

# Data Mining

## TP11 - Support Vector Machine

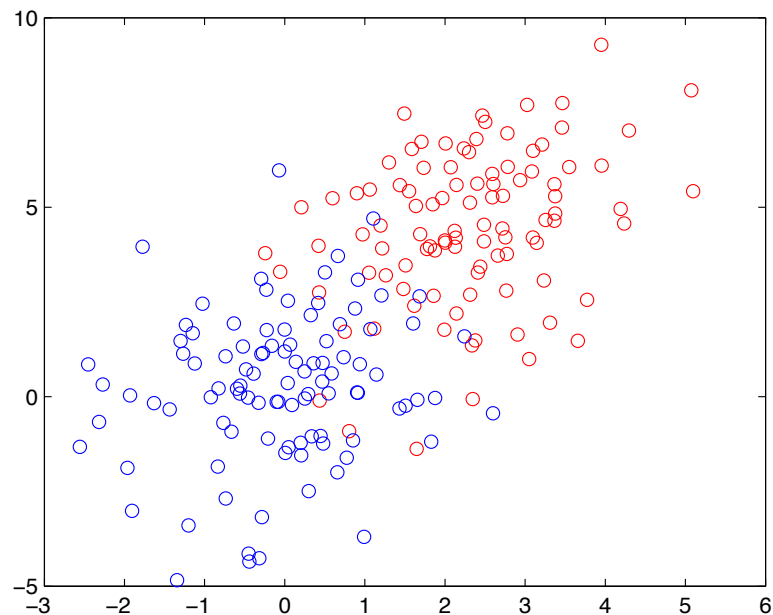
Rémi MUSSARD - Thomas ROBERT

### Partie 2.1

#### Q1. Dataset

Construction du jeu de données

```
1 n = 100;  
2  
3 X1 = randn(n, 2);  
4 X2 = randn(n, 2) + ones(n, 2)*2;  
5 S = [1 0.5; 0.5 4];  
6 X = [X1; X2]*S^(1/2);  
7 y = [ones(n,1); -ones(n,1)];  
8  
9 figure;  
10 plot(X(y==1,1), X(y==1,2), 'or');  
11 hold on  
12 plot(X(y==1,1), X(y==1, 2), 'ob');
```



#### Q2. Dual

Implémentation du problème de minimisation SVN dual sous CVX

```
1 C = 1000;
```

```
2
3 K = X*X';
4 Y = diag(y);
5 H = Y*K*Y;
6 q = ones(2*n, 1);
7
8 cvx_begin
9     variable a(2*n)
10
11     maximize(sum(a) - (1/2)*a'*H*a)
12
13     subject to
14
15         a>=0;
16         C*q>=a;
17         a'*y==0;
18
19 cvx_end
20
21 w = (a'*Y*X)'
```

```
w =

    -1.0586
    -0.5623
```

### Q3. Frontiere

On peut obtenir  $b$  en trouvant la frontière de décision grâce aux points supports et en trouvant donc directement  $b$ .

On peut également implémenter directement le problème primal sous CVX comme ci-dessous qui donne directement  $b$ .

```
1 cvx_begin
2     variable w(2)
3     variable b(1)
4     variable ksi(2*n)
5
6     minimize(1/2 * w'*w + C * sum(ksi) )
7
8     subject to
9
10         y .* (X * w + b) >= 1 - ksi;
11         ksi >= 0;
12
13 cvx_end
14
15 w
16 b
```

```
w =

    -1.0586
    -0.5623
```

```
b =

    2.4638
```

### Q3. Isocontours

```

1 % Isocontour pour  $f(x) = 0$ 
2
3 xFront = [min(X(:,1)) ; max(X(:,1))];
4 yFront0 = (-b-w(1)*xFront)/w(2);
5 plot(xFront, yFront0, '-g', 'LineWidth',2);
6
7 % Isocontour pour  $f(x) = 1$ 
8
9 yFront1 = (1 -b-w(1)*xFront)/w(2);
10 plot(xFront, yFront1, '-c');
11
12 % Isocontour pour  $f(x) = -1$ 
13
14 yFrontm1 = (-1 -b-w(1)*xFront)/w(2);
15 plot(xFront, yFrontm1, '-c');

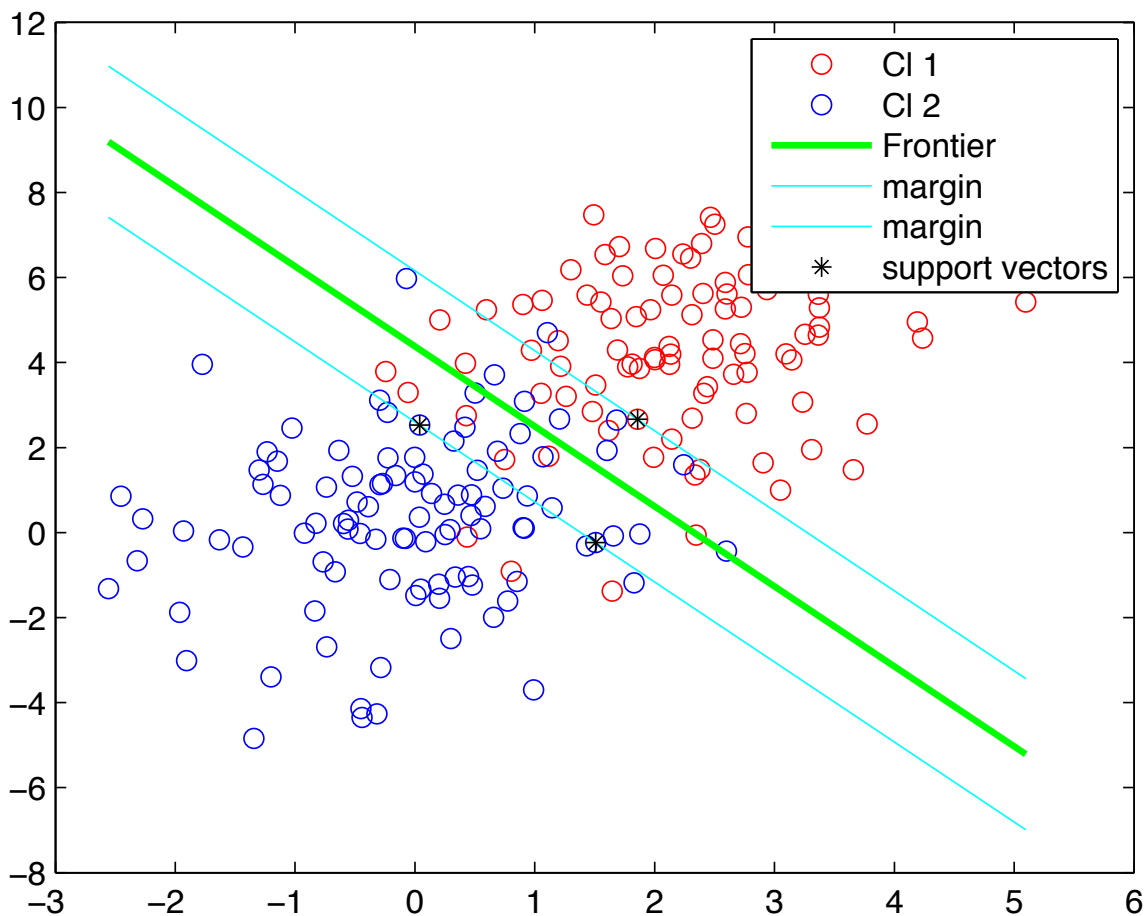
```

### Q5. Points supports

```

1 support = single(y.*(X*w + b));
2 plot(X(support==1,1), X(support==1, 2), '*k');
3
4 legend('Cl 1', 'Cl 2', 'Frontier', 'margin', 'margin', 'support vectors')

```



### Remarque

Avec cette méthode on se rend compte qu'il y a des erreurs de classification. Lorsqu'on a une donnée censé appartenir à une classe C1, avec la fonction `monsvmval`, on voit que cette donnée est affectée à la classe C2, car elle

se trouve du mauvais côté de la frontière de décision.

## Création des fonctions monsvclass.m et monsvval.m

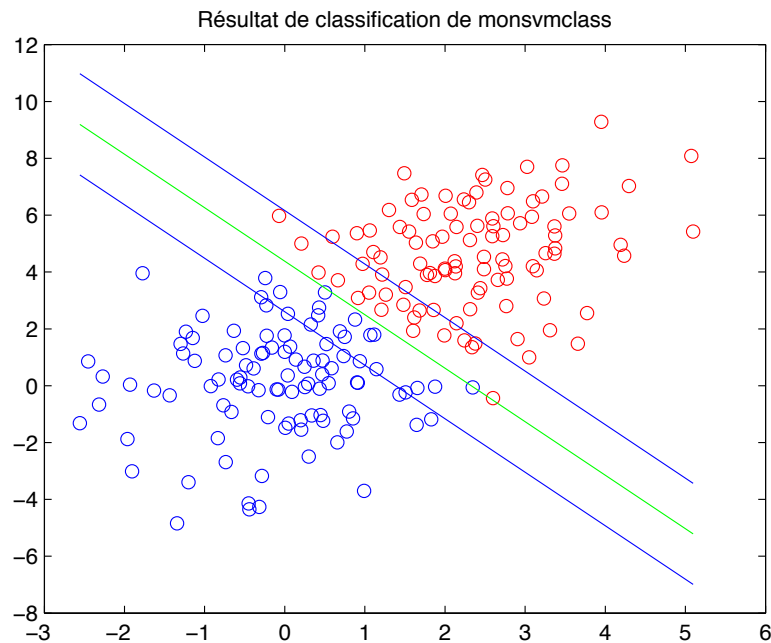
Voir les fichiers monsvclass.m et monsvval.m

Vérification des fonctions sur les données initiales

```

1 [w, b, alpha] = monsvclass(X, y, C);
2 w = w';
3 yReconst = monsvval(X, w, b);
4
5 figure();
6 plot(X(yReconst==-1,1), X(yReconst==-1,2), 'or');
7 hold on
8 plot(X(yReconst==1,1), X(yReconst==1, 2), 'ob');
9 title('Résultat de classification de monsvclass');
10
11 % Isocontour pour f(x) = 0
12
13 xFront = [min(X(:,1)) ; max(X(:,1))];
14 yFront0 = (-b-w(1)*xFront)/w(2);
15 plot(xFront, yFront0, '-g');
16
17 % Isocontour pour f(x) = 1
18
19 yFront1 = (1 -b-w(1)*xFront)/w(2);
20 plot(xFront, yFront1, '-b');
21
22 % Isocontour pour f(x) = - 1
23
24 yFrontm1 = (-1 -b-w(1)*xFront)/w(2);
25 plot(xFront, yFrontm1, '-b');

```



## Partie 2.2

```

1 clc
2 clear all

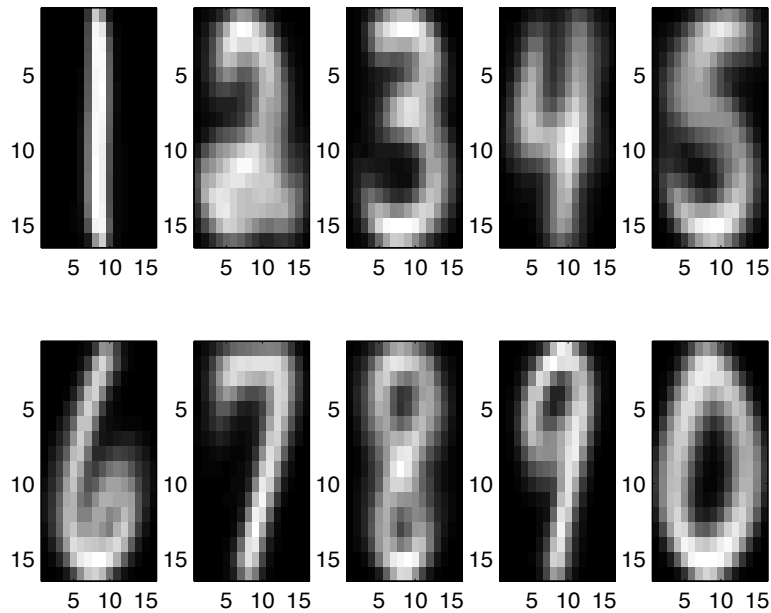
```

## Q2. chiffres moyens

```

1 m = size(x, 2);
2 Xmean = zeros(10, m);
3
4 for i=1:10
5     Xmean(i,:) = mean(x(y==i, :));
6 end
7
8
9 figure;
10 for i=1:10
11     subplot(2,5,i);
12     imagesc(reshape(Xmean(i,:),16,16)');
13     colormap gray;
14 end

```



## Q3. Discrimination des 1 et des 8

```

1 X = [x(y==1, :); x(y==8, :)];
2 Y = [y(y==1); y(y==8)];
3 Y(Y == 8) = -1;
4
5 [n, m] = size(X);
6
7 C = 0.001;
8
9 cvx_begin
10     variable w(m)
11     variable b(1)
12     variable ksi(n)
13
14     minimize(1/2 * w'*w + C * sum(ksi) )
15
16     subject to
17
18         Y .* (X * w + b) >= 1 - ksi;
19         ksi >= 0;
20
21 cvx_end
22
23 yReconst = monsvmval(X, w, b);
24 nbErreurs = sum(Y ~= yReconst)

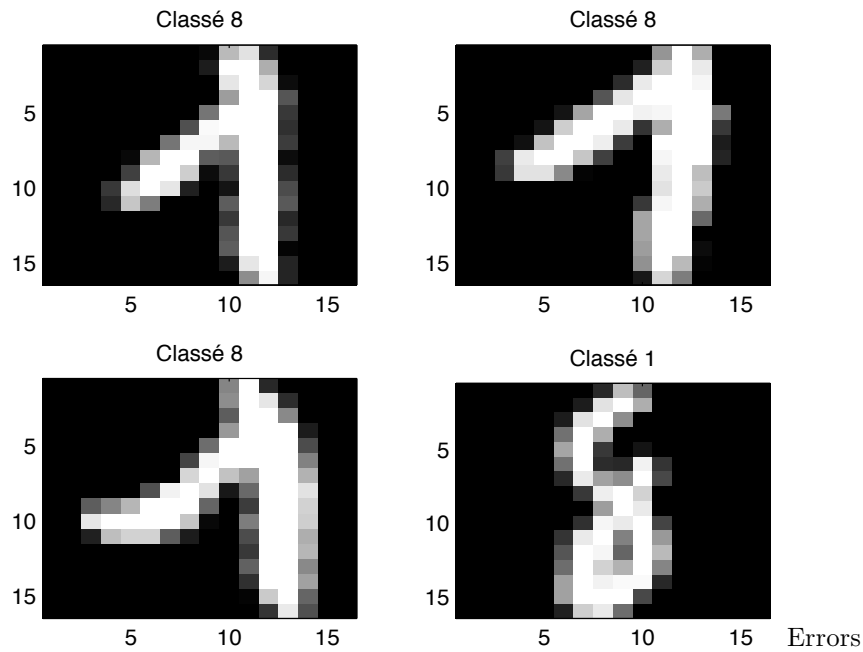
```

```

25
26 if (nbErreurs > 0)
27     Yerr = yReconst(Y ~= yReconst);
28     Yerr(Yerr == -1) = 8;
29     Xerr = X(Y ~= yReconst, :);
30
31     for k=1:min(4, size(Xerr,1))
32         subplot(2,2,k);
33         imagesc(reshape(Xerr(k,:),16,16)');
34         title(['Classé ' int2str(Yerr(k))]);
35     end
36 end

1 nbErreurs =
2
3     4

```



## Q5. Conclusion

Il n'y a quasiment aucune erreur. On obtient une erreur uniquement à partir de  $C < 0.001$ .