

Data Mining

TP3 - Classification Hiérarchique Ascendante

Rémi MUSSARD - Thomas ROBERT

1 Calcul de distance

```
1 function M = distance(X, param)
2 % X : R ^ N x d
3 % N : Nb points
4 % M : R ^ N x N distance entre les points
5
6     n = size(X,1);
7
8     if (param == 'euclid')
9         M = zeros(n,n);
10
11         %for i = 1:n
12         %     for j = i+1:n
13         %         dist = norm(X(i,:) - X(j,:));
14         %         M(i,j) = dist;
15         %         M(j,i) = dist;
16         %     end
17         %end
18
19         produits = X*X'; % produits scalaires
20         normes = diag(produits); % normes des vecteurs sur la diag
21         M = normes*ones(1,n) + ones(n,1)*normes' - 2* produits;
22     end
23 end
```

2 Fonctionnement d'aggclust

aggclust crée une hiérarchie ascendante des clusters. Il initialise le niveau 1 en créant 1 cluster par point. On boucle ensuite de 2 au nombre de points et en rassemblant à chaque fois les deux clusters les plus proches.

3 Fonction calc_dendro

Nous avons écrit une fonction `calc_dendro` qui se charge de calculer et d'afficher les deux dendrogrammes différents (un pour chaque méthode) à partir de données.

```
1 function [M, level, level_single] = calc_dendro(data, show)
2
3 M = distance(data, 'euclid');
4 M = M + diag(ones(1,size(M)));
5
6 level = aggclust(M, 'complete');
7 level_single = aggclust(M, 'single');
8
9 if (show)
```

```

10 subplot(1,2,1);
11 dendro(level);
12 title('Dendrogramme avec la méthode complete');
13
14 subplot(1,2,2);
15 dendro(level_single);
16 title('Dendrogramme avec la méthode single');
17 end

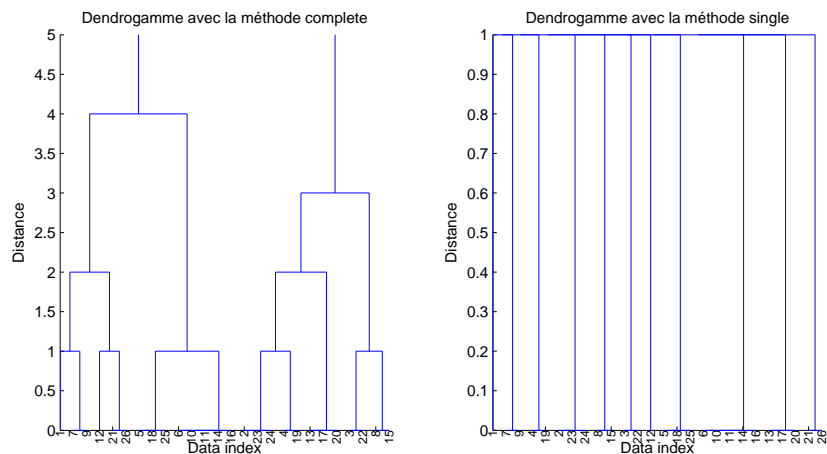
```

4 Classification ASI4

```

1 load('asi4.mat');
2 fig = figure();
3 calc_dendro(data, true);
4 set(fig, 'Position', [100 100 850 420]);

```



On voit que la méthode single n'arrive pas du tout à trouver de groupes alors que la méthode complete forme bien des groupes qui semblent (au vu du nombre d'étudiants par groupe) relativement bien répartis.

5 DS2

Afin de pouvoir visualiser graphiquement les clusters, on écrit une fonction `show_clusters(data, level, nbClust)` permettant d'afficher `nbClust` clusters.

```

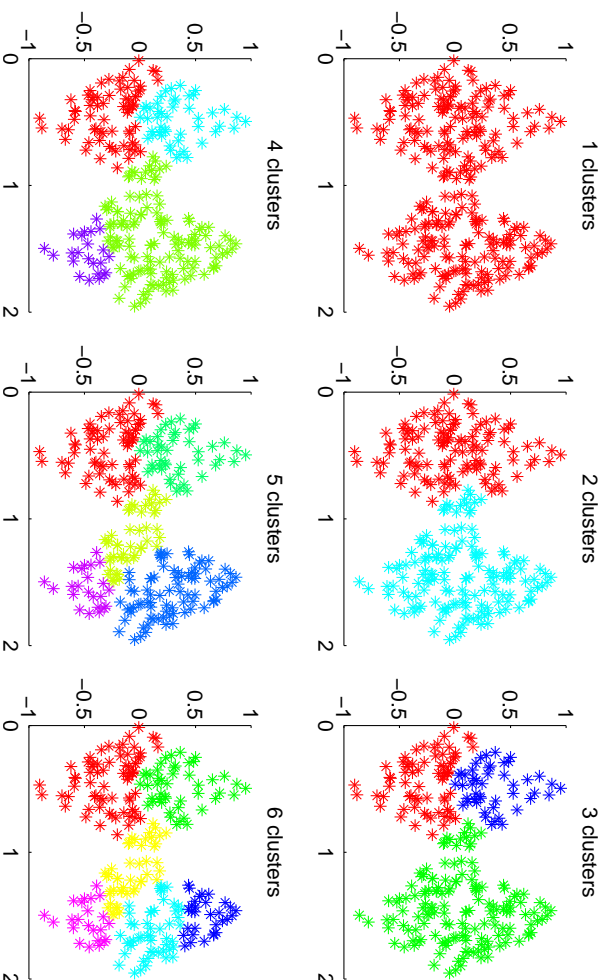
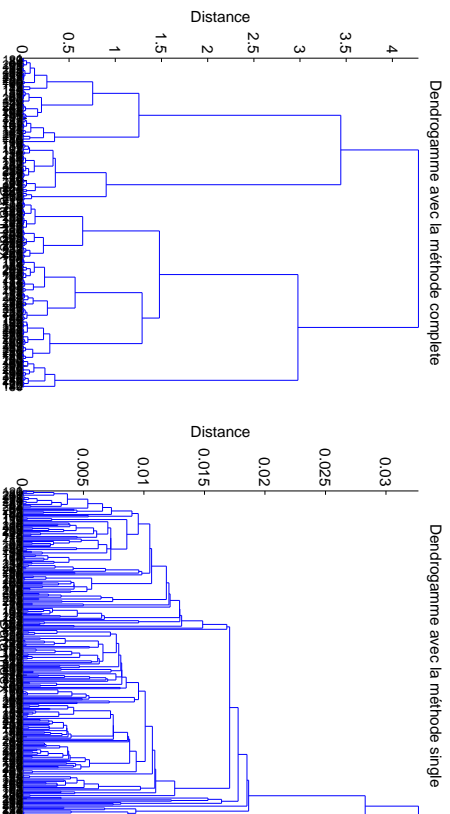
1 function show_clusters(data, level, nbClust)
2
3 colors = hsv(nbClust);
4
5 for i = 1:nbClust
6
7     inds = level(end - nbClust + 1).cluster{i};
8     dataClust = data(inds, :);
9     hold on;
10    plot(dataClust(:,1), dataClust(:,2), '*', 'color', colors(i, :));
11
12 end

```

```

1 load ds2.dat
2
3 data = mydownsampling(ds2, 30);
4
5 fig = figure();
6 [M, level, ~] = calc_dendro(data, true);
7 set(fig, 'Position', [100 100 850 420]);
8
9 % affichage
10
11 fig = figure();
12 for i=1:6
13     subplot(2,3,i);
14     show_clusters(data, level, i);
15     title([int2str(i) ' clusters']);
16 end
17 set(fig, 'Position', [100 100 750 420]);

```

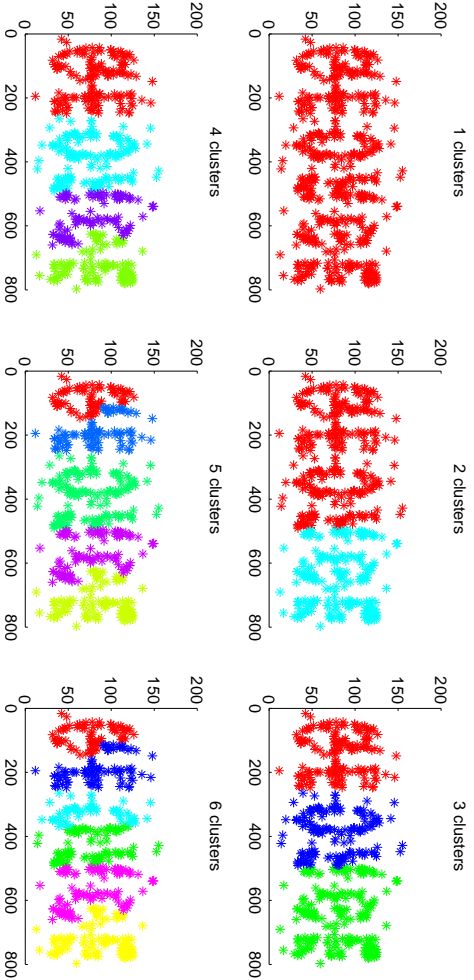
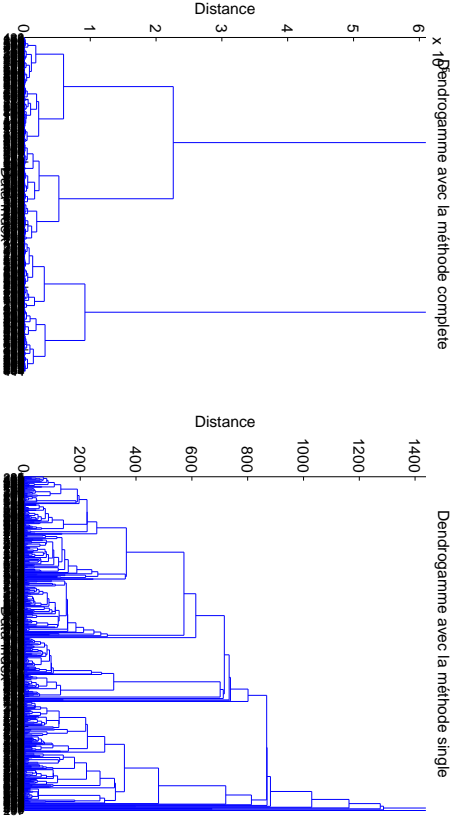


On constate que cette méthode ne donne pas toujours les deux clusters que l'on attendrait (un par losange) mais regroupe parfois les pointes (un cluster avec les pointes inférieures et un avec les pointes supérieures).

6 George

Avec les données "george", les clusters semblent plus proche de ce que l'on attends (1 cluster par lettre lorsque l'on choisit 6 clusters), et le résultat semble relativement stable vis à vis du sous-échantillonnage, contrairement à ce que l'on observait avec les données "ds2".

```
1 load george.dat
2
3 data = mydownsampling(george, 15);
4
5 fig = figure();
6 [M, level, ~] = calc_dendro(data, true);
7 set(fig, 'Position', [100 100 850 420]);
8
9 % affichage
10
11 fig = figure();
12 for i=1:6
13     subplot(2,3,i);
14     show_clusters(data, level, i);
15     title([int2str(i) ' clusters']);
16 end
17 set(fig, 'Position', [100 100 980 420]);
```



Data Mining

TP4 - Clustering K-means

Rémi MUSSARD - Thomas ROBERT

Voici les diverses fonctions codées

1 Calcul de distance

Cette fonction permet de calculer la distance entre chaque point des données de **X** par rapport à chaque centre **C** par la méthode **param**.

Elle retourne la matrice de distance **M**.

```
1 function M = distance(X, C, param)
2
3     N = size(X,1);
4     K = size(C,1);
5
6     % M = zeros(N,K);
7     %
8     % for i = 1:N
9     %     for k = 1:K
10    %         if (param == 'euclid')
11    %             M(i,k) = norm(X(i,:) - C(k,:));
12    %         end
13    %     end
14    % end
15
16    if (strcmp(param, 'euclid'))
17        normX = sum(X.^2, 2);
18        normC = sum(C.^2, 2);
19        ps = X*C';
20        M = repmat(normX, 1, K) + repmat(normC', N, 1) - 2*ps;
21    end
22
23    if (strcmp(param, 'mahal'))
24
25        M = zeros(N,K);
26
27        for k = 1:K
28            X2 = X + repmat(mean(X),N,1) - repmat(C(k,:),N,1);
29            M(:, k) = mahal(X2, X);
30        end
31    end
32 end
```

../distance.m

2 Affectation aux clusters

Cette fonction permet d'affecter chaque point au cluster dont le centre est le plus proche du point, grâce à la matrice de distance **M**.

Elle retourne les affectations sous forme d'une liste **liste** avec pour chaque point son numéro de cluster.

```
1 function liste = affectation(M)
2     [~, liste] = min(M,[], 2); % indices des minimums par ligne
3 end
```

../affectation.m

3 Coût

Cette fonction permet de calculer le critère d'inertie intra-cluster à partir de la matrice des distances **M** et des affectations **liste**.

Elle retourne les coût **Jw**.

```
1 function Jw = cout(M, liste)
2
3     [N, K] = size(M);
4
5     Jw = 0;
6
7     % pour chaque cluster
8     for k = 1:K
9         % distance intra-cluster du cluster k
10        Jw = Jw + sum(M(liste == k, k));
11    end
12
13    % Jw = 0;
14    % for i = 1:N
15    %     Jw = Jw + M(i, liste(i));
16    % end
17 end
```

../cout.m

4 Affectation et coût

Cette fonction combine les deux précédentes en une seule plus efficace.

```
1 function [Jw, liste] = affectation_cout(M)
2     [couts, liste] = min(M,[], 2); % indices des minimums par ligne
3     Jw = sum(couts);
4 end
```

../affectation_cout.m

5 Calcul des centres

Calcule les coordonnées des centres des clusters de **X** définis par la liste d'affectation **liste**.

Retourne les positions des centres **C**.

```
1 function C = nouveaux_centres(X, liste)
2     K = max(liste);
3     d = size(X, 2);
4
5     C = zeros(K, d);
6
7     % pour chaque cluster
```

```

8   for k = 1:K
9       % moyenne des points dans le cluster
10      C(k, :) = mean(X(liste == k, :));
11  end
12 end

```

../nouveaux_centres.m

6 Initialisation des centres

Calcule des coordonnées de **K** centres aléatoires dont les coordonnées sont incluses dans le nuage de points **X** (encadrées par les min et max de chaque variable).

Retourne les positions des centres initiaux **C**.

```

1 function C = init_centres(X, K)
2
3     d = size(X, 2);
4
5     % points centrés entre 0 et 1
6     C = rand(K, d) - repmat(0.5, K, d);
7
8     % étendue des données
9     etendue = max(X) - min(X);
10    centre = (max(X) + min(X)) / 2;
11
12    % multiplication par l'étendue
13    C = C .* repmat(etendue, K, 1);
14
15    % centrage de C sur X
16    C = C + repmat(centre, K, 1);
17 end

```

../init_centres.m

7 Calcul des clusters par la méthode des K-moyennes

Calcule les clusters de façon itérative jusqu'à stabilisation du critère de coût à un minimum local, sur les données **X**, à partir des centres initiaux **C0**, avec la méthode de calcul de distance **param**.

Retourne les centres finaux **C**, la liste des affectations finales **liste**, et l'évolution du critère de coût **Jw**.

```

1 function [C, liste, Jw] = k_moyennes(X, C0, param)
2     C = C0;
3
4     % init des couts
5     JwPrec = 1;
6     JwNew = 0;
7     Jw = [];
8
9     % tant que le cout varie de plus de 1e-3
10    while abs(JwPrec - JwNew) > 1e-3
11        JwPrec = JwNew;
12
13        % calcul des distances
14        M = distance(X, C, param);
15        % affectation et couts
16        [JwNew, liste] = affectation_cout(M);
17        Jw = [Jw; JwNew];
18        % calcul des nouveaux centres
19        C = nouveaux_centres(X, liste);
20    end
21 end

```

../k_moyennes.m

8 Calcul des formes fortes à partir des affectations

Découpe les clusters calculés en formes fortes à partir des multiples affectations **listes**. Chaque colonne représente les affectations pour une itération.

Retourne les formes fortes sous forme de liste d'affectation **newlist**.

```

1 function newlist = formes_fortes_from_listes(listes)
2
3     N = size(listes,1);
4
5     [~, ~, listeformesfortes]=unique(listes, 'rows') ;
6
7     Nff = max(listeformesfortes);
8     effectifs = zeros(Nff, 1);
9     newlist = zeros(N,1);
10
11     for c=1:Nff
12         ind = find(listeformesfortes==c) ;
13         effectifs(c)=length(ind);
14         newlist(ind)=c;
15     end
16
17 end

```

../formes_fortes_from_listes.m

9 Calcul des formes fortes

Calcule des formes fortes sur les données **X**, en appliquant **nbIte** fois la méthode k-moyennes avec **K** clusters et la fonctions de distance **param**.

Retourne les formes fortes sous forme de liste d'affectation **newlist**.

```

1 function newlist = formes_fortes(X, nbIte, param, K)
2
3     N = size(X,1);
4
5     listes = zeros(N,nbIte);
6
7     for j = 1:nbIte
8         C0 = init_centres(X, K);
9         [~, liste, ~] = k_moyennes(X, C0, param);
10        listes(:,j) = liste;
11    end
12
13    newlist = formes_fortes_from_listes(listes);
14
15 end

```

../formes_fortes.m

10 Test des méthodes de k-moyennes et formes fortes

Cette fonction permet de tester la méthode des k-moyennes et des formes fortes sur les données **X**, pour les nombres de clusters contenus dans le vecteur **Klist**, en réalisant **nbIte** à la fin desquelles est appliquée la méthode

des formes fortes. La méthode de calcul de distance **param** est utilisée.

Les résultats de clustering sont affichés sur la figure **figClusters**. Les variations de la fonction coût sont affichés dans la figure **figJw**.

```

1 function test_k_means(X, Klist, nbIte, param, figClusters, figJw)
2     subI = length(Klist);
3     subJ = nbIte + 1;
4
5     % pour chaque k
6     for i = 1:length(Klist)
7         K = Klist(i);
8         listes = zeros(size(X,1), nbIte);
9
10        % pour chaque itération
11        for j = 1:nbIte
12
13            % on initialise les K centres
14            C0 = init_centres(X, K);
15            % on calcule nos clusters
16            [C, liste, Jw] = k_moyennes(X, C0, param);
17            listes(:,j) = liste;
18
19            % affichage des clusters et des centres initiaux
20            figure(figClusters);
21            subplot(subI, subJ, (i - 1)*subJ + j);
22            show_clusters(X, liste);
23            title(['K = ' num2str(K) ' / Jw = ' num2str(Jw(end))]);
24            hold on
25            scatter(C(:,1),C(:,2), 50, 'k', 'd', 'fill');
26            scatter(C0(:,1),C0(:,2), 50, 'k', 'd');
27
28            % affichage de l'évolution de Jw
29            figure(figJw);
30            subplot(subI, subJ - 1, (i-1)*(subJ-1) + j);
31            plot(Jw, '-*');
32            title(['K = ' num2str(K)]);
33        end
34
35        % calcul des formes fortes à partir des itérations
36        newlist = formes_fortes_from_listes(listes);
37
38        % affichage formes fortes
39        figure(figClusters);
40        subplot(subI, subJ, i*subJ);
41        show_clusters(X, newlist);
42        title(['Formes fortes / Jw = ' num2str(cout(X, newlist))]);
43    end
44 end

```

../test_k_means.m

11 Lancement des tests

Ce script permet de lancer la procédure de test ci-dessus sur chaque jeu de données, avec les deux méthodes de calcul de distance.

```

1 % DM / TP4 / MUSSARD / ROBERT
2
3 clear all
4 close all
5
6 % Liste des méthodes à appliquer.
7 % Pour chaque méthode, liste des fichiers de données ainsi que des
8 % paramètres de test.
9 methodes = struct(
10     'methode',{ 'euclid'; 'mahal'},

```

```

11     'tests', {
12         struct(
13             'file', {'ds2.dat' ; 'ds3.dat' ; 'ds4.dat' ; 'ds5.dat' ; 'george.dat'}, ...
14             'Klist', {2:7 ; 2:6 ; 4:9 ; 3:8 ; 2:6},
15             'nbIte', {3 ; 3 ; 3 ; 3 ; 3});
16         struct(
17             'file', {'ds2.dat' ; 'ds3.dat' ; 'ds4.dat' ; 'ds5.dat' ; 'george.dat'}, ...
18             'Klist', {2:7 ; 2:6 ; 4:9 ; 3:8 ; 2:6},
19             'nbIte', {3 ; 3 ; 3 ; 3 ; 3})
20     });
21
22 % pour chaque méthode (euclid, mahal, ...)
23 for i = 1:length(methodes)
24
25     methode = methodes(i).methode;
26     tests = methodes(i).tests;
27
28 % pour chaque jeu de données
29 for j = 1:length(tests)
30
31     % init des données
32     X = load(tests(j).file);
33     Klist = tests(j).Klist;
34     nbIte = tests(j).nbIte;
35
36     % init des figures
37     figClusters = figure();
38     figJw = figure();
39
40     % execution du test
41     test_k_means(X, Klist, nbIte, methode, figClusters, figJw);
42
43     % sauvegarde des figures
44     set(figClusters, 'PaperUnits', 'points');
45     set(figClusters, 'PaperPosition', [0 0 1000 1000]);
46     saveas(figClusters, ['TP4_MUSSARD_ROBERT_' int2str(figClusters) '.png']);
47     set(figJw, 'PaperUnits', 'points');
48     set(figJw, 'PaperPosition', [0 0 500 600]);
49     saveas(figJw, ['TP4_MUSSARD_ROBERT_' int2str(figJw) '.png']);
50
51 end
52 end

```

../TP4_MUSSARD_ROBERT.m

12 Résultats

Ci-après sont présent les résultats des tests sur chaque figure, avec chaque méthode, avec plusieurs valeurs de K , plusieurs itérations suivies du calcul des formes fortes.

Des diamants indiquent les centres initiaux et finaux. Par ailleurs, un graphe permet de visualiser l'évolution de J_w .

On remarque que les clusters trouvés par cette méthode ne respectent pas les formes dessinées par les points, ce qui est logique puisque cette méthode cherche uniquement à créer des clusters de points les plus proches les uns des autres, alors que nos données comportent des groupes de points qui s'entremêlent et qui sont par ailleurs tous très proches les uns des autres. C'est pourquoi le résultat trouvé est souvent proche d'une simple partition de la "zone de dessin" en "parts" égales.

Concernant la distance de Mahalanobis, elle permet d'annuler la distorsion due à l'étalement des données. Mais pour certains jeu de données, l'étalement à une importance et une signification, et le fait de ne plus en tenir compte détériore le résultat. C'est par exemple le cas avec le jeu de données George où l'étalement des lettres n'est plus pris en compte et où les clusters sont plutôt en forme d'étoile autour du centre des données.

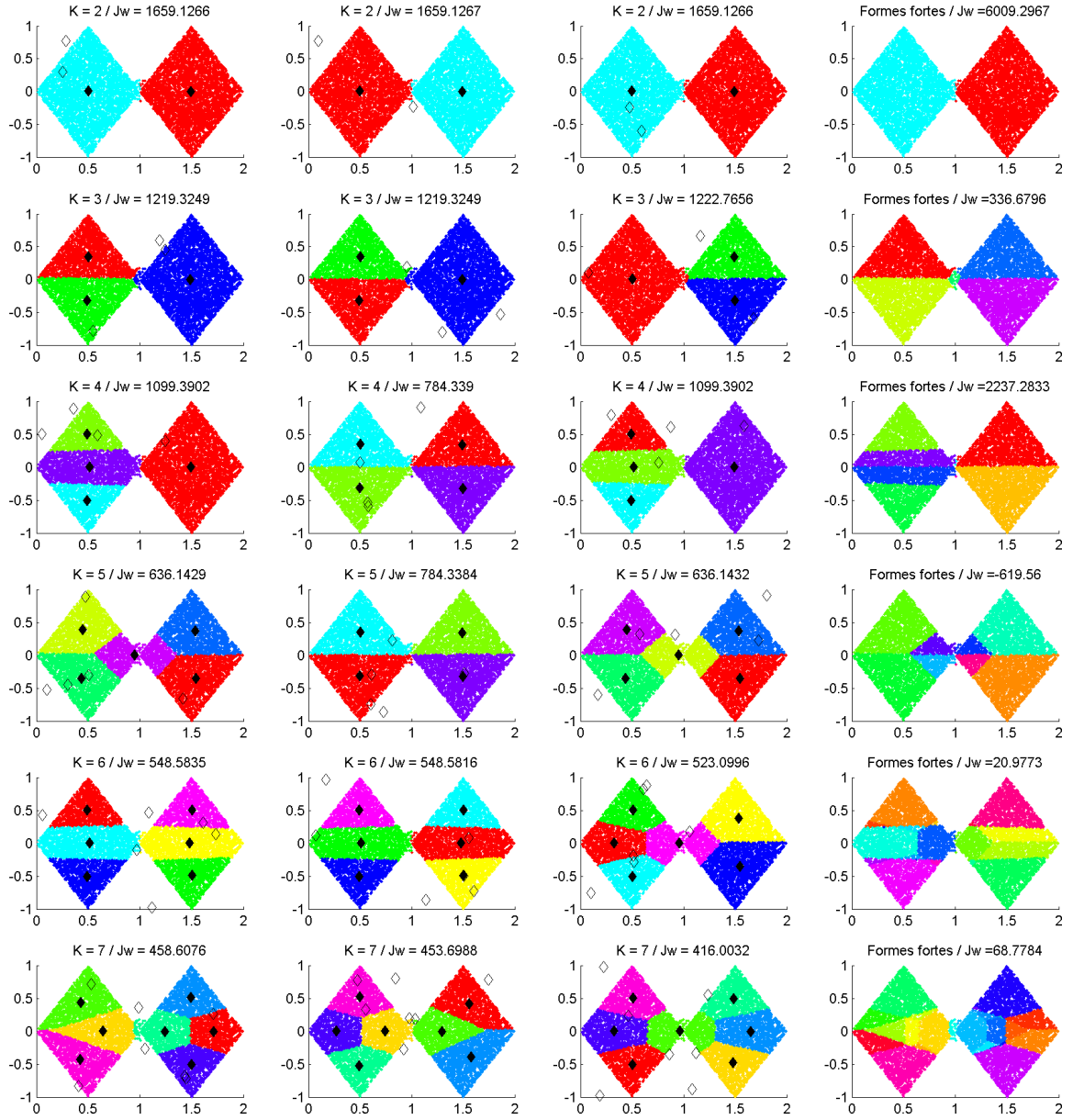


FIGURE 1 – ds2 - Distance euclidienne

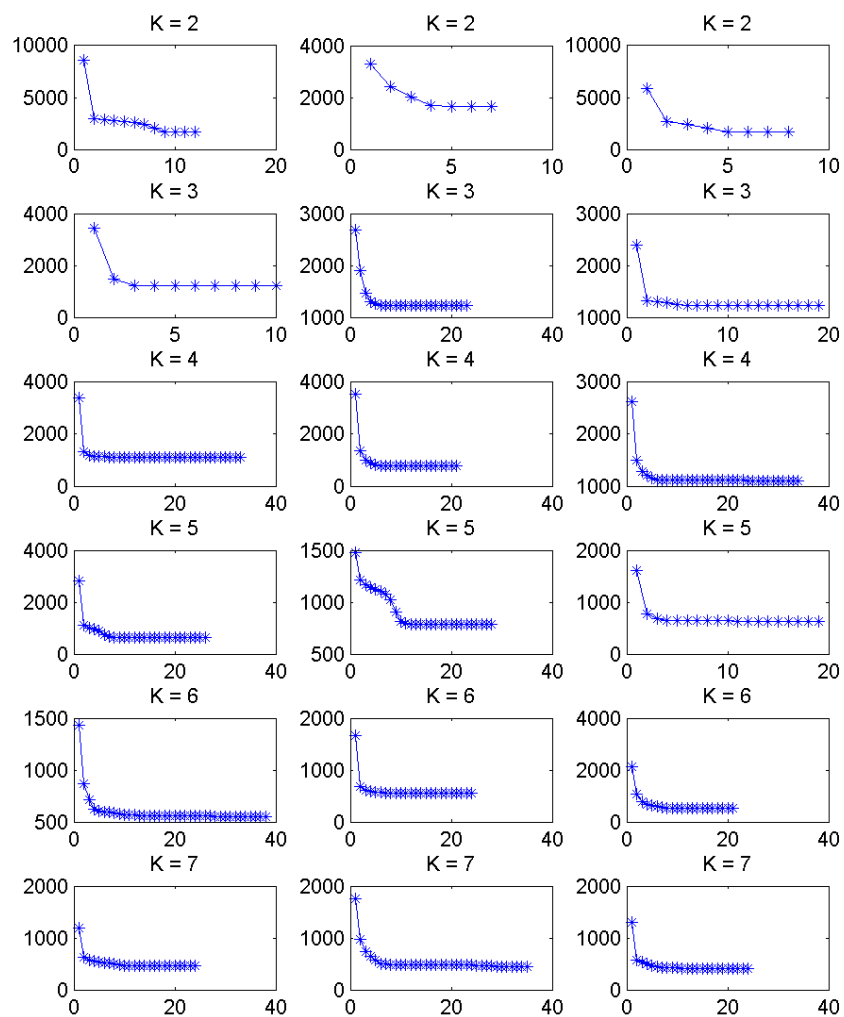


FIGURE 2 – DS2 - Distance euclidienne

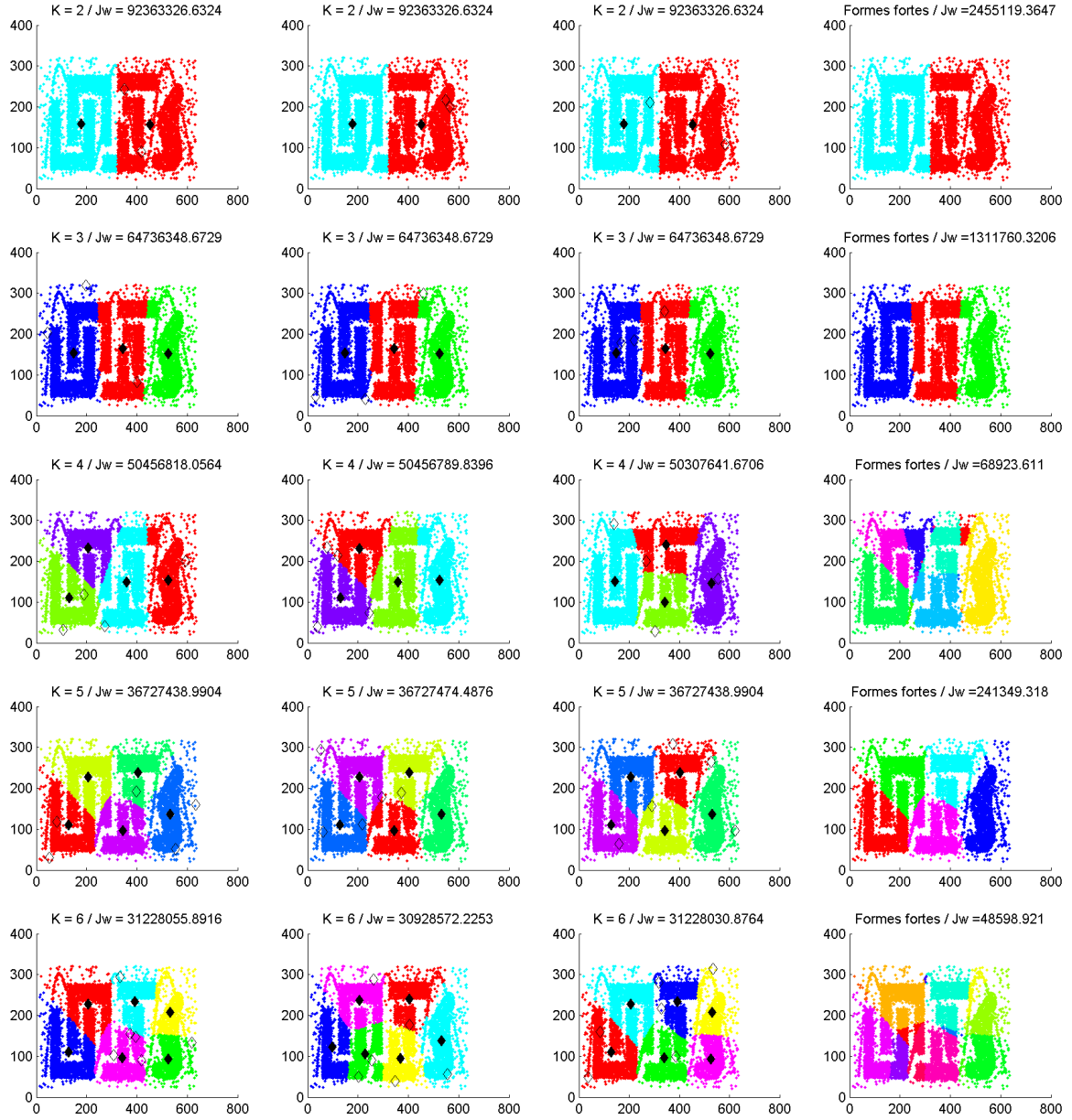


FIGURE 3 – DS3 - Distance euclidienne

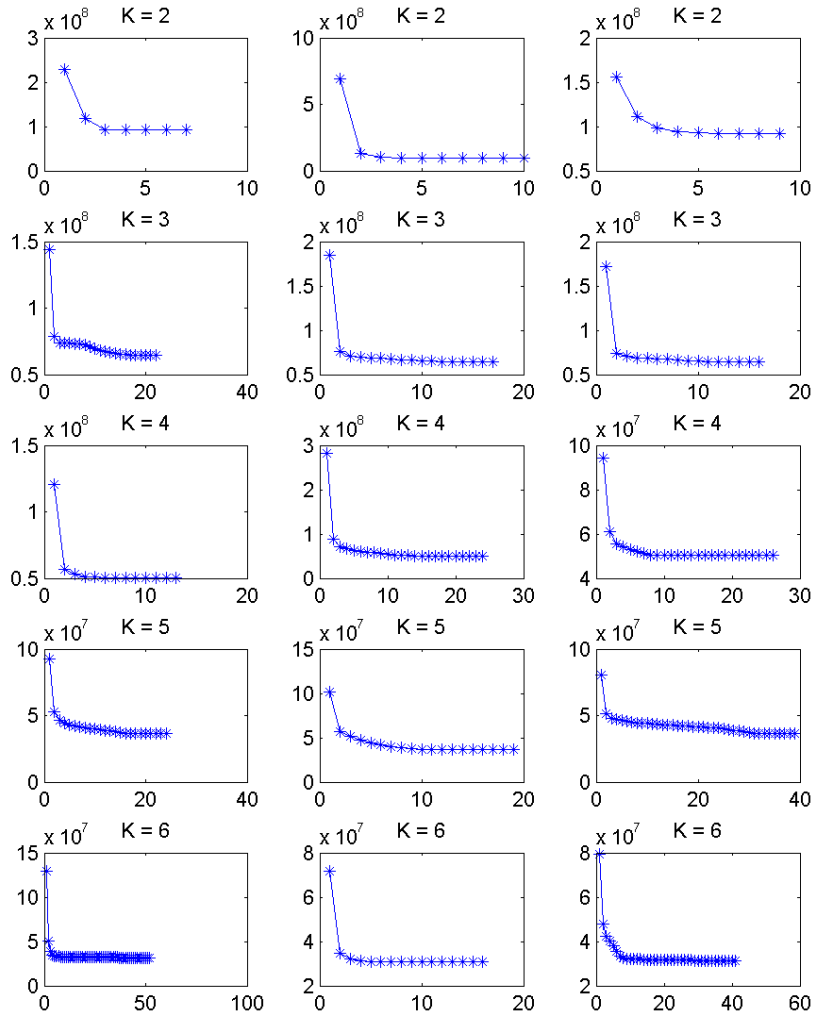


FIGURE 4 – DS3 - Distance euclidienne

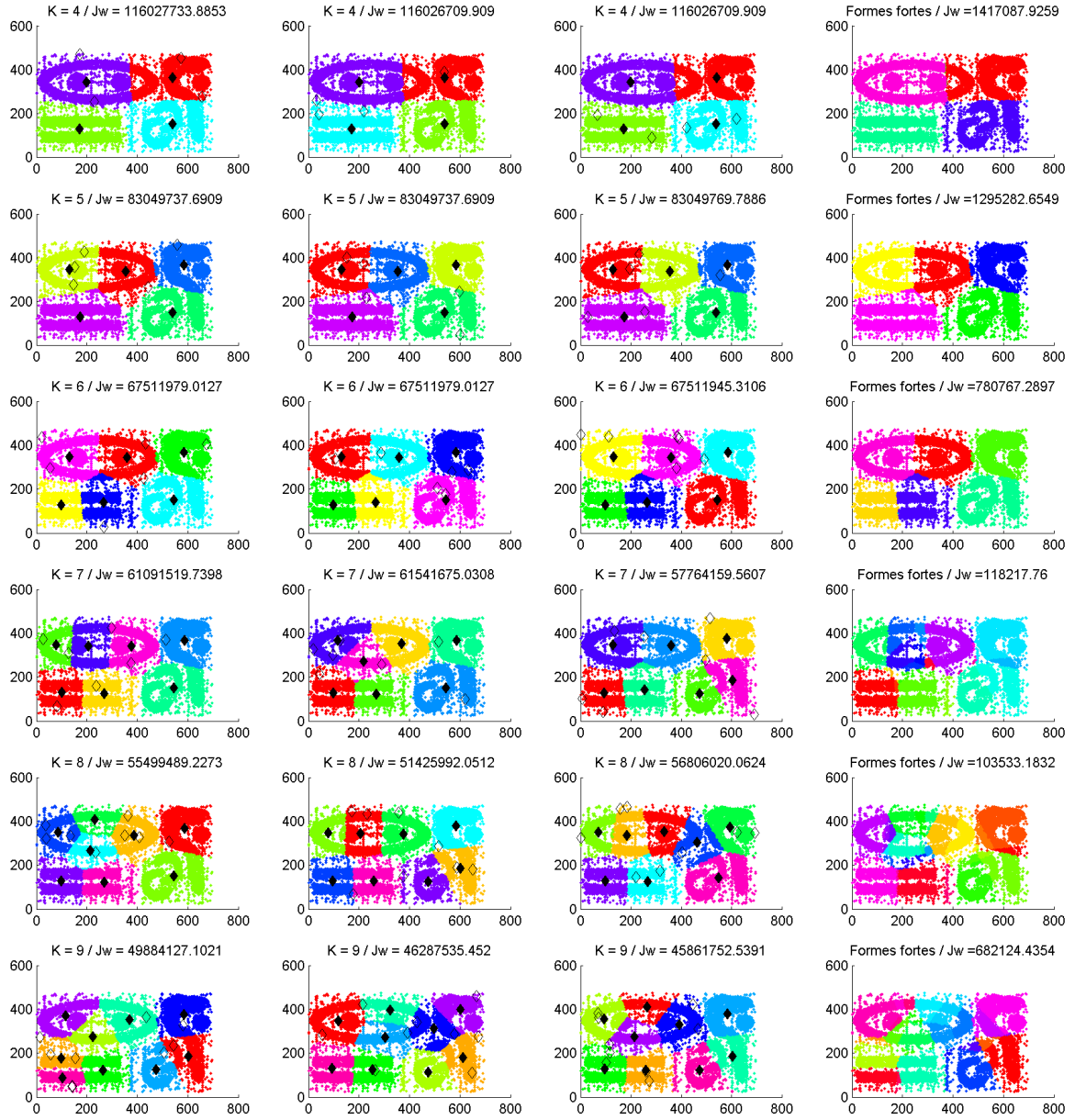


FIGURE 5 – DS4 - Distance euclidienne

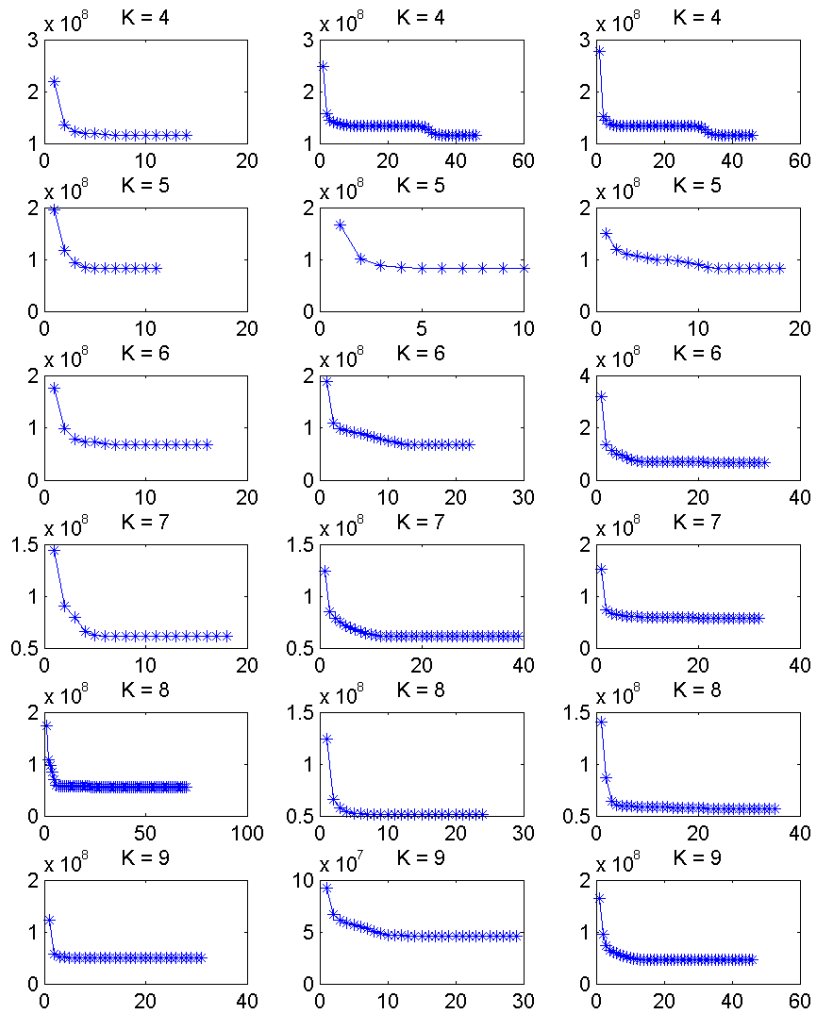


FIGURE 6 – DS4 - Distance euclidienne

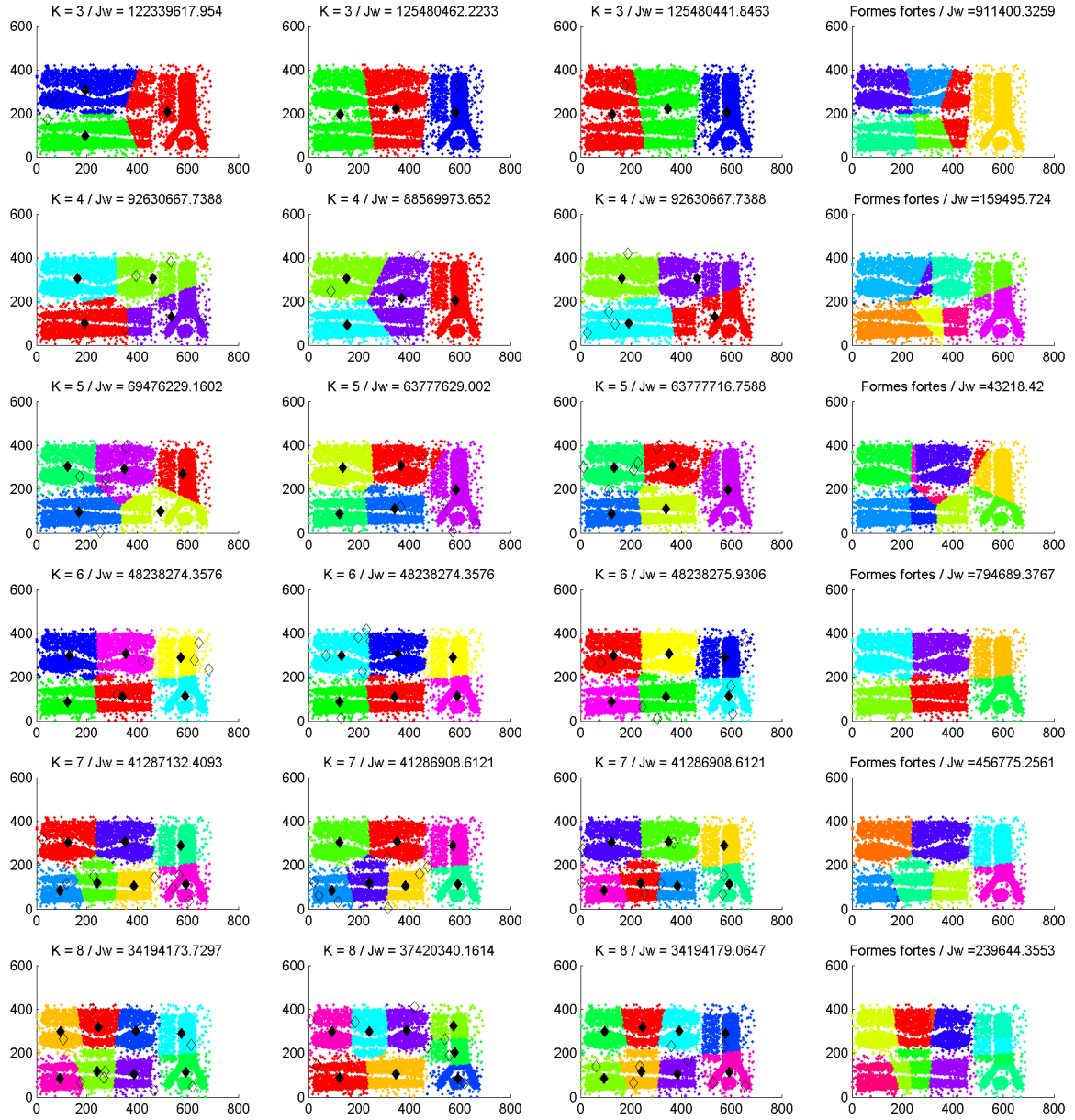


FIGURE 7 – DS5 - Distance euclidienne

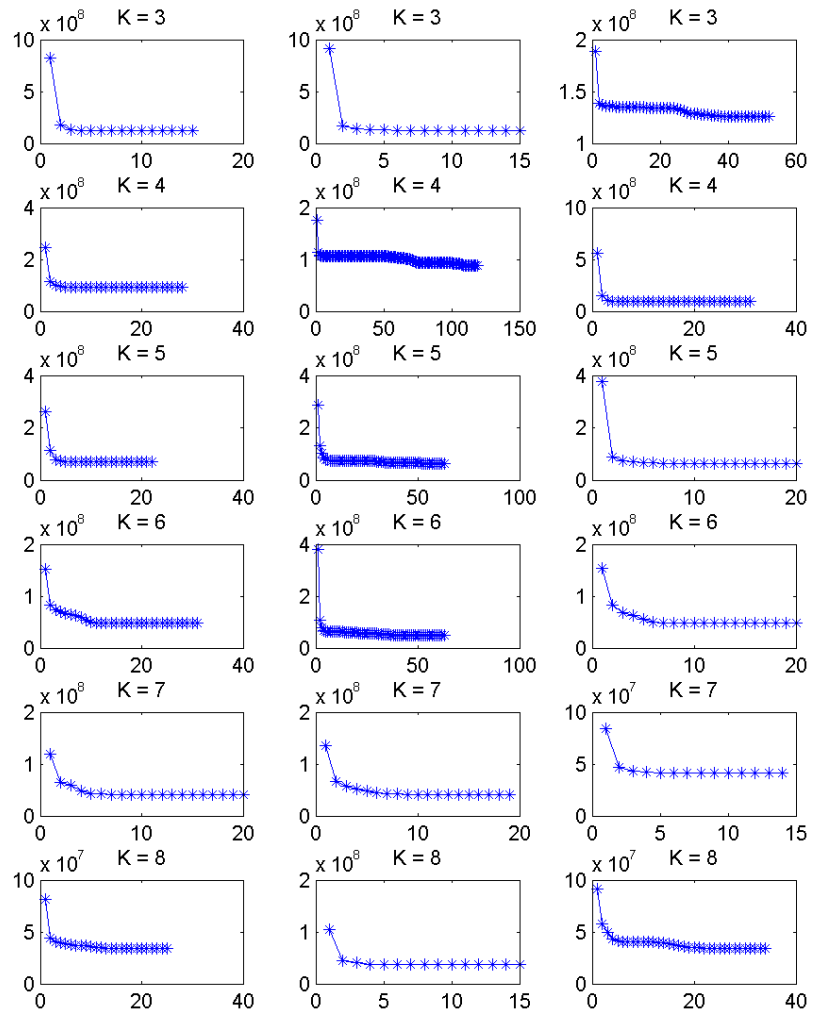


FIGURE 8 – DS5 - Distance euclidienne

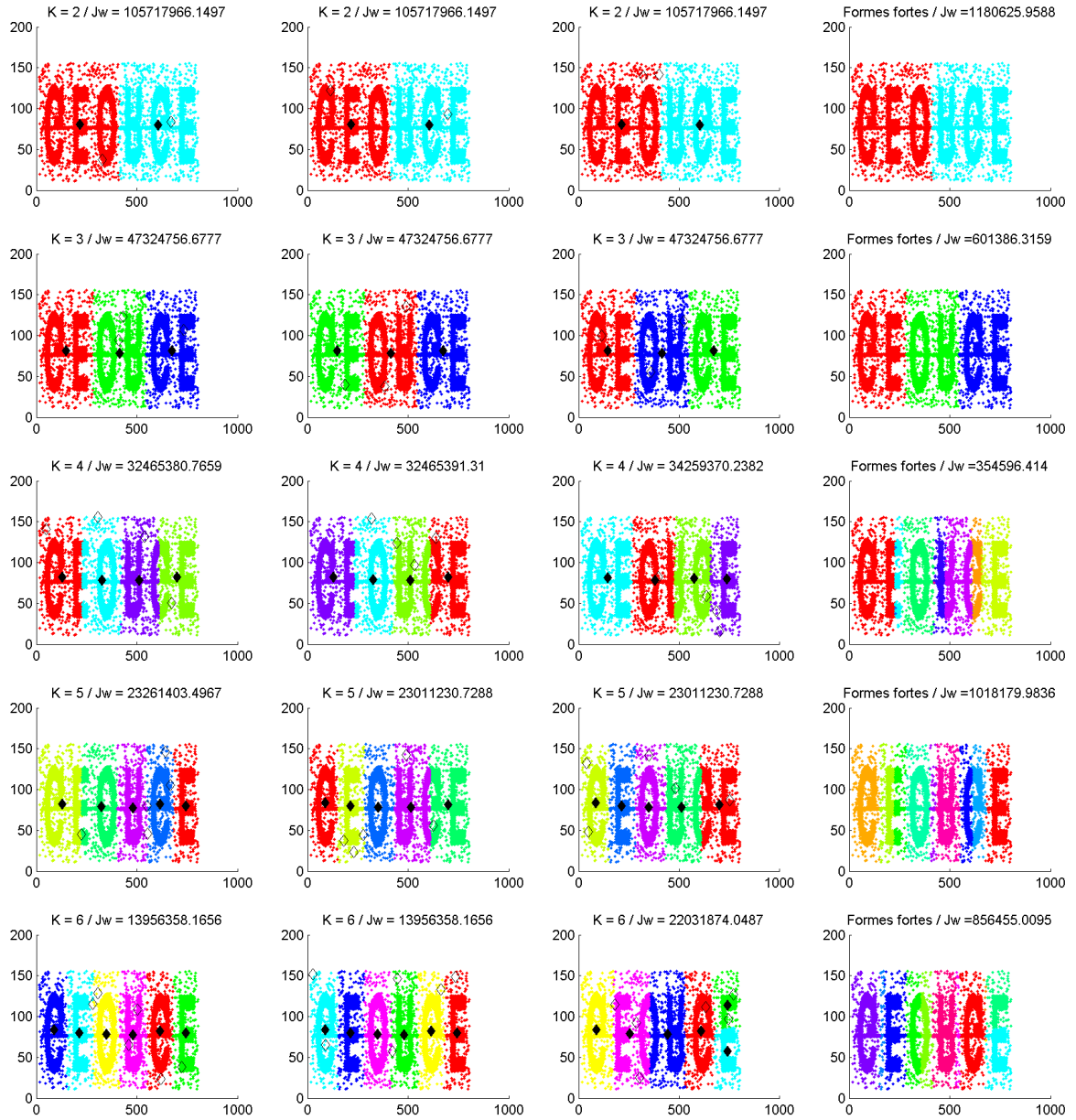


FIGURE 9 – George - Distance euclidienne

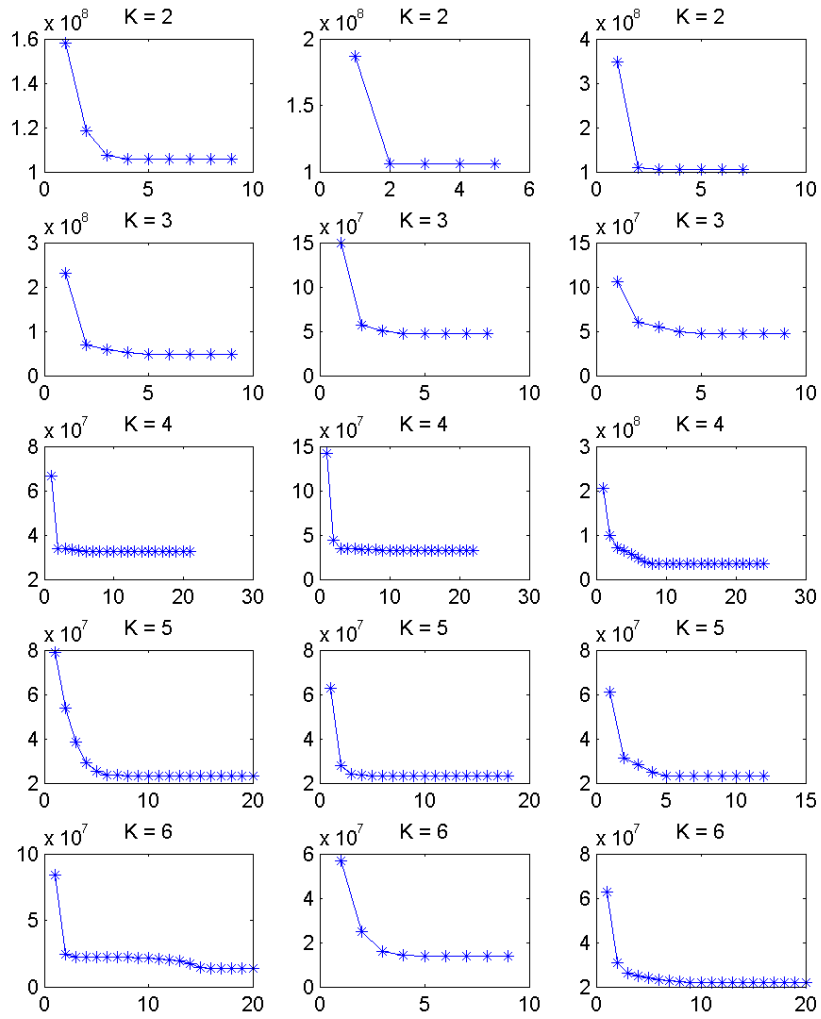


FIGURE 10 – George - Distance euclidienne

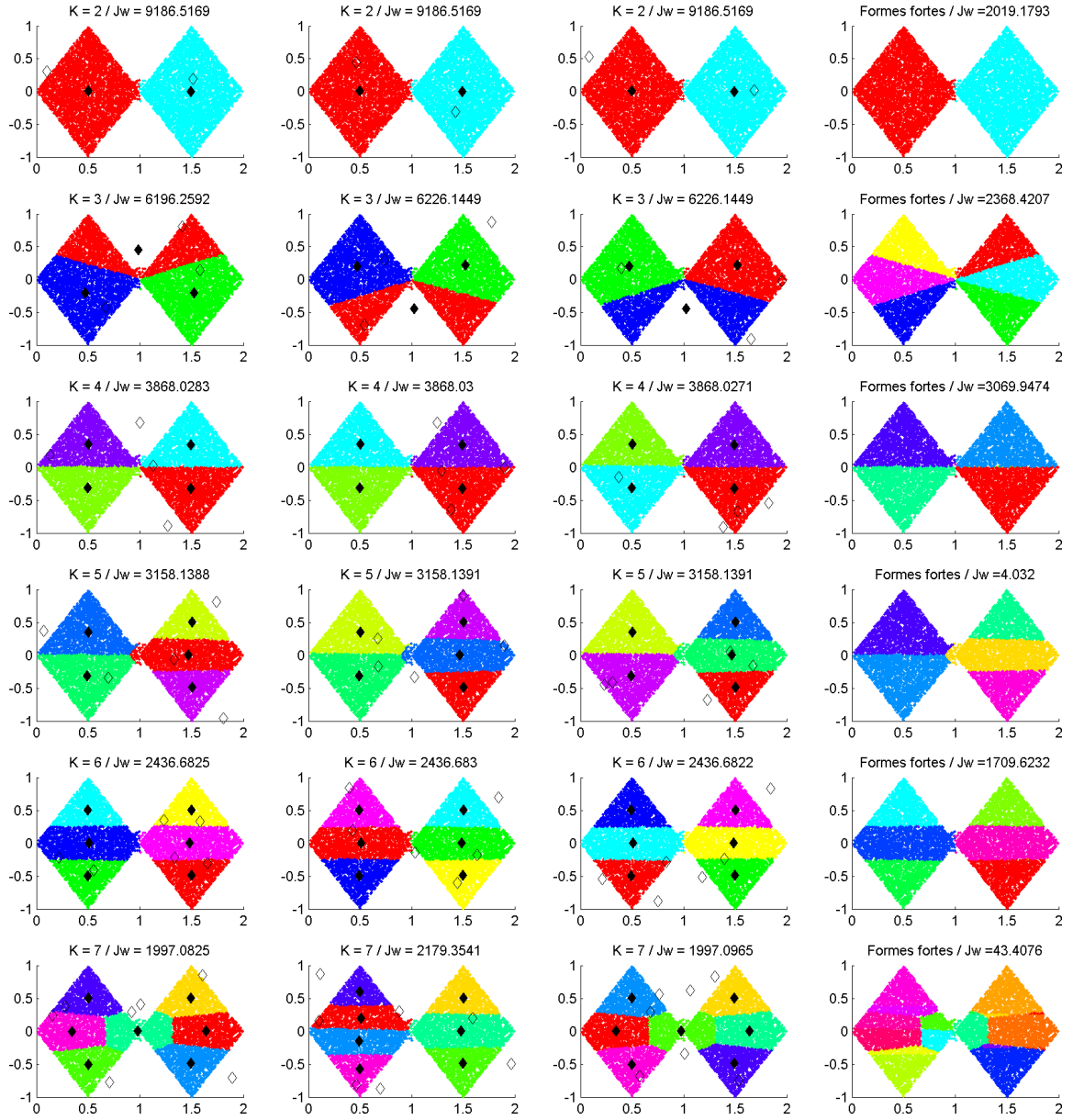


FIGURE 11 – DS2 - Distance de Mahalanobis

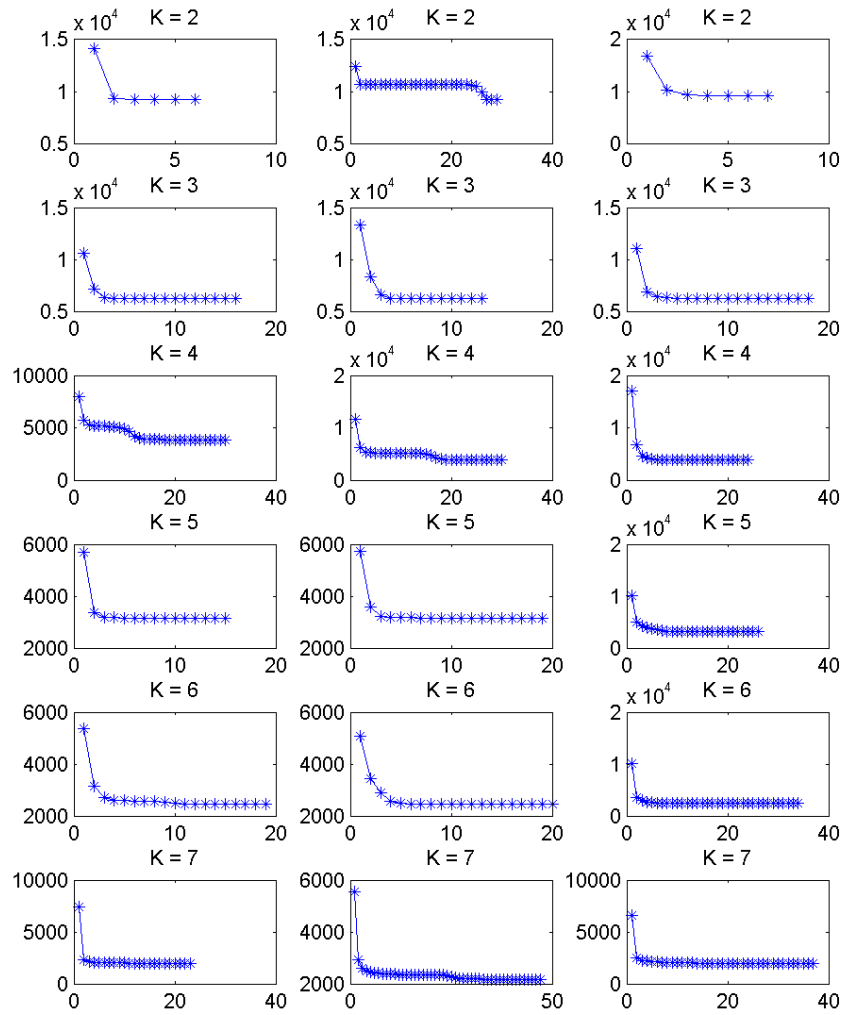


FIGURE 12 – DS2 - Distance de Mahalanobis

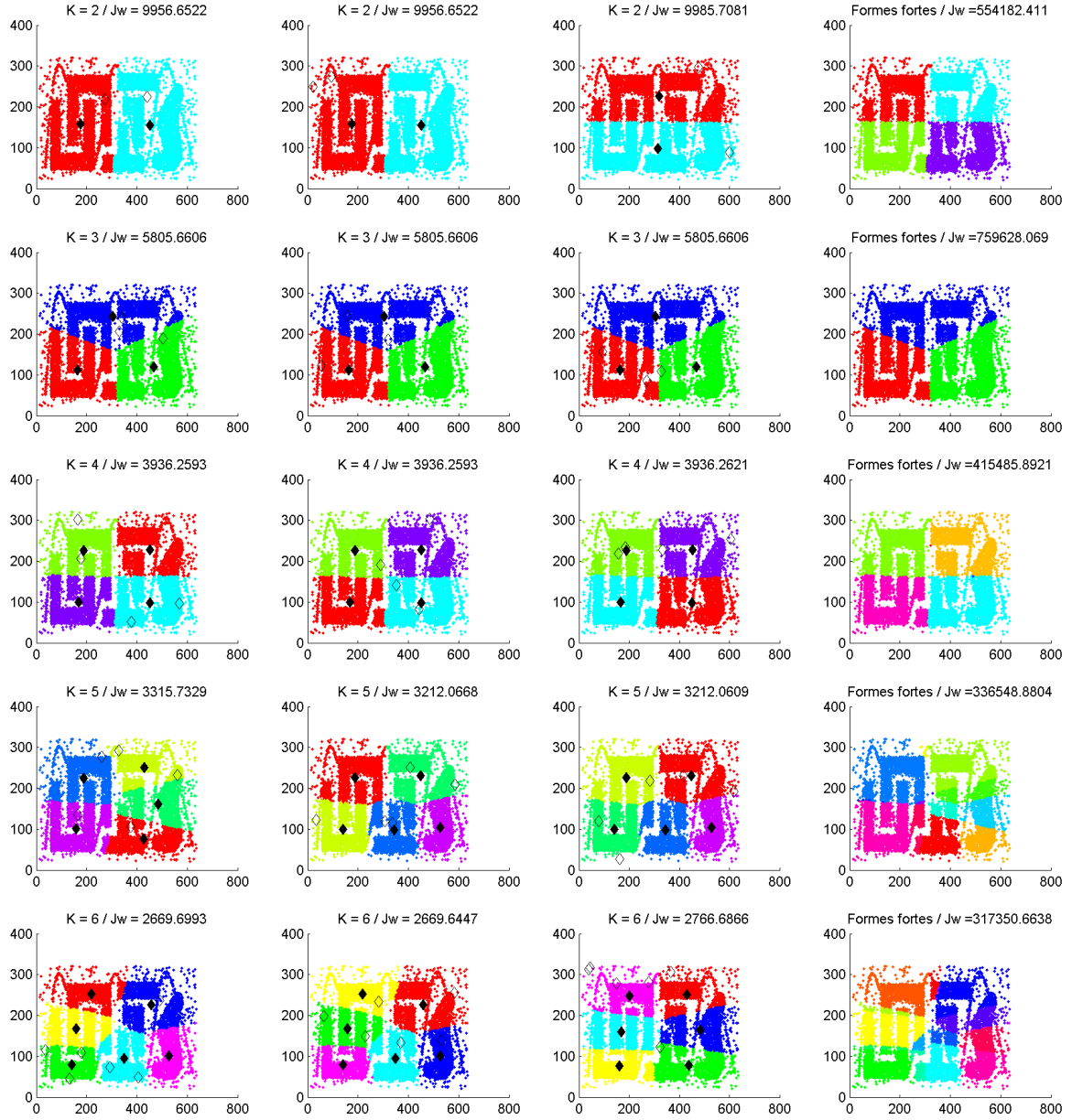


FIGURE 13 – DS3 - Distance de Mahalanobis

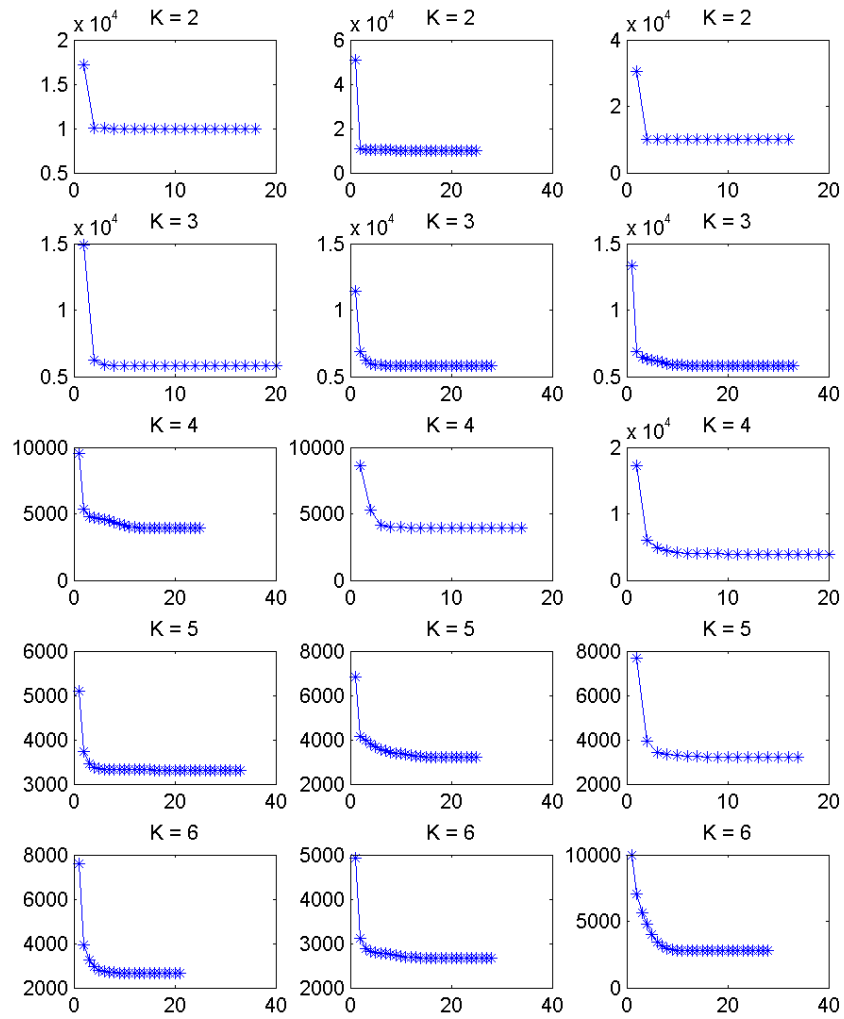


FIGURE 14 – DS3 - Distance de Mahalanobis

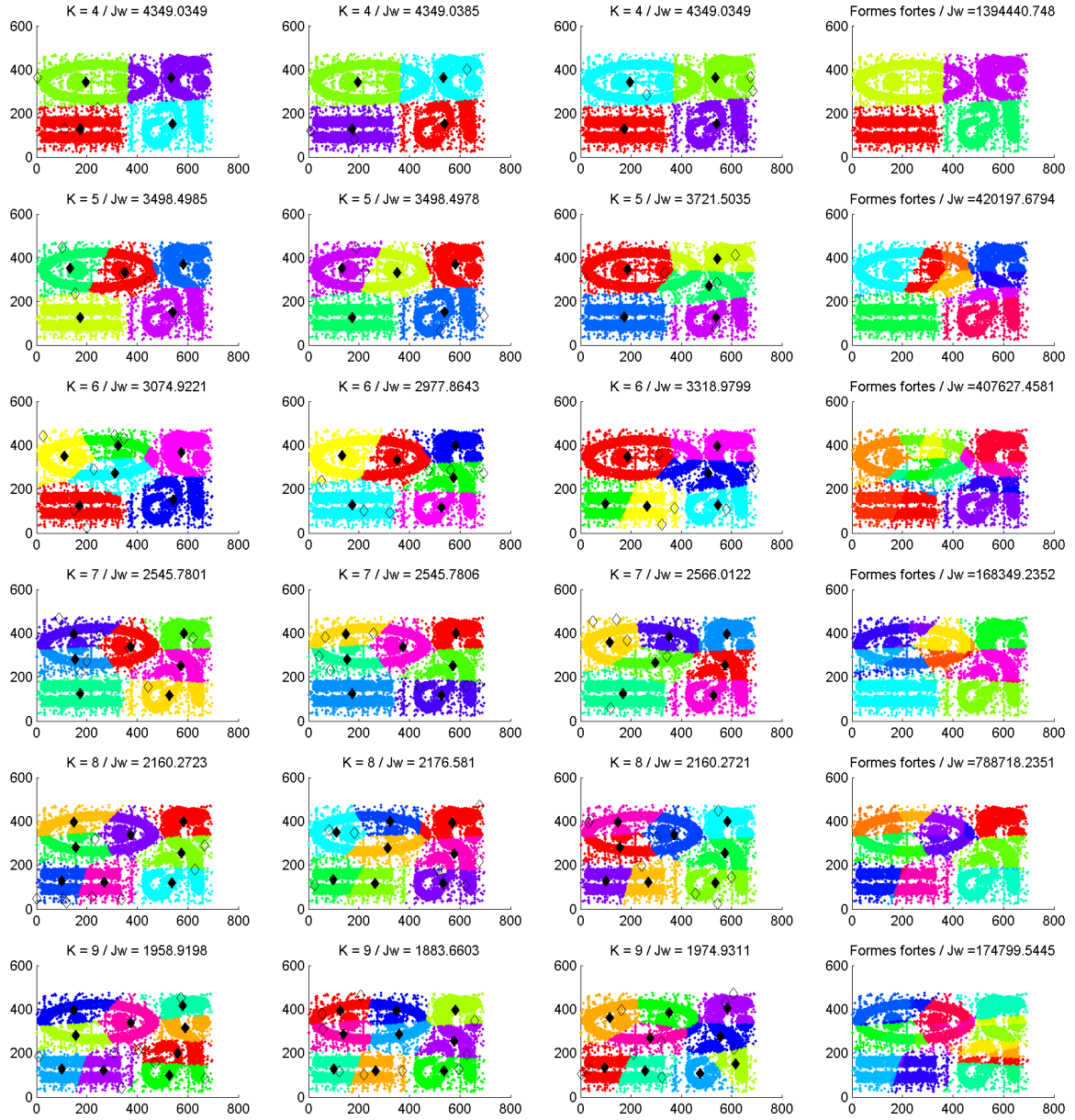


FIGURE 15 – DS4 - Distance de Mahalanobis

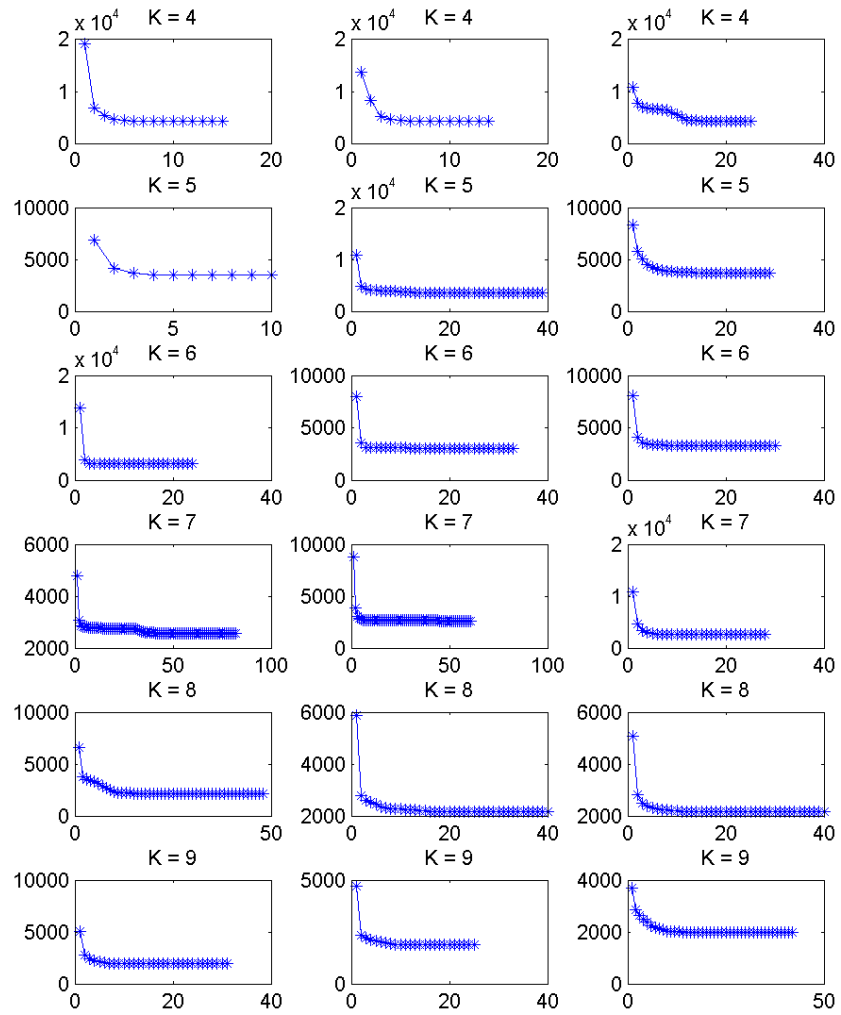


FIGURE 16 – DS4 - Distance de Mahalanobis

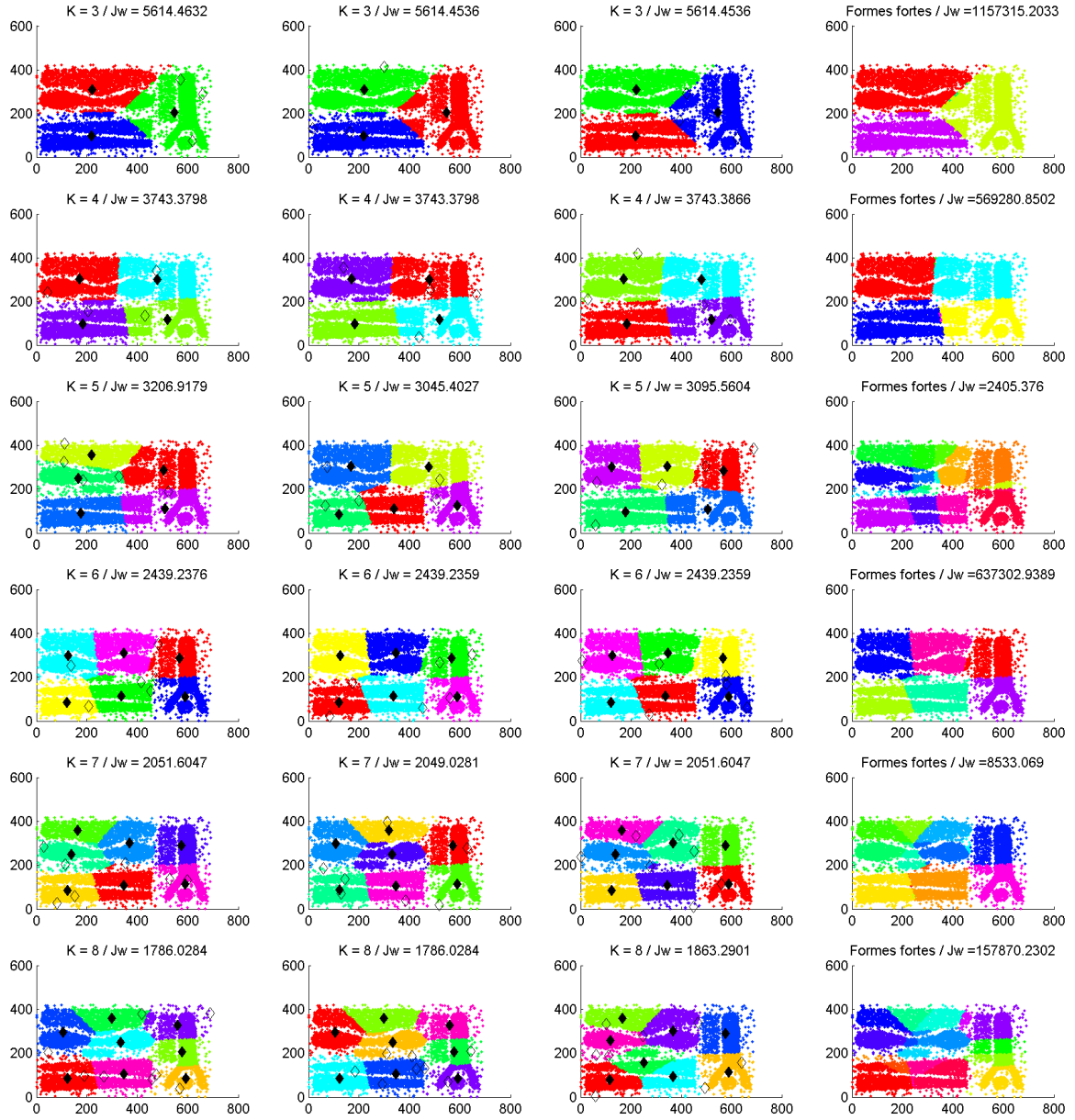


FIGURE 17 – DS5 - Distance de Mahalanobis

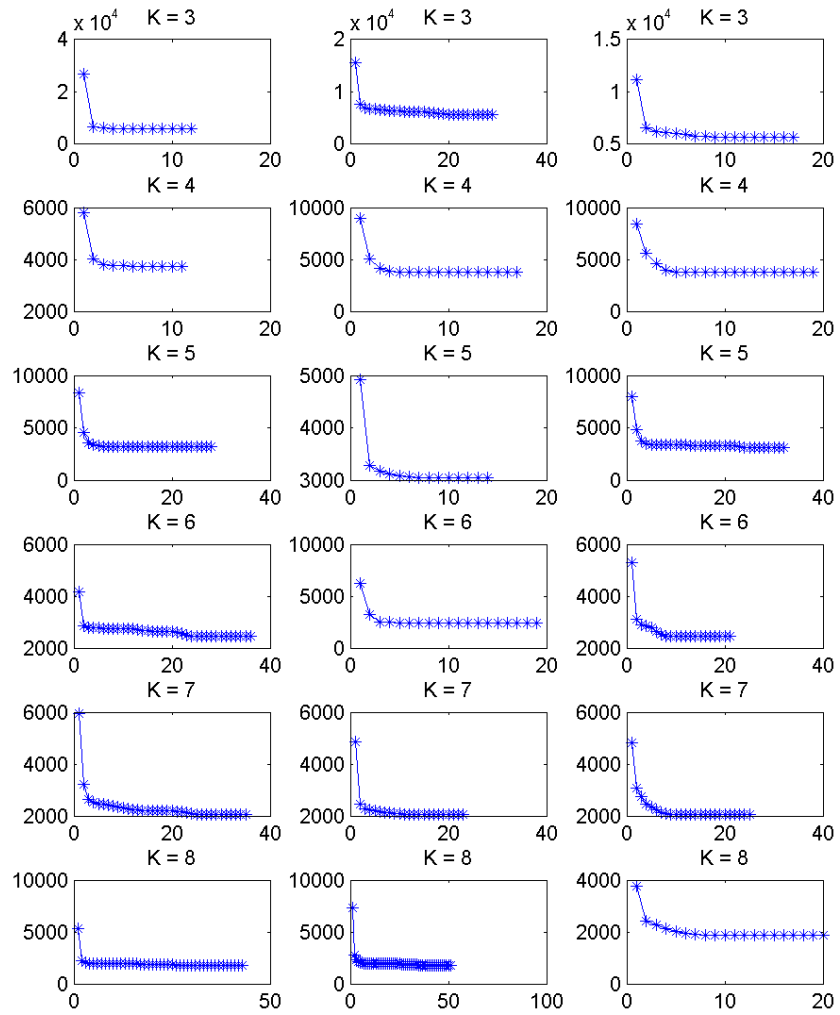


FIGURE 18 – DS5 - Distance de Mahalanobis

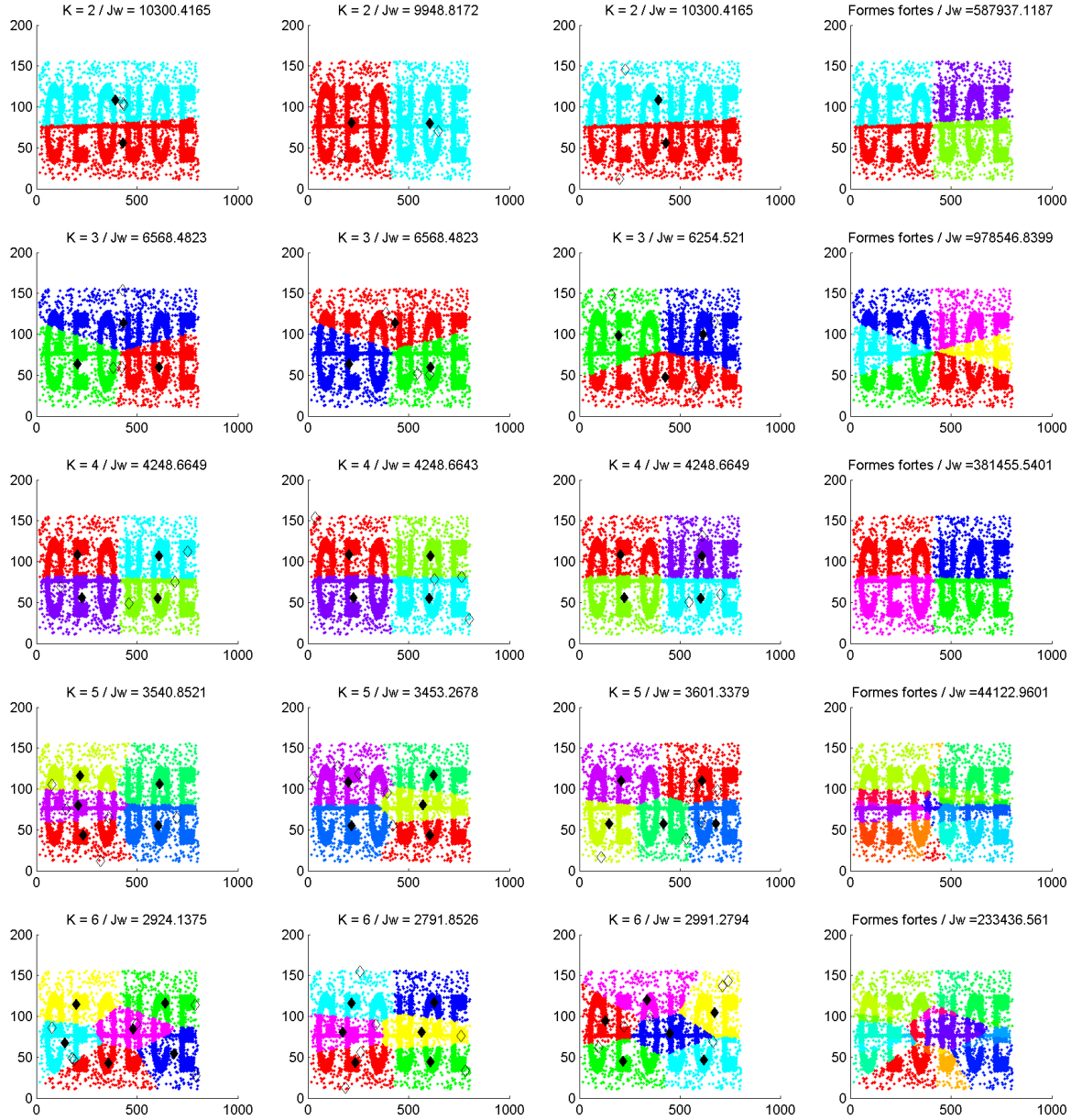


FIGURE 19 – George - Distance de Mahalanobis

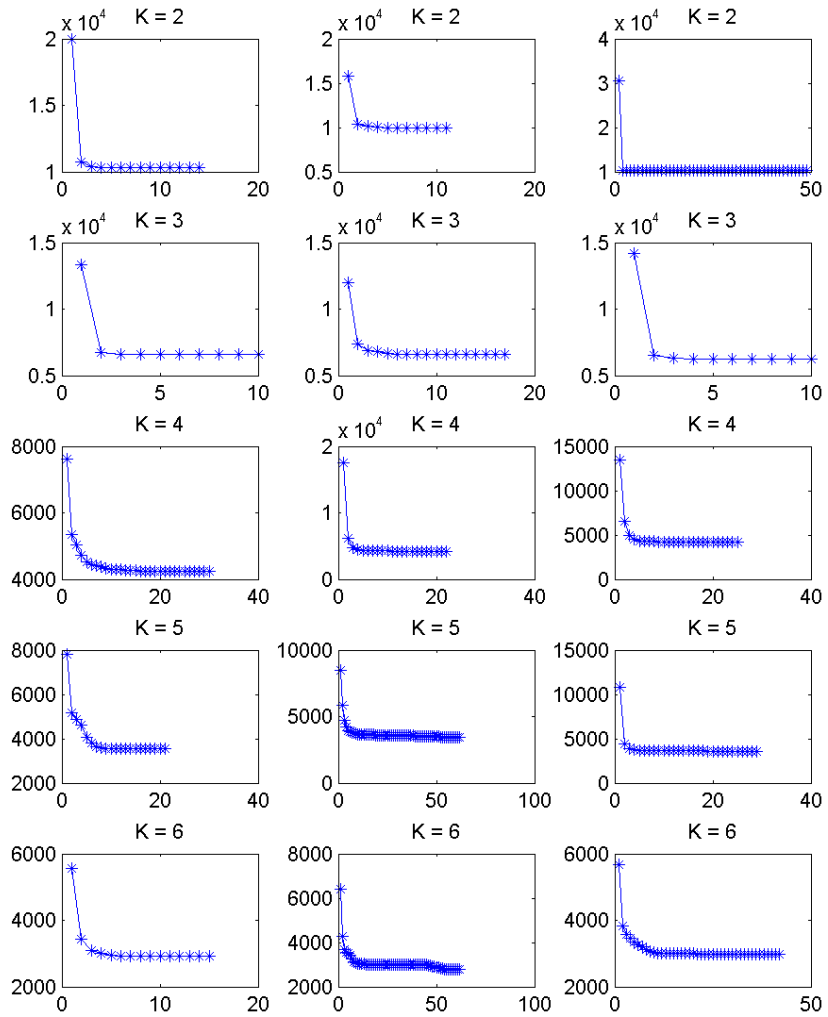


FIGURE 20 – George - Distance de Mahalanobis

Data Mining

TP5 - Clustering par mélange de gaussiennes

Rémi MUSSARD - Thomas ROBERT

1 Code

Afin de réaliser un clustering par la méthode du mélange de gaussiennes, on code l'algorithme EM. On retrouve dans ce code les deux étapes délimitées par des commentaires. Le code en commentaire est le code avant optimisation de plusieurs étapes de calcul.

```
1 %%
2 close all;
3 clc;
4 clear all;
5
6 load gauss3.mat;
7 %X = data;
8 X = load('george.dat');
9 X = mydownsampling(X, 7);
10
11 [N, d] = size(X);
12 K = 6;
13
14 figure();
15 hold on;
16
17 % Initialisation
18 mu = init_centres(X, K);
19 S = {};
20 for j = 1:K
21     S{j} = eye(d)*var(X(:,j));
22
23 end
24 pi = ones(K,1)*1/K;
25
26 mu_prec = mu + 1;
27 while abs(sum(mu-mu_prec)) > 1e-6
28     mu_prec = mu;
29
30     % Etape E
31
32     P = zeros(N,K);
33     for j = 1:K
34         P(:,j) = pi(j) * mvnpdf(X, mu(j,:), S{j});
35     end
36     P = P./repmat(sum(P,2), 1, K);
37
38     % Etape M
39
40     mu = (X'*P./repmat(sum(P),d,1))';
41     %mu = zeros(K, d);
42     %for j = 1:K
43     %     mu(j,:) = sum(repmat(P(:,j), 1, d).*X)/sum(P(:,j));
44     %end
45
46     pi = sum(P)/N;
47     %pi = zeros(K, 1);
48     %for j = 1:K
```



```

49  %    pi(j) = sum(P(:,j))/N;
50  %end
51
52  for j = 1:K
53      S{j} = (X - repmat(mu(j,:),N,1))'*diag(P(:,j))*(X - repmat(mu(j,:),N,1)) / sum(P(:,j));
54
55      %S{j} = zeros(d);
56      %for i = 1:N
57      %    S{j} = S{j} + P(i,j)*(X(i,:)-mu(j,:))'*(X(i,:)-mu(j,:));
58      %end
59      %S{j} = S{j} / sum(P(:,j));
60  end
61
62  plot(mu(:,1), mu(:,2), '+', 'Color',[0.5,0.5,0.5]);
63
64 end
65
66 % extraction de la liste des clusters
67 [~, liste] = max(P');
68 % affichage des clusters
69 show_clusters(X, liste);
70 % affichage des centres
71 scatter(mu(:,1), mu(:,2), 50, 'd', 'k', 'fill');

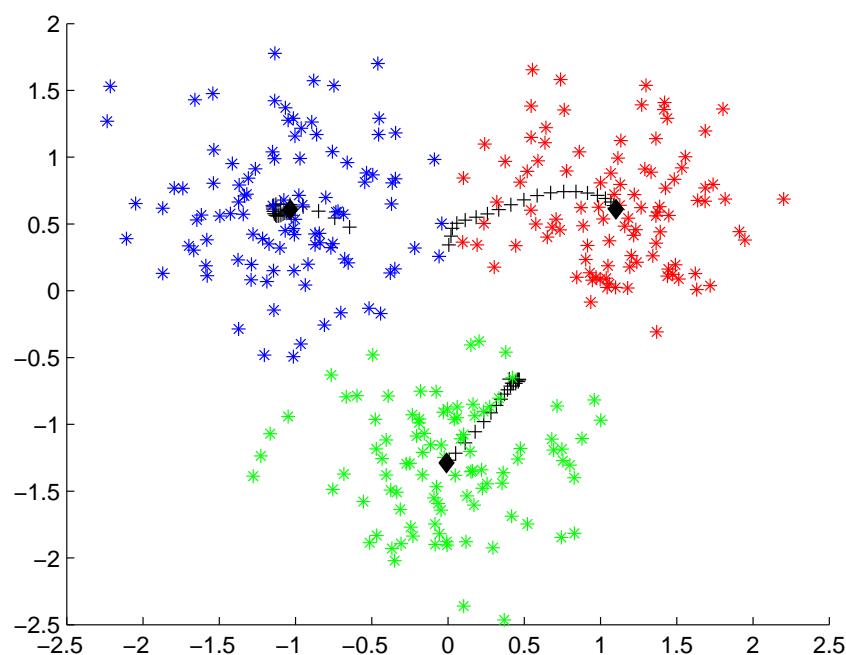
```

../TP5_MUSSARD_ROBERT.m

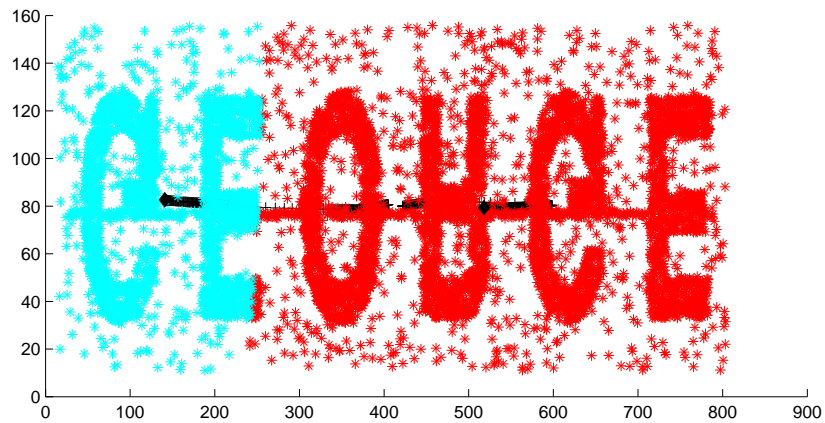
2 Résultats et commentaires

On voit que sur un jeu de données appropriés (gauss3.mat), la méthode fonctionne bien et permet de déterminer les centres des points répartis selon une loi normale.

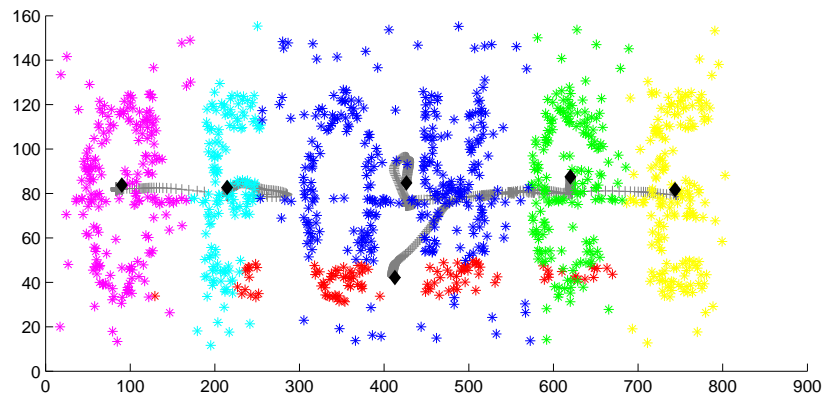
En revanche, sur un jeu de données où les points ne respectent pas une loi normale, comme par exemple george.dat où les données représentent des lettres, la méthode a beaucoup de mal à converger et ne donne pas un résultat satisfaisant.



Données gauss avec 3 clusters



Données george avec 2 clusters



Données george avec 6 clusters

3 Conclusion

Après avoir testé diverses méthodes de clustering, on a pu constater que les diverses méthodes fonctionnent globalement bien. Cependant, les résultats fournis par les diverses méthodes regroupent dans un clusters des points les plus proches, ce qui ne correspond pas toujours vraiment aux données fournies, qui contiennent plutôt des formes et non des groupes distincts de points.

La méthode CHAmélion présentée dans le cours semble pouvoir résoudre plus ou moins ce problème d'après les exemples fournis.