

# Data Mining

## TP8 - Optimisation sans contraintes

Rémi MUSSARD - Thomas ROBERT

Le but du TP est de faire de la régression logistique afin de faire de classification par apprentissage supervisé.

## 1 Régression logistique binomiale

On commence par une régression logistique binomiale, c'est à dire avec deux classes.

### 1.1 Fonctions codées

On code pour cela plusieurs fonctions :

#### 1.1.1 EncoderLabel01

Cette fonction permet de passer d'un label -1 / 1 en un label 0 / 1.

```
1 function z = EncoderLabel01(y)
2     z = (y+1)/2;
3 end
```

#### 1.1.2 probaAPosteriori

Cette fonction permet de calculer les probabilités à posteriori. Elle fonctionne pour 2 classes ou plus.

```
1 function p = probaAPosteriori(theta , phi)
2
3     p = exp(phi*theta);
4     p = p./repmat((1+sum(p,2)), 1, size(theta,2));
5
6 end
```

#### 1.1.3 ma\_reg\_log

Cette fonction calcule la régression logistique.

```
1 function [theta,L,p] = ma_reg_log(phi , z)
2
3     theta = zeros(size(phi,2) , 1);
4     theta_old = theta + 5;
5     i=1;
6
7     L = [];
8
9     while norm(theta_old - theta) > 1e-1
10         theta_old = theta;
```

```
11     p = probaAPosteriori(theta , phi);
12     l = -sum(z.*log(p) + (1 - z).*log(1-p));
13     L = [L l];
14     W = diag(p.*(1-p));
15     r = phi*theta + W\ (z - p);
16     theta = (phi'*W*phi)\(phi'*W*r);
17     i = i + 1;
18 end
19 end
```

## 1.2 Application aux données clowns

Après avoir séparé les données en un ensemble de test et un ensemble d'apprentissage (80% en test / 20% en apprentissage), on appliquera les fonctions vues ci-dessus.

```
1 load clownsv7.mat
2
3 % Labels réencodés de -1 / 1 à 0 / 1
4 z = EncoderLabel01(y);
5
6 % Partage des données en deux : Apprentissage et Test
7 [xapp, zapp, xtest, ztest] = splitdata(X, z, 0.2);
8 n = size(X, 1);
9 napp = size(xapp, 1);
10 ntest = size(xtest, 1);
```

## 1.3 Frontière de décision linéaire

On calcul d'abord la régression sur une matrice  $\phi = [1 \ x_1 \ x_2]$  afin d'obtenir une frontière de décision linéaire.

On obtient alors une erreur d'environ 18% avec cette frontière linéaire. Cette erreur est raisonnable mais reste assez importante, d'autant qu'on pourrait sans doute beaucoup l'améliorer en utilisant une frontière de décision quadratique plus complexe et plus adaptée aux données.

Lorsque l'on cherche à classer les données de test à partir des paramètres de décisions obtenus avec les données d'apprentissage, on se rend compte que la différence d'erreurs est assez faible entre les données d'apprentissage et de tests.

```
1 % matrices de données
2 phiApp = [ones(napp,1) xapp];
3 phiTest = [ones(ntest,1) xtest];
4
5 % application de la regression logistique
6 [theta,L] = ma_reg_log(phiApp, zapp);
7
8 % calcul des classes
9 zappExp = round(probaAPosteriori(theta, phiApp));
10 ztestExp = round(probaAPosteriori(theta, phiTest));
11
12 % erreur
13 errApp = sum(zappExp ~= zapp)/napp;
14 errTest = sum(ztestExp ~= ztest)/ntest;
15
16 fprintf('Erreur en apprentissage : %i %%\n', round(errApp*100));
17 fprintf('Erreur en test : %i %%\n', round(errTest*100));
18
19 % calcul de la frontière de décision
20 xFront = [min(X(:,1)) max(X(:,1))]';
21 yFront = -(theta(1) + theta(2)*xFront)/theta(3);
22
23 % affichage
24 figure;
25 plot(xapp(zapp == 1,1), xapp(zapp == 1,2), 'r'); hold on;
26 plot(xapp(zapp == 0,1), xapp(zapp == 0,2), 'b');
```

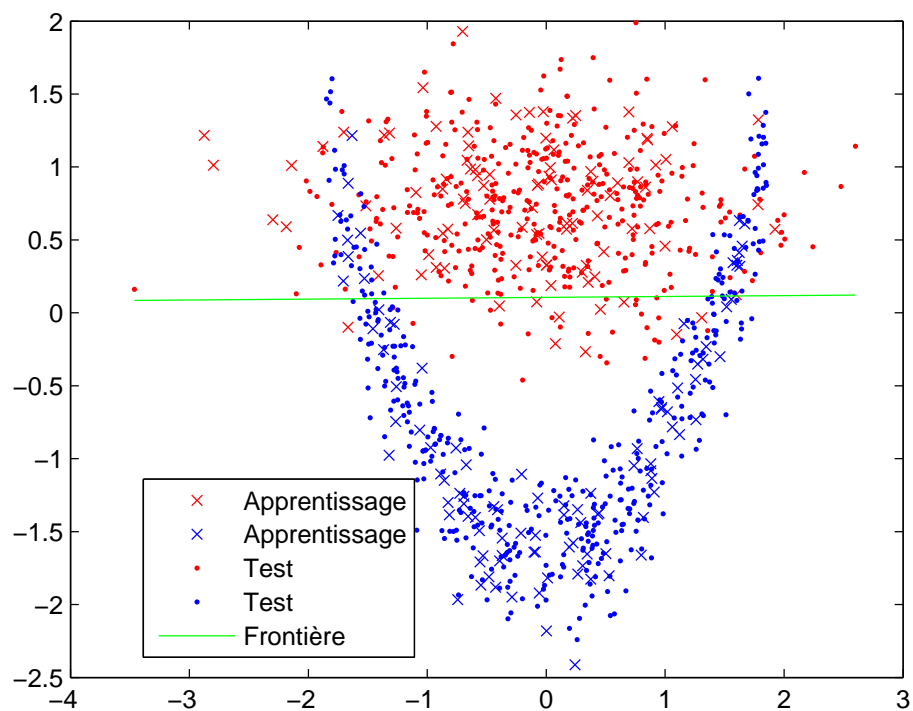
```

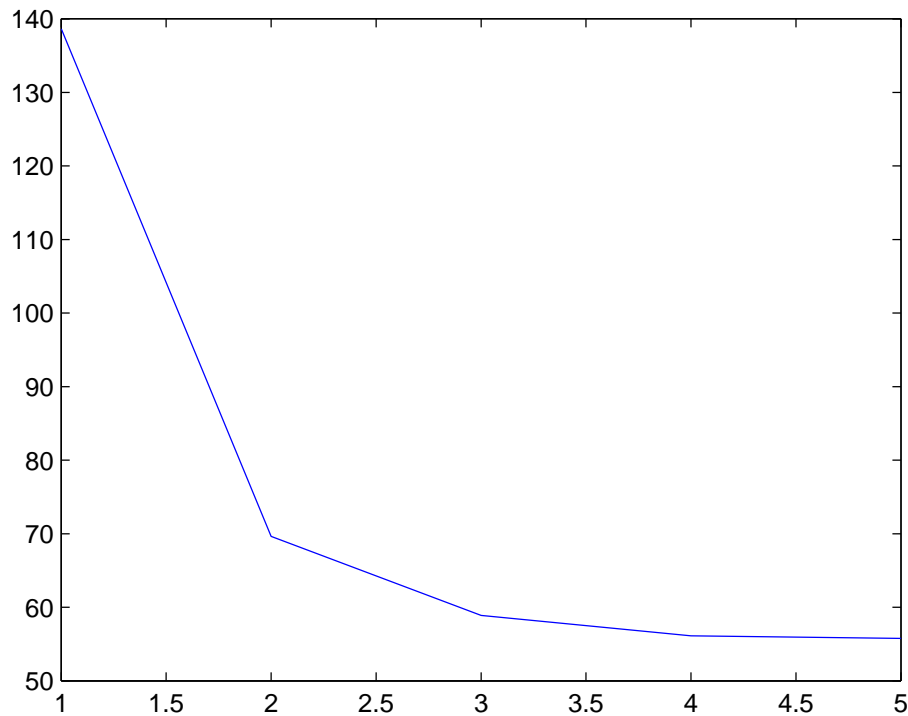
27 plot(xtest(ztest == 1,1),xtest(ztest == 1,2),'xr');
28 plot(xtest(ztest == 0,1),xtest(ztest == 0,2),'xb');
29 plot(xFront, yFront, '-g');
30
31 legend('Apprentissage', 'Apprentissage', 'Test', 'Test', 'Frontière', 'Location', 'Best');
32
33 % evolution critere de convergence
34 figure;
35 plot(L);

```

Erreur en apprentissage : 11 %

Erreur en test : 17 %





Evolution du maximum de vraisemblance

## 1.4 Frontière de décision quadratique

On réalise donc une frontière de décision quadratique afin d'améliorer la performance de l'algorithme.

On calcule pour cela la régression logistique avec la matrice  $\phi = [1 \ x_1 \ x_2 \ x_1 x_2 \ x_1^2 \ x_2^2]$ .

Le reste fonctionne de la même façon que pour une frontière linéaire, on a simplement changé d'espace. La seule différence est que la frontière est quadratique, donc à tracer avec la fonction `contour` de Matlab.

L'erreur tombe à 10%, soit presque la moitié de l'erreur avec une frontière linéaire. Cette erreur reste assez important car les données sont très mélangées autour de la frontière, il semble donc difficile de mieux séparer les données le plus excentrées des milieux de classes.

La meilleure solution pour diminuer le taux d'erreur serait de rejeter les points incertains trop proche de la frontière de décision. Cependant, l'inconvénient de cette méthode est qu'on classe moins de points.

```

1 % matrices de données
2 phiApp = [ones(napp,1) xapp xapp(:,1).*xapp(:,2) xapp.^2];
3 phiTest = [ones(ntest,1) xtest xtest(:,1).*xtest(:,2) xtest.^2];
4
5 % application de la regression logistique
6 [theta,L] = ma_reg_log(phiApp, zapp);
7
8 % calcul des classes
9 zappExp = round(probaAPosteriori(theta, phiApp));
10 ztestExp = round(probaAPosteriori(theta, phiTest));
11
12 % erreur
13 errApp = sum(zappExp ~= zapp)/napp;
14 errTest = sum(ztestExp ~= ztest)/ntest;
15
16 fprintf('Erreur en apprentissage : %i %%\n', round(errApp*100));
17 fprintf('Erreur en test : %i %%\n', round(errTest*100));
18
19 % calcul de la frontière de décision
20 [xx yy] = meshgrid(-4:0.1:4, -4:0.1:4);

```

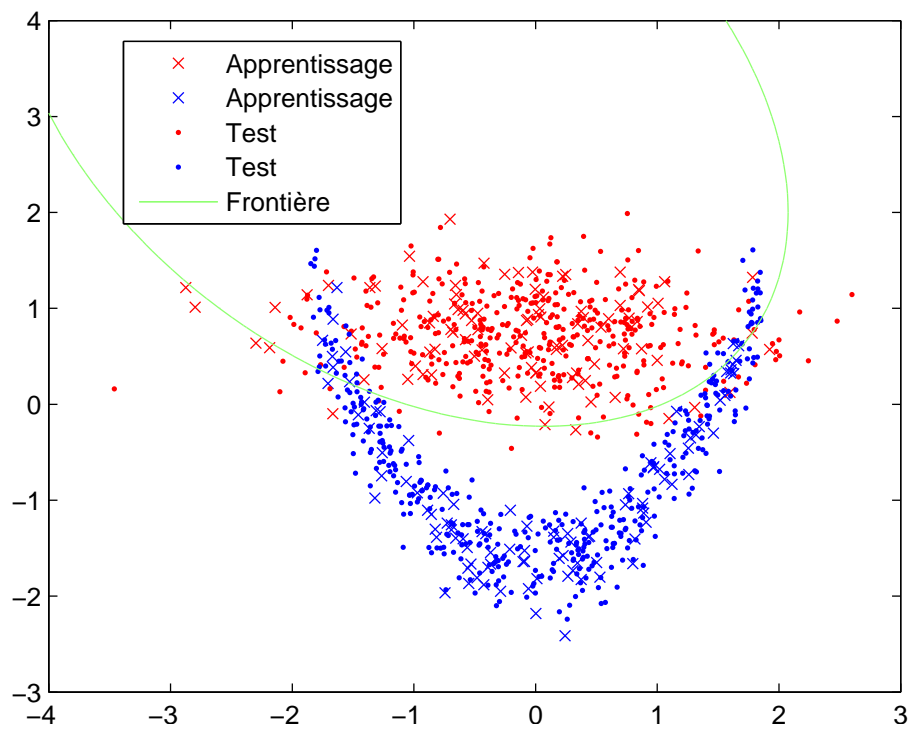
```

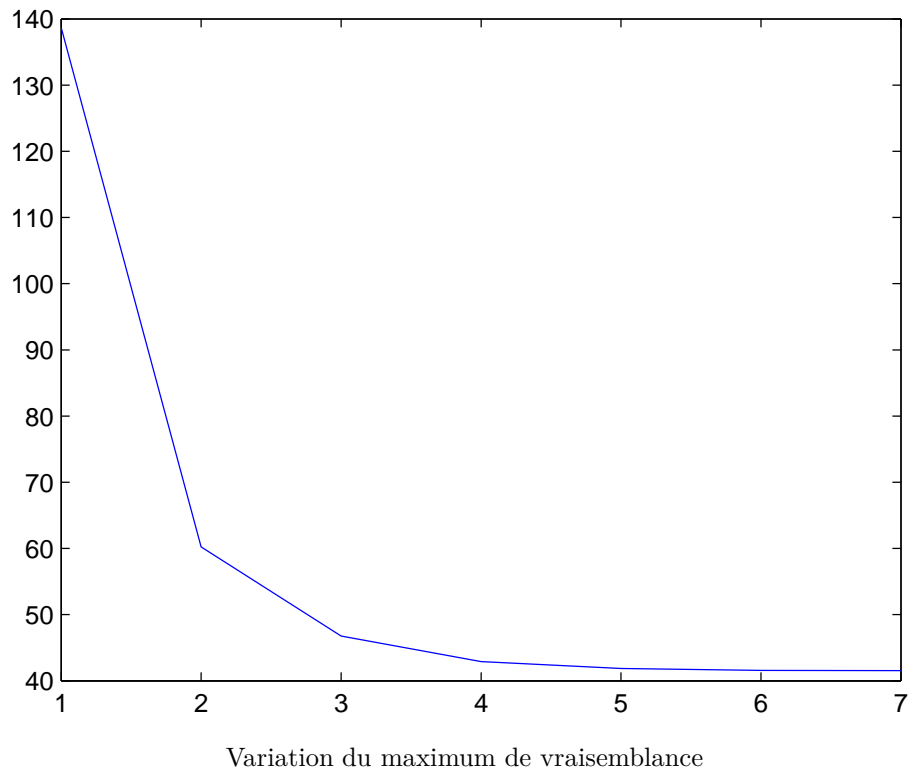
21 zz = theta(1) + theta(2)*xx + theta(3)*yy + theta(4)*(xx.*yy) + ...
22     theta(5) * (xx.^2) + theta(6) * (yy.^2);
23
24 % affichage
25 figure;
26 plot(xapp(zapp == 1,1),xapp(zapp == 1,2),'.r'); hold on;
27 plot(xapp(zapp == 0,1),xapp(zapp == 0,2),'.b');
28 plot(xtest(ztest == 1,1),xtest(ztest == 1,2),'xr');
29 plot(xtest(ztest == 0,1),xtest(ztest == 0,2),'xb');
30 contour(xx, yy, zz, [0 0]);
31
32 legend('Apprentissage', 'Apprentissage', 'Test', 'Test', 'Frontière', 'Location', 'Best');
33
34 % evolution critere de convergence
35 figure;
36 plot(L);

```

Erreur en apprentissage : 10 %

Erreur en test : 12 %





## 2 Régression logistique multimodale

Pour adapter le programme à 3 classes, il faudrait répéter l'opération de manière à séparer chaque groupe les uns des autres. Donc, pour l'ajout d'un  $(k+1)$ ème groupe, on aura  $k$  frontières de plus à celles déjà existantes.

Il faut donc une matrice d'appartenance aux clusters  $Z \in \{0, 1\}^{n \times k}$  qui se calcule grâce à la fonctions :

```

1 function [z, nbClasse] = EncoderLabel01_Multi(y)
2
3     vals = unique(y);
4     nbClasse = length(vals);
5
6     z = zeros(size(y, 1), nbClasse - 1);
7
8     for i = 1:nbClasse - 1
9         z(y == vals(i), i) = 1;
10    end
11
12 end

```

En notant  $\Theta = [\theta_1 \ \theta_2 \ \dots \ \theta_k]$  la matrice avec chaque paramètre de frontière en colonne, on peut simplement calculer les probabilités à posteriori. On aura dans la matrice de probabilités  $P \in \mathbb{R}^{n \times k}$  contenant la probabilité pour chaque frontière de décision  $P_i$  dans chaque colonne.

$$P_i = \frac{e^{\theta_i \phi}}{1 + \sum_{j=1}^k e^{\theta_j \phi}}$$

La fonction de calcul de probabilité a été généralisée et est la même que celle disponible dans la partie 1.

L'algorithme itératif de calcul de la régression logistique multimodale est plus compliquée. Nous n'avons pas réussi à faire fonctionner cet algorithme correctement.

La méthode que nous avons appliqué était de calculer  $W$ ,  $r$  et  $\theta_i$  pour chaque colonne de la matrice  $\Theta$ . Malheureusement, cette méthode ne semble pas fonctionner correctement.

Nous avons finalement le code suivant :

```
1 function [theta,L,p] = ma_reg_log_multi(phi, z)
2
3     nbFront = size(z, 2);
4
5     theta = zeros(size(phi,2), nbFront)+1;
6     theta_old = theta + 5;
7     i=1;
8
9     L = [];
10
11     %while norm(theta_old - theta) > 1e-1
12     for i = 1:100
13         theta_old = theta;
14         p = probaAPosteriori(theta, phi);
15         %l = -sum(z.*log(p) + (1 - z).*log(1-p));
16         %L = [L 1];
17         for j = 1:nbFront
18             W = diag(p(:,j) .* (1 - p(:,j)));
19             r = phi*theta(:,j) + W\ (z(:,j) - p(:,j));
20             theta(:,j) = (phi'*W*phi)\(phi'*W*r);
21         end
22         %i = i + 1;
23     end
24 end
```