et test
16 [Xapp, Yapp, Xtest, Ytest] = splitdata(X, Y, 0.1);

```

### 4.3 Création du réseau

```

1 % valeurs à tester
2 n_vals = 1:10;
3
4 % Stockage pour affichage
5 errApp = zeros(length(n_vals),1);
6 errVal = zeros(length(n_vals),1);
7 nbErrApp = zeros(length(n_vals),1);
8 nbErrVal = zeros(length(n_vals),1);
9
10 % Pour chaque nb de neurones dans la couche
    cachée
11 for i = 1:length(n_vals)
12     n = n_vals(i);
13
14     net=ASINETfactory(2,[n 1],{'tanh','tanh'});
15     [netout, learningErr, valError]=
16         ASINETonlinegrad(net,Xapp,Yapp,0.01,250,
17             'mse', true, Xtest, Ytest);
18     YE=ASINETforward(netout,Xapp);
19     YE2=ASINETforward(netout,Xtest);
20     errApp(i) = learningErr(end);
21     errVal(i) = valError(end);
22     nbErrApp(i) = sum(round(YE) ~= Yapp);
23     nbErrVal(i) = sum(round(YE2) ~= Ytest);
24
25 end
26
27 figure;
28 subplot(1,2,1);
29 plot(n_vals, errApp, 'r'); hold on
30 plot(n_vals, errVal, 'b');
31 title('Valeur du critère en fin de calcul');
32 legend('Apprentissage', 'Validation');
33 subplot(1,2,2);
34 plot(n_vals, nbErrApp/length(Xapp), 'r'); hold
    on
35 plot(n_vals, nbErrVal/length(Xtest), 'b');
36 title('Taux de classification en erreur');
37 legend('Apprentissage', 'Validation');

```