

Data Mining

TP6 - Théorie bayésienne de la décision

Rémi MUSSARD - Thomas ROBERT

1 Question 1

Préparation des données

```
1 % taille des échantillons
2 n1 = 50;
3 n2 = n1;
4 n = n1 + n2;
5
6 % paramètres lois
7 mu1 = [0 0];
8 mu2 = [2 2];
9 mu = [mu1; mu2];
10 S = [1 0.5 ; 0.5 4];
11
12 % probas à priori
13 p = [n1/n, n2/n];
14
15 % échantillons aléatoires
16 X1 = randn(n1,2)*S^(1/2) + repmat(mu1, n1, 1);
17 X2 = randn(n2,2)*S^(1/2) + repmat(mu2, n2, 1);
18
19 % données
20 X = [X1;X2];
```

1.a. Frontière de décision théorique

On calcule et on trace la frontière de décision théorique d'après les moyennes et variances utilisées pour générer les points. On trace cette frontière.

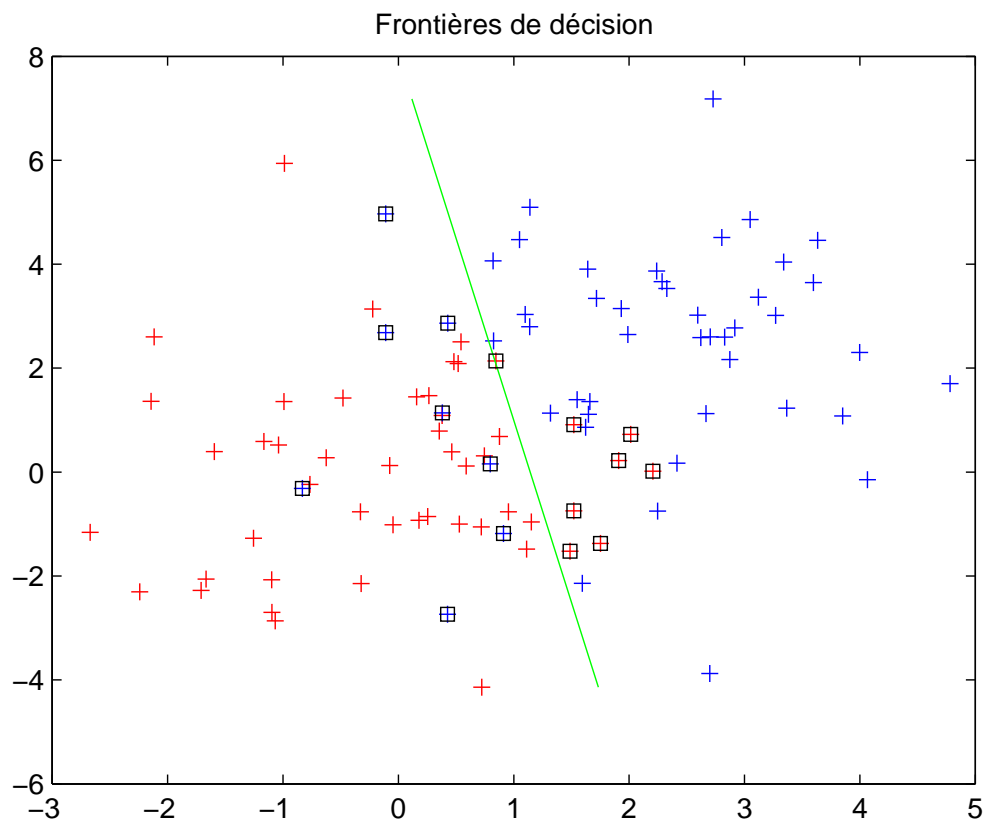
```
1 % affichage des points
2 close all;
3 plot(X1(:,1), X1(:,2), 'r');
4 hold on
5 plot(X2(:,1), X2(:,2), 'b');
6 title('Frontières de décision');
7
8 % frontière de décision d'après la formule du cours (la 2ème partie de la
9 % formule disparaît car les probas à priori sont identiques)
10 w = S\'(mu1'-mu2');
11 x0 = (mu1+mu2)/2;
12
13 y = max(X(:,2));
14 x = x0(1) + w(2)*(x0(2) - y)/w(1);
15 y2 = min(X(:,2));
16 x2 = x0(1) + w(2)*(x0(2) - y2)/w(1);
17
18 plot([x;x2],[y;y2], 'g');
```

```

19
20 % nombre d'erreurs
21 inds1 = find(w'*(X1-repmat(x0,n1,1))'<0);
22 inds2 = find(w'*(X2-repmat(x0,n1,1))'>0);
23 erreursX1 = length(inds1);
24 erreursX2 = length(inds2);
25 erreursTheo = erreursX1 + erreursX2
26
27 plot([X1(inds1,1);X2(inds2,1)],[X1(inds1,2);X2(inds2,2)], 'sk');

```

```
erreursTheo = 16
```



1.b. estimation des moyennes et matrice de covariance

On estime les paramètres selon le maximum de vraisemblance et on utilise les paramètres estimés pour calculer une frontière de décision estimée.

Sur les données d'apprentissage, on fait moins d'erreur en moyenne avec la frontière de décision estimée qu'avec la frontière théorique. Cependant, sur des données de tests, il est probable que ce soit l'inverse puisque la frontière théorique est la frontière idéale pour un jeu de données de taille infinie.

```

1 % estimation
2 mulhat = mean(X1)
3 mu2hat = mean(X2)
4 Shat = (cov(X1) + cov(X2)) / 2
5
6 % frontière de décision estimée
7 w = Shat \ (mulhat'-mu2hat');
8 x0 = (mulhat+mu2hat)/2;

```

```

9
10 y = max(X(:,2));
11 x = x0(1) + w(2)*(x0(2) - y)/w(1);
12 y2 = min(X(:,2));
13 x2 = x0(1) + w(2)*(x0(2) - y2)/w(1);
14
15 plot([x;x2],[y;y2], '—c');
16
17 % nombre d'erreurs
18 inds1 = find(w'*(X1-repmat(x0,n1,1))'<0);
19 inds2 = find(w'*(X2-repmat(x0,n1,1))'>0);
20 erreursX1 = length(inds1);
21 erreursX2 = length(inds2);
22 erreursEstim = erreursX1 + erreursX2
23
24 plot([X1(inds1,1);X2(inds2,1)],[X1(inds1,2);X2(inds2,2)], 'ok');
25
26 % légende
27 legend('X1','X2', ...
28        'FdD théorique', [num2str(erreursTheo) ' erreurs théo'],...
29        'FdD estimée', [num2str(erreursEstim) ' erreurs estim']);

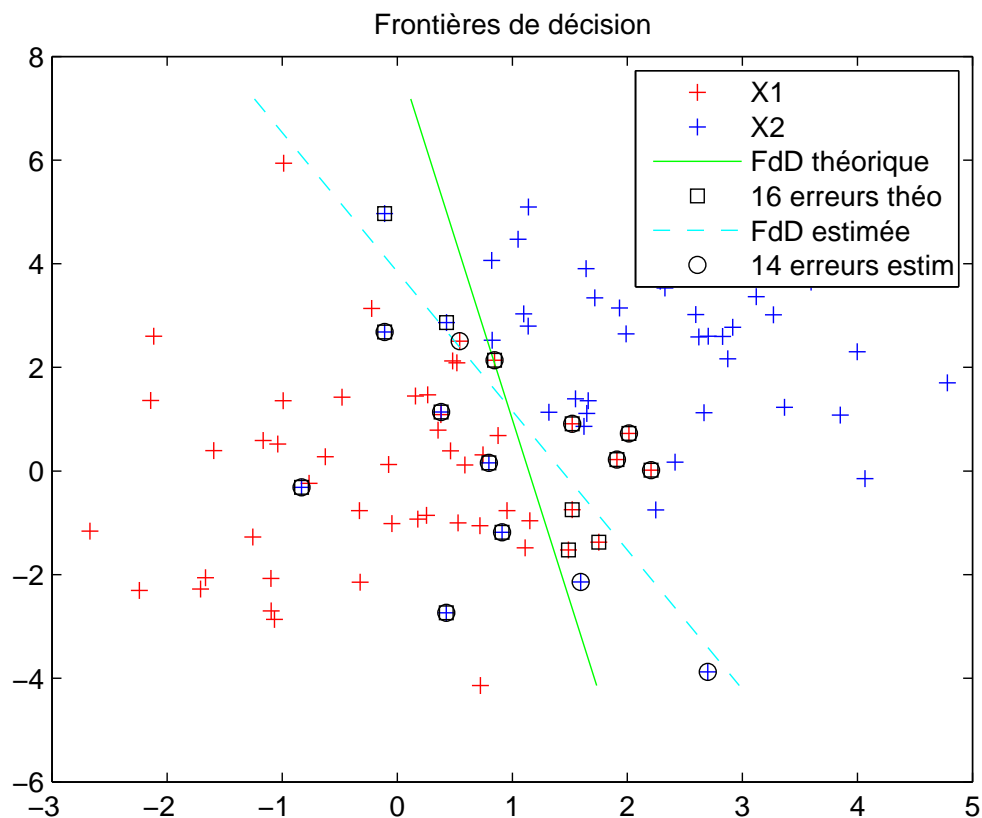
```

```
mu1hat = -0.0293 -0.0191
```

```
mu2hat = 2.0732 2.2202
```

```
Shat = 1.4714 0.2050
      0.2050 3.8767
```

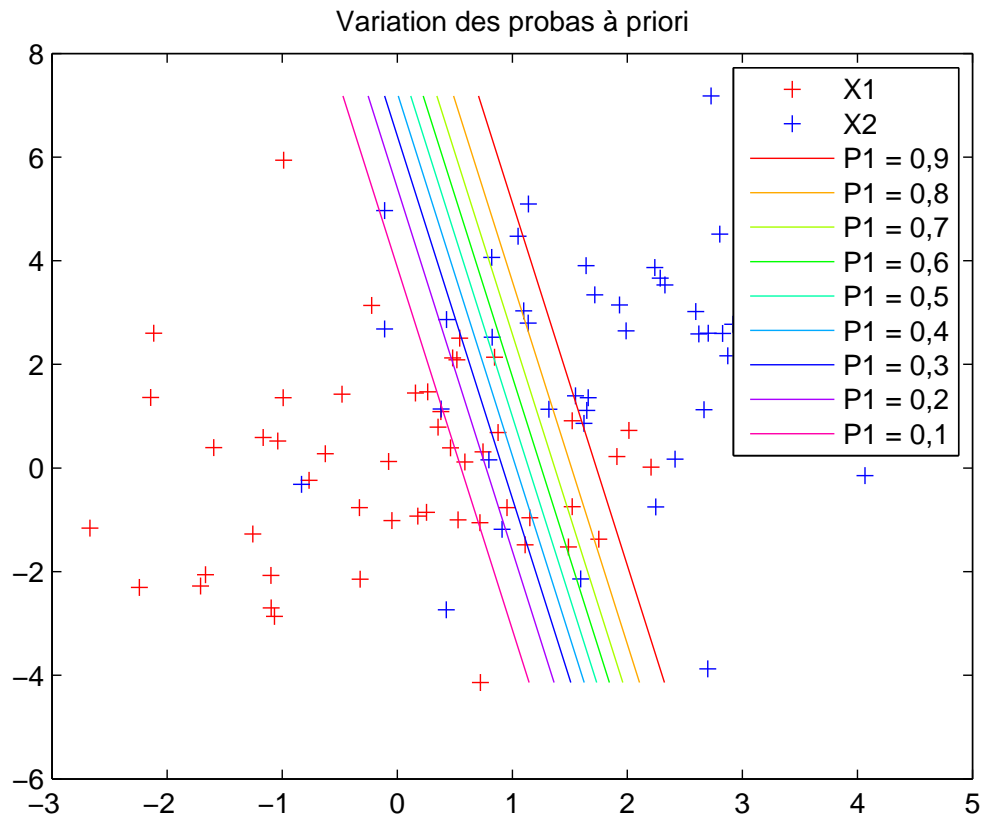
```
erreursEstim = 14
```



1.c. Variation des probabilités à priori

On décide de faire varier les probabilités à priori afin de voir l'impact que cela à.

```
1 figure;
2 plot(X1(:,1), X1(:,2), 'r');
3 hold on
4 plot(X2(:,1), X2(:,2), 'b');
5 title('Variation des probas à priori');
6
7 w=S\((mu1'-mu2')');
8
9 P1list = 0.1:0.1:0.9;
10 colors = hsv(9);
11
12 for i=1:length(P1list)
13     P1 = P1list(i);
14     P2 = 1 - P1;
15
16     x0 = 1/2*(mu1'+mu2')'- repmat(log(P1/P2)/((mu1'-mu2')'*w),1,2);
17     y = max(X(:,2));
18     x = x0(1) + w(2)*(x0(2) - y)/w(1);
19     y2 = min(X(:,2));
20     x2 = x0(1) + w(2)*(x0(2) - y2)/w(1);
21
22     plot([x;x2],[y;y2], '—', 'Color', colors(i,:));
23 end
24
25 legend('X1', 'X2', ...
26     'P1 = 0,9', ...
27     'P1 = 0,8', ...
28     'P1 = 0,7', ...
29     'P1 = 0,6', ...
30     'P1 = 0,5', ...
31     'P1 = 0,4', ...
32     'P1 = 0,3', ...
33     'P1 = 0,2', ...
34     'P1 = 0,1');
```



On remarque que (logiquement) cela déplace la frontière de décision vers le cluster ayant la plus grande probabilité à priori.

1.d. Rejet si ambiguïté

On a remarqué que les clusters étaient très proches l'un de l'autre. On décide donc de rejeter les points trop proche de la frontière de décision afin de réduire le nombre d'erreurs de classification.

Après 500 itérations pour chaque valeur de α , on calcule le nombre moyen de points rejetés, le nombre de rejets inutile (nombre de points rejetés alors qu'ils auraient bien été catégorisés), et les taux de classifications erronés avec et sans rejet.

On obtient les résultats suivants :

α	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.40	0.45	0.50
Nombre de rejets moyen	78	73	70	66	63	61	58	53	54	0
Rejets inutiles	65	49	39	31	25	19	14	9	5	0
Taux d'erreur avant rejet	15%	15%	15%	15%	15%	15%	15%	15%	15%	15%
Taux d'erreur après rejet	2%	4%	5%	6%	8%	9%	11%	12%	13%	15%

Le choix d'un α dépendra est donc un compromis entre le fait de rejeter beaucoup de point et le fait de conserver des erreurs. Un α de 0,20 nous semble un bon compromis.

Note : le code en italique permet de générer les résultats de tests des différentes valeurs de α visibles ci-dessus.

```

1 % seuil de rejet
2 alpha = 0.2
3
4 % statistiques
5 % stats = [];
6 % for alpha = 0.05:0.05:0.50
7 % rejetsInutile = [];

```

```

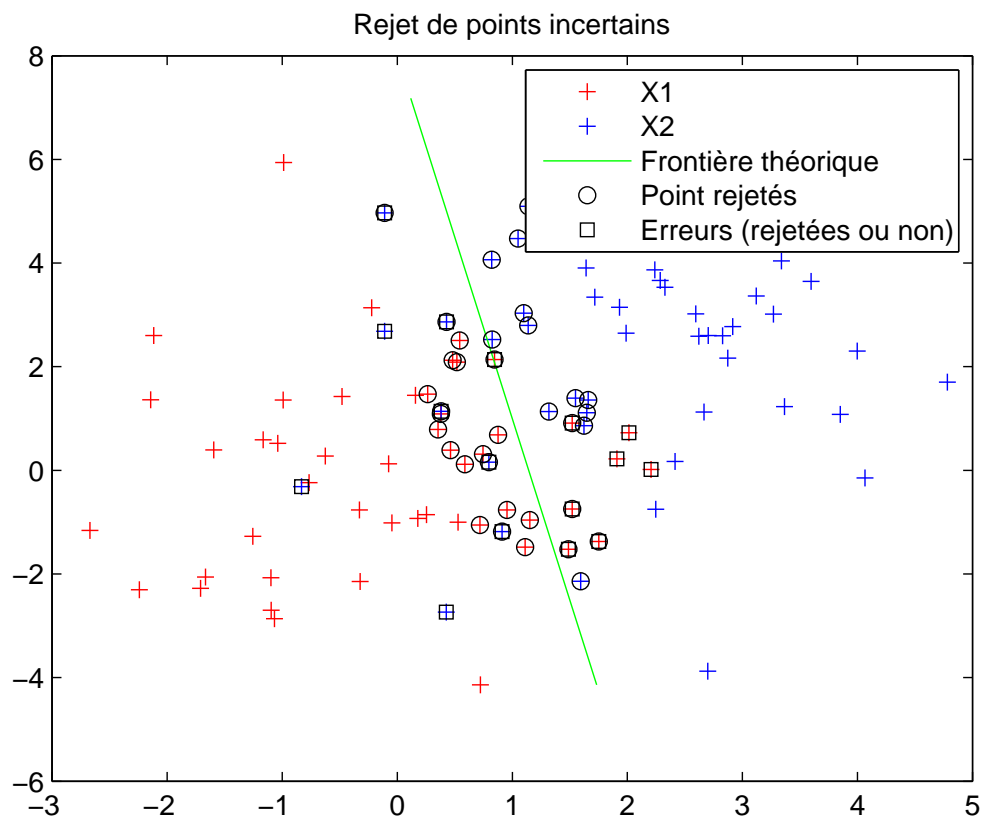
8 % nbsErreursAvantRejet = [];
9 % nbsErreursApresRejet = [];
10 % nbsRejets = [];
11 % for i = 1:500
12 % X1 = randn(n1,2)*S^(1/2) + repmat(mu1, n1, 1);
13 % X2 = randn(n2,2)*S^(1/2) + repmat(mu2, n2, 1);
14 % X = [X1;X2];
15
16 % probas à priori
17 P_1 = n1/n;
18 P_2 = n2/n;
19
20 % probabilités conditionnelles
21 P_x_1 = mvnpdf(X, mu1, S);
22 P_x_2 = mvnpdf(X, mu2, S);
23
24 % loi marginale de X
25 P_x = P_1 * P_x_1 + P_2 * P_x_2;
26
27 % probabilités à posteriori
28 P_1_x = (P_x_1 * P_1) ./ P_x;
29 P_2_x = 1 - P_1_x;
30
31 % affichage
32 figure;
33 plot(X1(:,1), X1(:,2), 'r');
34 hold on
35 plot(X2(:,1), X2(:,2), 'b');
36 title('Rejet de points incertains');
37
38 % frontière de décision
39 w = S \ (mu1' - mu2');
40 x0 = (mu1 + mu2) / 2;
41
42 y = max(X(:,2));
43 x = x0(1) + w(2) * (x0(2) - y) / w(1);
44 y2 = min(X(:,2));
45 x2 = x0(1) + w(2) * (x0(2) - y2) / w(1);
46
47 plot([x;x2],[y;y2], 'g');
48
49 % indices des points à rejeter
50 inds = ((P_1_x - P_2_x < 0) | (P_1_x < 1 - alpha)) & ...
51         ((P_2_x - P_1_x < 0) | (P_2_x < 1 - alpha));
52
53 % affichage des points à rejeter
54 plot(X(inds,1), X(inds,2), 'ok');
55
56 % nombre d'erreurs
57 inds1 = (w' * (X1 - repmat(x0, n1, 1)))' < 0;
58 inds2 = (w' * (X2 - repmat(x0, n1, 1)))' > 0;
59 inds1avecRejet = inds1' & ~inds(1:length(X1));
60 inds2avecRejet = inds2' & ~inds(length(X1)+1:end);
61
62 % affichage
63 plot([X1(inds1,1);X2(inds2,1)], [X1(inds1,2);X2(inds2,2)], 'sk');
64
65 legend('X1', 'X2', 'Frontière théorique', 'Point rejetés', 'Erreurs (rejetées ou non)');
66
67 % statistiques
68 nbRejets = sum(inds);
69 nbErreursAvantRejet = sum(inds1) + sum(inds2);
70 erreursX1 = sum(inds1avecRejet);
71 erreursX2 = sum(inds2avecRejet);
72 nbErreursApresRejet = erreursX1 + erreursX2;
73 nbRejetsUtiles = nbErreursAvantRejet - nbErreursApresRejet;
74 nbRejetsInutiles = nbRejets - nbRejetsUtiles;
75 if (nbRejets > 0)
76     rejetInutile = nbRejetsInutiles / nbRejets * 100;
77 else
78     rejetInutile = 0;

```

```

79 end
80
81 % statistiques
82 % nbsErreursAvantRejet = [nbsErreursAvantRejet nbErreursAvantRejet];
83 % nbsErreursApresRejet = [nbsErreursApresRejet nbErreursApresRejet];
84 % nbsRejets = [nbsRejets nbRejets];
85 % rejetsInutile = [rejetsInutile rejetInutile];
86 % end
87 % stats = [stats [alpha ;
88 %     mean(rejetsInutile);
89 %     mean(nbsRejets);
90 %     mean(nbsErreursAvantRejet)/n*100;
91 %     mean(nbsErreursApresRejet)/(n-mean(nbsRejets))*100;
92 %     ]];
93 % end
94 %

```



2 Question 2

2.a Chargement des données

Partage des données en données d'apprentissage et données de test grâce à la fonction `splitdata` disponible sur Moodle.

```

1 load('clownsv7.mat');
2
3 [xApp, yApp, xTest, yTest] = splitdata(X, y, 0.5);

```

2.b. et 2.c. Frontières et erreurs en LDA

On calcule la frontière de décisions pour la LDA et on calcule le nombre d'erreurs de classification.

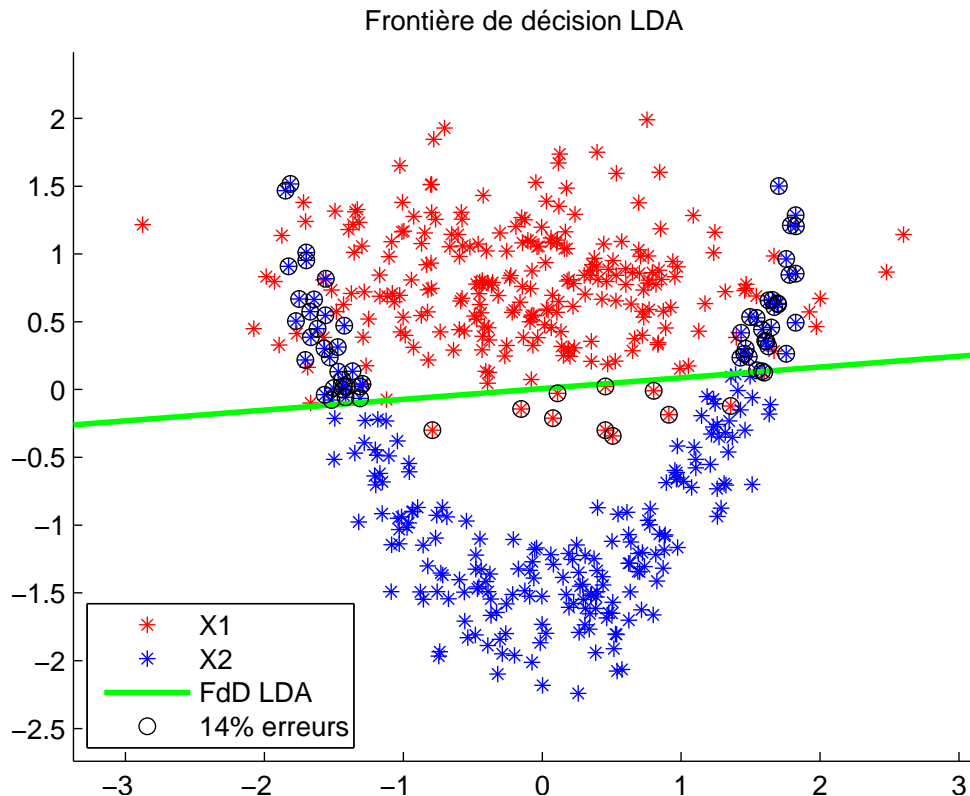
```

1 x1 = xApp(yApp > 0, :);
2 x2 = xApp(yApp < 0, :);
3 x12 = [x1; x2];
4
5 mu1 = mean(x1);
6 mu2 = mean(x2);
7
8 S = (cov(x1) + cov(x2)) / 2;
9
10 % calcul de la frontière de décision en LDA
11
12 w = S \ (mu1' - mu2');
13 x0 = (mu1 + mu2) / 2;
14
15 Y1 = max(x12(:, 2));
16 X1 = x0(1) + w(2) * (x0(2) - Y1) / w(1);
17 Y2 = min(x12(:, 2));
18 X2 = x0(1) + w(2) * (x0(2) - Y2) / w(1);
19
20 % Erreurs LDA
21 [n1, m1] = size(x1);
22 [n2, m2] = size(x2);
23
24 inds1 = find(w' * (x1 - repmat(x0, n1, 1))' < 0);
25 inds2 = find(w' * (x2 - repmat(x0, n2, 1))' > 0);
26 erreursLDA = (length(inds1) + length(inds2));
27 ratioErreursLDA = erreursLDA / length(x12);
28
29 % affichage
30 figure; hold on;
31 plot(x1(:, 1), x1(:, 2), 'r');
32 plot(x2(:, 1), x2(:, 2), 'b');
33 plot([X1; X2], [Y1; Y2], 'g', 'linewidth', 2);
34 plot([x1(inds1, 1); x2(inds2, 1)], [x1(inds1, 2); x2(inds2, 2)], 'ok');
35 axis([min(x12(:, 1)) - 0.5, max(x12(:, 1)) + 0.5, min(x12(:, 2)) - 0.5, max(x12(:, 2)) + 0.5]);
36 title('Frontière de décision LDA');
37 leg = legend('X1', 'X2', 'FdD LDA', ...
38             [int2str(round(ratioErreursLDA * 100)) '% erreurs']);
39 set(leg, 'Location', 'SouthWest')

```

erreursLDA = 70

ratioErreursLDA = 14%



A la vue des données, tracer une frontière de décision via une LDA ne semble pas convenir. En effet, la frontière entre les données ne semble pas linéaire. Il serait plus judicieux de faire une frontière de décision via une QDA, étant donné le caractère hyperbolique apparent de la frontière.

Lorsque l'on calcule l'erreur que nous avons avec la LDA, on obtient un environ 14% d'erreurs.

2.b. et 2.c. Frontières et erreurs en QDA

On calcule maintenant la frontière de décisions pour la QDA et on calcule le nombre d'erreurs de classification pour cette nouvelle méthode.

```

1 % calcul de la QDA
2 Wj = inv(S)/2;
3 wj = 2*Wj*mul';
4 wj0 = (mul*wj)/2 - (log(norm(2*Wj)))/2 + log(length(x1)/length(xApp));
5
6 xQDA = sort(x12(:,1));
7 yQDA = wj0 + wj(1)*xQDA + Wj(1,1)*xQDA.^2;
8
9 % Erreurs QDA
10 x = x1(:,1);
11 yQDA1 = wj0+wj(1)*x+Wj(1,1)*x.^2;
12 x = x2(:,1);
13 yQDA2 = wj0+wj(1)*x+Wj(1,1)*x.^2;
14
15 inds1 = find((x1(:,2)-yQDA1)<0);
16 inds2 = find((x2(:,2)-yQDA2)>0);
17 ratioErreursQDA = (length(inds1) + length(inds2))/length(x12)
18
19 % affichage
20 figure;
21 plot(x1(:,1), x1(:,2), 'r');
22 hold on
23 plot(x2(:,1), x2(:,2), 'b');

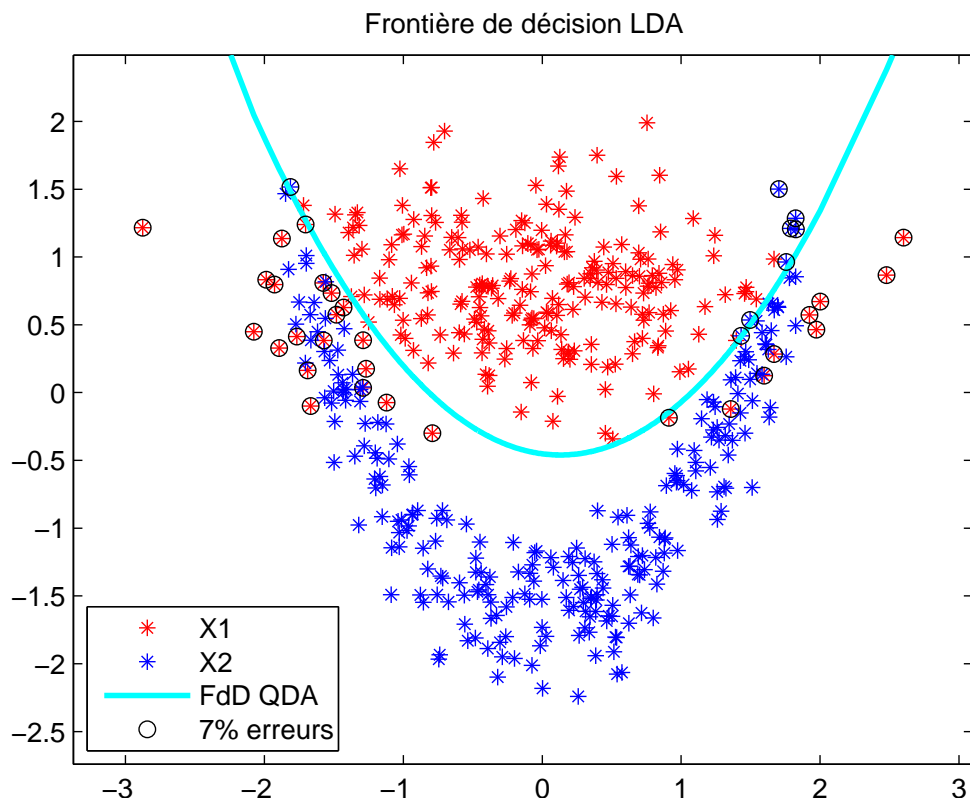
```

```

24 plot(xQDA,yQDA,'-c', 'linewidth', 2);
25 plot([x1(inds1,1) ; x2(inds2,1)], [x1(inds1,2) ; x2(inds2,2)], 'ok');
26 axis([min(x12(:,1))-0.5 max(x12(:,1))+0.5 min(x12(:,2))-0.5 max(x12(:,2))+0.5]);
27 title('Frontière de décision LDA');
28 leg = legend('X1', 'X2', 'FdD QDA', ...
29 [int2str(round(ratioErreursQDA*100)) '% erreurs']);
30 set(leg, 'Location', 'SouthWest')

```

ratioErreursQDA = 7,4 %



La frontière de décision réalisée avec la QDA semble plus appropriée. On obtient une frontière de décision de forme hyperbolique qui sépare relativement bien les données en deux classes.

En calculant l'erreur de classification sur QDA, on obtient un taux d'erreur d'environ 8%, ce qui est moins élevé que le taux d'erreur de la LDA.

2.d. Frontière de décision pour LDA dans l'espace $\{x_1, x_2, x_1^2, x_2^2, x_1 \cdot x_2\}$

On trace cette fois la frontière LDA dans l'espace $\{x_1, x_2, x_1^2, x_2^2, x_1 \cdot x_2\}$.

```

1 x1carre = x12(:,1).*x12(:,1);
2 x2carre = x12(:,2).*x12(:,2);
3 x1x2 = x12(:,1).*x12(:,2);
4
5 xfinal = [x12 x1carre x2carre x1x2];
6 xnew1 = xfinal(yApp > 0, :);
7 xnew2 = xfinal(yApp < 0, :);
8

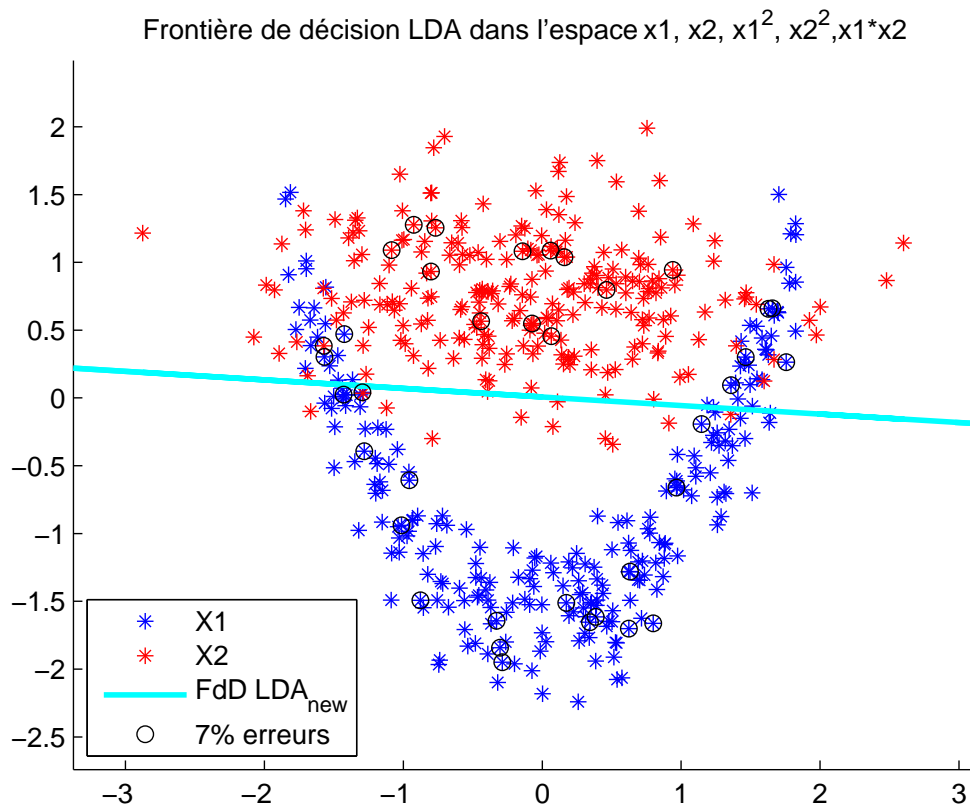
```

```

9 % estimation paramètres
10 mu1 = mean(xnew1);
11 mu2 = mean(xnew2);
12 S = (cov(xnew1) + cov(xnew2)) / 2;
13
14 % calcul FdD
15 w = S \ (mu1' - mu2');
16 x0 = (mu1 + mu2) / 2;
17
18 X1 = min(xfinal);
19 X2 = max(xfinal);
20 Y1 = x0(2) - (w(1)*(x0(1)-X1))/w(2);
21 Y2 = x0(2) - (w(1)*(x0(1)-X2))/w(2);
22
23 % erreurs
24 inds1 = find(w'*(xnew1 - repmat(x0, size(xnew1,1), 1))' < 0);
25 inds2 = find(w'*(xnew2 - repmat(x0, size(xnew2,1), 1))' > 0);
26 ratioErreursLDAnew = (length(inds1) + length(inds2))/length(xfinal)
27
28 % affichage
29 figure; hold on;
30 plot(xnew1(:,1), xnew1(:,2), '*b');
31 plot(xnew2(:,1), xnew2(:,2), '*r');
32 plot([X1 X2], [Y1 Y2], '-c', 'linewidth', 2);
33 plot([x1(inds1,1) ; x2(inds2,1)], [x1(inds1,2) ; x2(inds2,2)], 'ok');
34 axis([min(x12(:,1))-0.5 max(x12(:,1))+0.5 min(x12(:,2))-0.5 max(x12(:,2))+0.5]);
35 title('Frontière de décision LDA dans l'espace {x1, x2, x1^2, x2^2, x1*x2}');
36 leg = legend('X1', 'X2', 'FdD LDA_{new}', ...
37 [int2str(round(ratioErreursLDAnew*100)) '% erreurs']);
38 set(leg, 'Location', 'SouthWest')

```

ratioErreursLDAnew = 7.4 %



Il est donc impossible de visualiser cette frontière dans cet espace de \mathbb{R}^5 , on se contente donc de le visualiser sur les deux premières composantes de cet espace. Cependant, il est important de noter qu'il ne s'agit que donc que d'une partie des composantes des points et que la frontière est en fait un hyperplan. On ne peut donc pas visualiser quels points sont de tel ou tel côté de la frontière directement. Ces points en erreur ont tout de même été marqués.

Le taux d'erreur est ici d'environ 8%.

2.e. Comparaison des résultats

Appliquer la LDA dans un espace de plus grande dimension permet d'améliorer la précision de la LDA. On se retrouve ici avec un taux d'erreur très proche de celui de la QDA réalisée précédemment.

Cependant, le fait de devoir augmenter le nombre de dimensions de la LDA la rend sans doute plus gourmande en ressource qu'une QDA.