

# DATA MINING 2

## TP 2.1 SVM & Kernel

Thomas ROBERT

L'objectif du TP est de tester la méthode du SVM avec kernel qui permet de résoudre un problème non linéaire avec un SVM linéaire en changeant les données d'espace.

Durant le TP, le même code (fourni) a été testé et modifié afin de tester différents paramètres de la méthode. Les résultats de ces modifications seront présentés, le code est présenté en annexe.

### 1 Temps de calcul

La première version du code avait pour but de comparer les temps de calcul.

On voit bien que les méthodes de calcul optimisées sont plus rapide que les méthodes de calcul généralistes / manuelles.

Par exemple, le calcul du kernel manuel (déjà légèrement optimisé) prend 0,18s alors que la méthode `svkernel` obtient le même résultat en 0,018s, c'est à dire 10 fois plus rapidement.

Par ailleurs, la résolution du problème de minimisation par `cvx` prend 20.5s contre 0.075s pour `monqp` soit presque 300 fois plus vite. On voit ici que `cvx` est une toolbox très pratique pour la souplesse qu'elle offre dans le formulation du problème de minimisation, mais qu'elle est (et c'est logique) beaucoup moins rapide qu'une méthode de résolution optimisée.

Pour avoir des temps de calcul raisonnable quand on augmente le nombre de données, il faut donc retravailler les problèmes optimisation pour revenir à une forme standard permettant d'utiliser des solveurs optimisés.

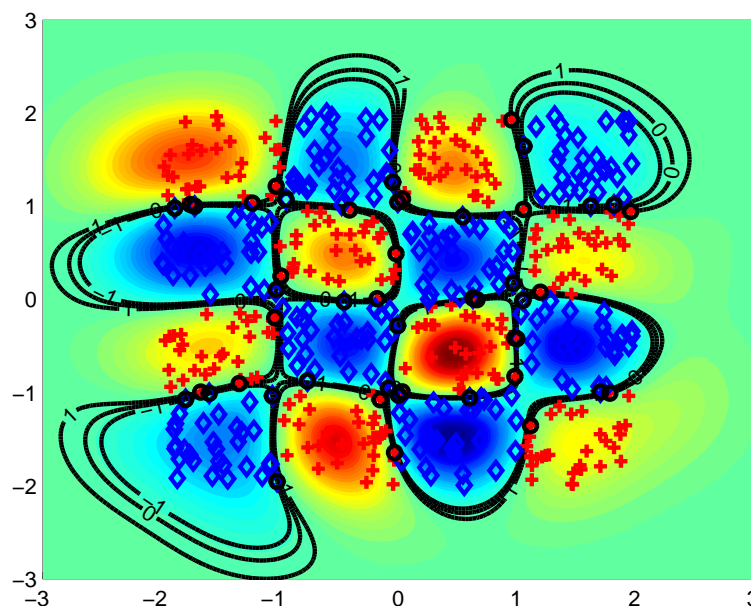


FIGURE 1 – Résultat du calcul avec 500 données

## 2 Test avec 10 000 données

On teste cette fois le même code (en ne gardant que les méthodes optimisées) avec 10000 données. On voit que le SVM arrive à bien résoudre le problème de classification et que les frontières tracent bien le quadrillage prévu.

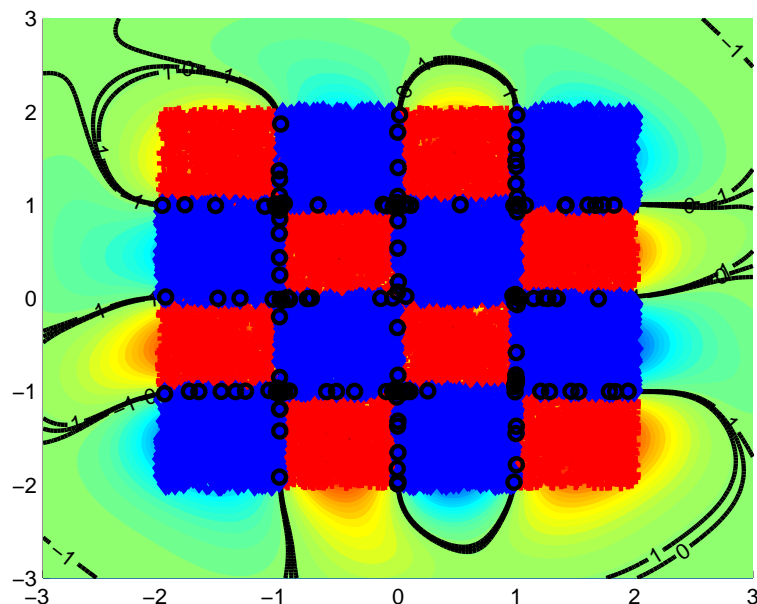


FIGURE 2 – Résultat du calcul avec 10 000 données

## 3 Modification de la largeur de bande

Cette fois ci, on modifie la largeur de bande (paramètre du noyau) en la divisant par 10.

On voit que la *bandwidth* règle bien la dispersion des données donnant des classes plus petites (et plus vraiment représentatives de la distribution des données). On est ici trop proche des points et donc mauvais en généralisation.

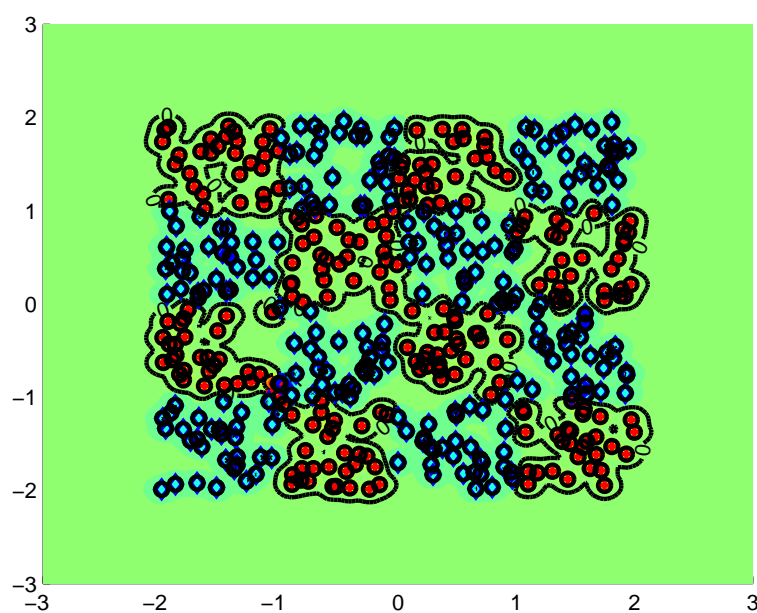


FIGURE 3 – Résultat avec bandwidth divisée par 10

## 4 Utilisation d'un kernel polynomial

On essaye maintenant un kernel polynomial au lieu de kernel gaussien utilisé précédemment.

Les degrés de 2 à 5 n'ont pas permis d'obtenir de résultats satisfaisants. Par contre, les degrés 6 et 7 ont permis d'obtenir des résultats satisfaisants. Les degrés supérieurs à 8 ont ensuite "divergés" et n'ont pas permis d'avoir de bons résultats.

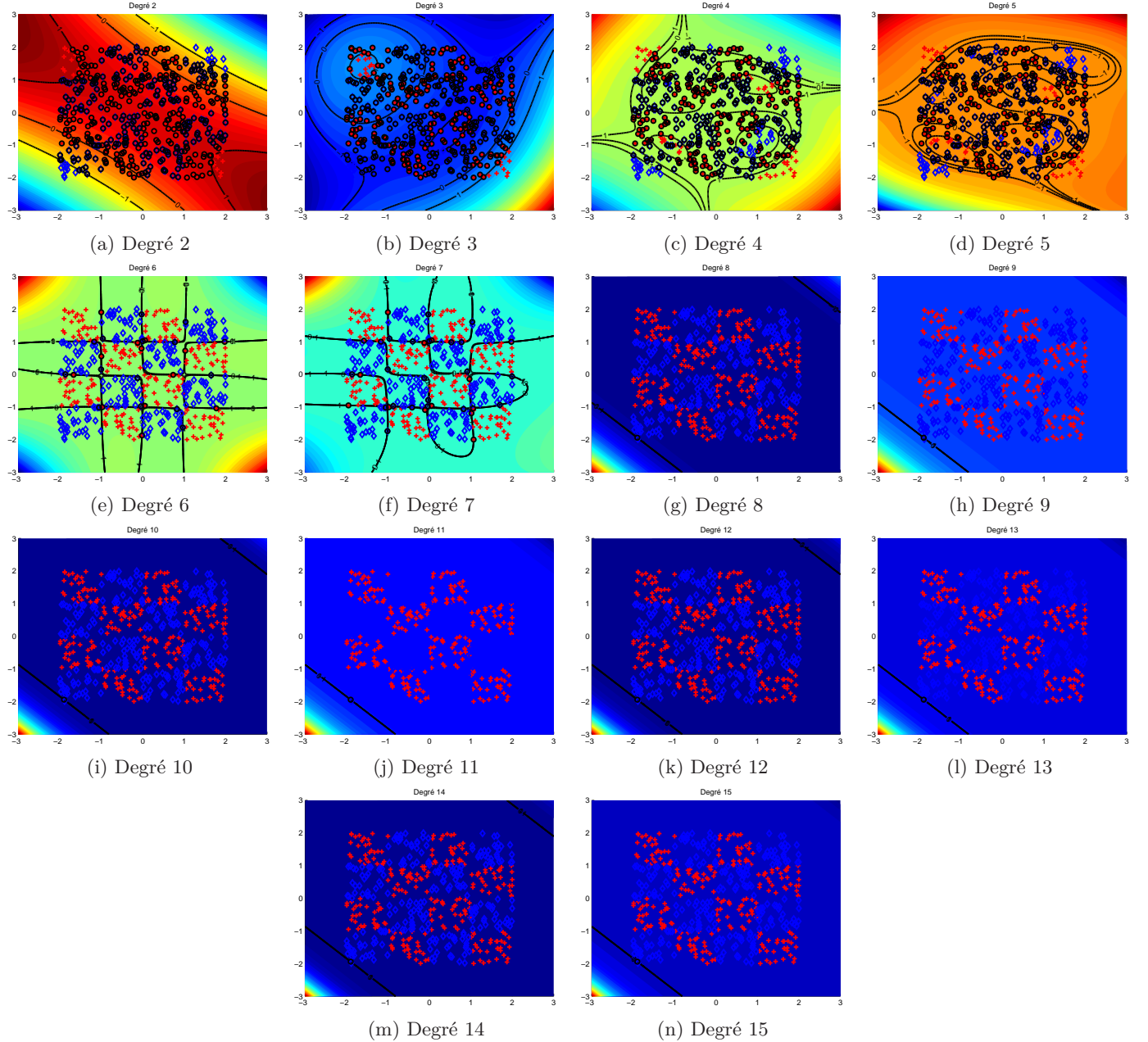


FIGURE 4 – Résultats avec kernel polynomial

## 5 SVM L2

Enfin, après un peu de maths pour trouver le dual du SVM L2, on a testé ce problème.

La différence de résultat avec le SVM L1 est faible et à l'œil nu, on voit peu de différences. On constate tout de même quelques variations pour un même jeu de données.

Le SVM L2 ayant moins de contraintes, il est plus rapide à résoudre par CVX que le SVM L1 : 5s contre 20s. Pour `monqp` par contre, aucune différence puisque cet algorithme ne "voit" de toute façon pas la disparition de la contrainte.

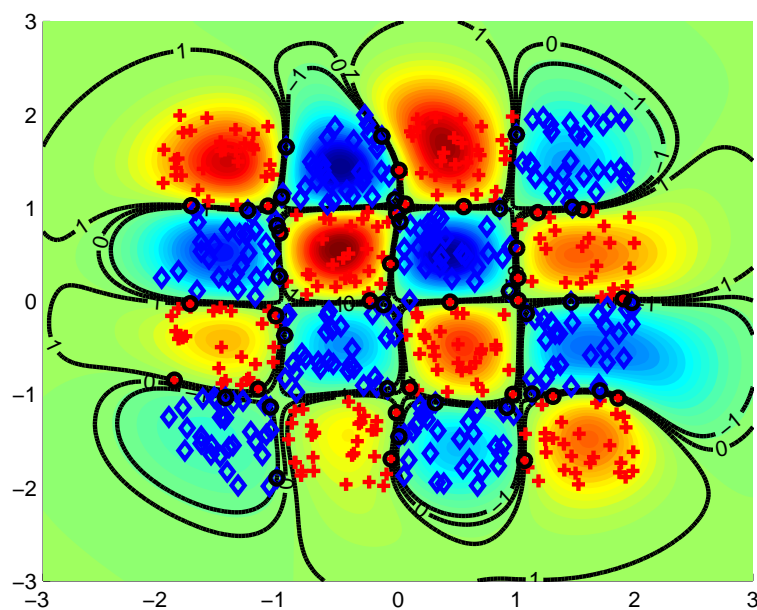


FIGURE 5 – Résultat de SVM L2

## 6 Annexes (code Matlab)

### 6.1 Temps de calcul

```

1 %% Question 1
2 %
3 % Ce code génère un jeu de données en forme de grille d'échec et l'affiche.
4
5 n = 500;
6 sigma=1.4;
7 [Xapp,yapp,Xtest ,ytest]=dataset_KM('checkers',n,n^2 ,sigma) ;
8 [n,p] = size(Xapp) ;
9 figure(1) ;
10 clf;
11 set(gcf,'Color',[1 ,1 ,1])
12 hold on
13 h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , 'r') ;
14 set(h1,'LineWidth',2) ;
15 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db') ;
16 set(h2,'LineWidth',2) ;
17
18 %% Question 2
19 %
20 % On calcule le kernel gaussien manuellement puis avec la fonction
21 % |svkernel|. On notera le paramètre |kerneloption| qui est l'écart type
22 % dans la formule du kernel gaussien.
23
24 % Compute a gaussian kernel and the matrix on your data with kerneloption =
25 % .5.
26 tic
27 D = (Xapp * Xapp') ; % produit scalaire
28 N = diag(D) ; % normes
29 D = -2*D + N*ones(1, n) + ones(n,1) *N'; %  $D_{ij} = ||x_i - x_j||^2 = -2 x_i' * x_j + ||x_i||^2 + ||x_j||^2$ 
30 kerneloption = .5;
31 s = 2 * kerneloption^2;
32 K = exp(-D/s);
33 G = (yapp*yapp') .* K;
34 toc
35
36 % Compute the same gaussian kernel using the svmkernel function of the
37 % SVMKM toolbox
38 kernel = 'gaussian';
39 kerneloption = .5;
40 tic
41 K=svkernel(Xapp, kernel, kerneloption);
42 G = (yapp*yapp') .* K;
43 toc
44
45 %% 2. c
46 %
47 % On résoud le problème de SVM dual avec cvx
48
49 e = ones(n,1);
50 C = 10000;
51 cvx_begin
52     variable a(n)
53     dual variables dp dC
54     minimize( 1/2*a'*G*a - e'*a )
55     subject to
56         de : yapp'*a == 0;
57         dp : a >= 0;
58         dC : a <= C;
59 cvx_end
60
61 %% 2.d
62 %
63 % On résoud cette fois le problème avec le solveur de problème quadratique
64 % |monqp|
65
66 tic
67 lambda = eps^.5;
68 [alpha ,b,pos] = monqp(G,e,yapp ,0 ,C,lambda ,0) ;
69 toc
70

```

```

71 %% Question 3
72 %
73 % Affichage du résultat avec un |meshgrid|
74
75 [xtest1 xtest2] = meshgrid([ -1:0.01:1]*3 , [ -1:0.01:1]*3) ;
76
77 nn = length(xtest1);
78 Xgrid = [reshape(xtest1 , nn*nn,1) reshape(xtest2 ,nn*nn,1) ] ;
79 Kgrid = svmkernel(Xgrid ,kernel ,kerneloption ,Xapp(pos ,:) ) ;
80 ypred = Kgrid*(yapp(pos) .*alpha) + b;
81 ypred = reshape(ypred ,nn,nn);
82 contourf(xtest1 ,xtest2 ,ypred ,50) ; shading flat;
83 hold on;
84 [cc,hh]=contour(xtest1 ,xtest2 ,ypred ,[ -1 0 1] , 'k') ;
85 clabel(cc,hh) ;
86 set(hh, 'LineWidth', 2) ;
87 h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , '+r' , 'LineWidth' ,2) ;
88 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' , 'LineWidth' ,2) ;
89 xsup = Xapp(pos ,:) ;
90 h3=plot(xsup(:,1) ,xsup(:,2) , 'ok' , 'LineWidth',2) ;
91 axis([ -3 3 -3 3]) ;
92
93
94 %% Conclusion sur les temps de calcul
95 %
96 % On voit bien que les méthodes de calcul optimisées sont plus rapide que
97 % les méthodes de calcul généralistes / manuelles.
98 %
99 % Par exemple, le calcul du kernel manuel (déjà légèrement optimisé) prend
100 % 0,18s alors que la méthode |svkernel| obtient le même résultat en 0,018s,
101 % c'est à dire 10 fois plus rapidement.
102 %
103 % Par ailleurs, la résolution du problème de minimisation par |cvx| prend
104 % 20.5s contre 0.075s pour |monqp|. On voit ici que |cvx| est une toolbox
105 % très pratique pour la souplesse qu'elle offre dans le formulation du
106 % problème de minimisation, mais qu'elle (et c'est logique) beaucoup moins
107 % rapide qu'une méthode de résolution optimisée.
108 %
109 % Pour avoir des temps de calcul raisonnable quand on augmente le nombre de
110 % données, il faut donc retravailler les problèmes optimisation pour
111 % revenir à une forme standard permettant d'utiliser des solveurs
112 % optimisés.

```

## 6.2 Test avec 10 000 données

```

1 %% Question 1
2 %
3 % Ce code génère un jeu de données en forme de grille d'échec et l'affiche.
4
5 n = 10000;
6 sigma=1.4;
7 [Xapp,yapp,Xtest ,ytest]=dataset_KM('checkers',n,n^2 ,sigma) ;
8 [n,p] = size(Xapp) ;
9 figure(1) ;
10 clf;
11 set(gcf,'Color',[1 ,1 ,1])
12 hold on
13 h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , '+r' ) ;
14 set(h1,'LineWidth',2) ;
15 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' ) ;
16 set(h2,'LineWidth',2) ;
17
18 %% Question 2
19 %
20 % On calcule le kernel gaussien manuellement puis avec la fonction
21 % |svkernel|. On notera le paramètre |kerneloption| qui est l'écart type
22 % dans la formule du kernel gaussien.
23
24 % Compute the same gaussian kernel using the svmkernel function of the
25 % SVMKM toolbox
26 kernel = 'gaussian';
27 kerneloption = .5;
28 K=svkernel(Xapp, kernel , kerneloption);
29 G = (yapp*yapp') .* K;
30

```

```

31 %% 2.d
32 %
33 % On résoud cette fois le problème avec le solveur de problème quadratique
34 % |monqp|
35
36 e = ones(n,1);
37 C = 10000;
38 lambda = eps^.5;
39 [alpha ,b,pos] = monqp(G,e,yapp ,0 ,C,lambda ,0) ;
40
41 %% Question 3
42 %
43 % Affichage du résultat avec un |meshgrid|
44
45 [xtest1 xtest2] = meshgrid([ -1:.01:1]*3 ,[ -1:0.01:1]*3) ;
46
47 nn = length(xtest1);
48 Xgrid = [reshape(xtest1 , nn*nn,1) reshape(xtest2 ,nn*nn,1) ];
49 Kgrid = svmkernel(Xgrid ,kernel ,kerneloption ,Xapp(pos ,:) ) ;
50 ypred = Kgrid*(yapp(pos) .*alpha) + b;
51 ypred = reshape(ypred ,nn,nn);
52 contourf(xtest1 ,xtest2 ,ypred ,50) ; shading flat;
53 hold on;
54 [cc,hh]=contour(xtest1 ,xtest2 ,ypred ,[ -1 0 1] , 'k') ;
55 clabel(cc,hh) ;
56 set(hh, 'LineWidth', 2) ;
57 h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , '+r' , 'LineWidth' ,2) ;
58 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' , 'LineWidth' ,2) ;
59 xsup = Xapp(pos ,:) ;
60 h3=plot(xsup(: ,1) ,xsup(: ,2) , 'ok' , 'LineWidth' ,2) ;
61 axis([ -3 3 -3 3]) ;
62
63 %%
64 % On utilise les méthodes les plus rapides uniquement compte tenu du nombre
65 % de données conséquent.

```

### 6.3 Modification de la largeur de bande

```

1 %% Question 1
2 %
3 % Ce code génère un jeu de données en forme de grille d'échec et l'affiche.
4
5 n = 500;
6 sigma=1.4;
7 [Xapp,yapp,Xtest ,ytest]=dataset_KM('checkers',n,n^2 ,sigma) ;
8 [n,p] = size(Xapp) ;
9 figure(1) ;
10 clf;
11 set(gcf,'Color',[1 ,1 ,1])
12 hold on
13 h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , '+r' ) ;
14 set(h1, 'LineWidth',2) ;
15 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' ) ;
16 set(h2, 'LineWidth',2) ;
17
18 %% Question 2
19 %
20 % On calcule le kernel gaussien manuellement puis avec la fonction
21 % |svmkernel|. On notera le paramètre |kerneloption| qui est l'écart type
22 % dans la formule du kernel gaussien.
23
24 % Compute the same gaussian kernel using the svmkernel function of the
25 % SVMKM toolbox
26 kernel = 'gaussian';
27 kerneloption = .05;
28 K=svmkernel(Xapp, kernel , kerneloption);
29 G = (yapp*yapp') .* K;
30
31 %% 2.d
32 %
33 % On résoud cette fois le problème avec le solveur de problème quadratique
34 % |monqp|
35
36 e = ones(n,1);
37 C = 10000;

```



```

38 lambda = eps^.5;
39 [alpha ,b,pos] = monqp(G,e,yapp ,0 ,C,lambda ,0) ;
40
41 %% Question 3
42 %
43 % Affichage du résultat avec un |meshgrid|
44
45 [xtest1 xtest2] = meshgrid([ -1:.01:1]*3 ,[ -1:0.01:1]*3) ;
46
47 nn = length(xtest1);
48 Xgrid = [reshape(xtest1 , nn*nn,1) reshape(xtest2 ,nn*nn,1) ];
49 Kgrid = svmkernel(Xgrid ,kernel ,kerneloption ,Xapp(pos ,:) ) ;
50 ypred = Kgrid*(yapp(pos) .*alpha) + b;
51 ypred = reshape(ypred,nn,nn);
52 contourf(xtest1 ,xtest2 ,ypred ,50) ; shading flat;
53 hold on;
54 [cc,hh]=contour(xtest1 ,xtest2 ,ypred ,[ -1 0 1] , 'k') ;
55 clabel(cc,hh) ;
56 set(hh, 'LineWidth', 2) ;
57 h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , 'r' , 'LineWidth' ,2) ;
58 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' , 'LineWidth' ,2) ;
59 xsup = Xapp(pos ,:) ;
60 h3=plot(xsup(: ,1) ,xsup(: ,2) , 'ok' , 'LineWidth' ,2) ;
61 axis([ -3 3 -3 3]) ;
62
63 %%
64 %
65 % On voit que la /bandwidth/ règle bien la dispersion des données donnant
66 % des classes plus petites (et plus vraiment représentatives de la
67 % distribution des données).

```

## 6.4 Utilisation d'un kernel polynomial

```

1 %% Question 1
2 %
3 % Ce code génère un jeu de données en forme de grille d'échec et l'affiche.
4
5 n = 500;
6 sigma=1.4;
7 [Xapp,yapp,Xtest ,ytest]=dataset_KM('checkers',n,n^2 ,sigma) ;
8 [n,p] = size(Xapp) ;
9
10
11 for kerneloption = 2:15
12
13     % Affichage
14     figure ;
15     clf;
16     set(gcf,'Color',[1 ,1 ,1])
17     hold on
18     h1=plot(Xapp(yapp==1 ,1) ,Xapp(yapp==1 ,2) , 'r' ) ;
19     set(h1,'LineWidth' ,2) ;
20     h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' ) ;
21     set(h2,'LineWidth' ,2) ;
22     title(['Degré ' int2str(kerneloption)]);
23
24     %% Question 2
25     %
26     % On calcule le kernel gaussien manuellement puis avec la fonction
27     % |svmkernel|. On notera le paramètre |kerneloption| qui est l'écart type
28     % dans la formule du kernel gaussien.
29
30     % Compute a gaussian kernel and the matrix on your data with kerneloption =
31     % .5.
32     % D = (Xapp * Xapp') ; % produit scalaire
33     % N = diag(D); % normes
34     % D = -2*D + N*ones(1, n) + ones(n,1) *N'; % Dij = ||t-s||^2 = -2 xi'*xj + ||xi||^2 + ||xj||^2
35     % kerneloption = .5;
36     % s = 2 * kerneloption^2;
37     % K = exp(-D/s);
38     % G = (yapp*yapp') .* K;
39
40     % Compute the same gaussian kernel using the svmkernel function of the
41     % SVMKM toolbox
42     kernel = 'poly';

```



```

43 K=svkernel(Xapp, kernel, kerneloption);
44 G = (yapp*yapp') .* K;
45
46 %% 2.d
47 %
48 % On résoud cette fois le problème avec le solveur de problème quadratique
49 % |monqp|
50
51 lambda = eps^.5;
52 [alpha, b, pos] = monqp(G, e, yapp, 0, C, lambda, 0);
53
54 %% Question 3
55 %
56 % Affichage du résultat avec un |meshgrid|
57
58 [xtest1, xtest2] = meshgrid([-1:0.01:1]*3, [-1:0.01:1]*3);
59
60 nn = length(xtest1);
61 Xgrid = [reshape(xtest1, nn*nn, 1), reshape(xtest2, nn*nn, 1)];
62 Kgrid = svkernel(Xgrid, kernel, kerneloption, Xapp(pos, :));
63 ypred = Kgrid*(yapp(pos) .* alpha) + b;
64 ypred = reshape(ypred, nn, nn);
65 contourf(xtest1, xtest2, ypred, 50); shading flat;
66 hold on;
67 [cc, hh] = contour(xtest1, xtest2, ypred, [-1 0 1], 'k');
68 clabel(cc, hh);
69 set(hh, 'LineWidth', 2);
70 h1 = plot(Xapp(yapp==1, 1), Xapp(yapp==1, 2), '+r', 'LineWidth', 2);
71 h2 = plot(Xapp(yapp== -1, 1), Xapp(yapp== -1, 2), 'db', 'LineWidth', 2);
72 xsup = Xapp(pos, :);
73 h3 = plot(xsup(:, 1), xsup(:, 2), 'ok', 'LineWidth', 2);
74 axis([-3 3 -3 3]);
75
76 end
77
78 %%
79 %
80 % Les degrés de 2 à 5 n'ont pas permis d'obtenir de résultats satisfaisants.
81 % Par contre, les degrés 6 et 7 ont permis d'obtenir des résultats
82 % satisfaisants. Les degrés supérieurs à 8 ont ensuite "divergés" et n'ont
83 % pas permis d'avoir de bons résultats.

```

## 6.5 SVM L2

```

1 %% Question 1
2 %
3 % Ce code génère un jeu de données en forme de grille d'échec et l'affiche.
4
5 n = 500;
6 sigma = 1.4;
7 [Xapp, yapp, Xtest, ytest] = dataset_KM('checkers', n, n^2, sigma);
8 [n, p] = size(Xapp);
9 figure(1);
10 clf;
11 set(gcf, 'Color', [1 1 1]);
12 hold on
13 h1 = plot(Xapp(yapp==1, 1), Xapp(yapp==1, 2), '+r');
14 set(h1, 'LineWidth', 2);
15 h2 = plot(Xapp(yapp== -1, 1), Xapp(yapp== -1, 2), 'db');
16 set(h2, 'LineWidth', 2);
17
18 %% Question 2
19
20 % Compute the same gaussian kernel using the svmkernel function of the
21 % SVMKM toolbox
22 kernel = 'gaussian';
23 kerneloption = .5;
24 K=svkernel(Xapp, kernel, kerneloption);
25 G = (yapp*yapp') .* K;
26
27 %% 2. c
28 %
29 % On résoud le problème de SVM dual avec cvx
30
31 e = ones(n, 1);

```

```

32 C = 10000;
33 cvx_begin
34     variable a(n)
35     dual variables de dp dC
36     minimize( 1/2*a'*(G+1/C)*a - e'*a )
37     subject to
38         de : yapp'*a == 0;
39         dp : a >= 0;
40 cvx_end
41
42 %% 2.d
43 %
44 % On résoud cette fois le problème avec le solveur de problème quadratique
45 % |monqp|
46
47 e = ones(n,1);
48 C = 10000;
49 lambda = eps^.5;
50 [alpha ,b,pos] = monqp(G,e,yapp ,0 ,C,lambda ,0) ;
51
52 %% Question 3
53 %
54 % Affichage du résultat avec un |meshgrid|
55
56 [xtest1 xtest2] = meshgrid([ -1:.01:1]*3 , [ -1:0.01:1]*3) ;
57
58 nn = length(xtest1);
59 Xgrid = [reshape(xtest1 , nn*nn,1) reshape(xtest2 ,nn*nn,1) ] ;
60 Kgrid = svmkernel(Xgrid ,kernel ,kerneloption ,Xapp(pos ,:) ) ;
61 ypred = Kgrid*(yapp(pos) .*alpha) + b;
62 ypred = reshape(ypred,nn,nn);
63 contourf(xtest1 ,xtest2 ,ypred ,50) ; shading flat;
64 hold on;
65 [cc,hh]=contour(xtest1 ,xtest2 ,ypred ,[ -1 0 1] , 'k') ;
66 clabel(cc,hh) ;
67 set(hh, 'LineWidth', 2) ;
68 h1=plot(Xapp(yapp==1 ,1) , Xapp(yapp==1 ,2) , 'r' , 'LineWidth' ,2) ;
69 h2=plot(Xapp(yapp== -1 ,1) ,Xapp(yapp== -1 ,2) , 'db' , 'LineWidth' ,2) ;
70 xsup = Xapp(pos ,:) ;
71 h3=plot(xsup(: ,1) ,xsup(: ,2) , 'ok' , 'LineWidth' ,2) ;
72 axis([ -3 3 -3 3]) ;
73
74 %% Conclusion
75 %
76 % En faisant du L2 SVM, on gagne passe de 20s à 5s de temps de calcul par
77 % CVX car on est moins contraint.

```