# Pattern Recognition Practical 5

Group 24:    Maikel Withagen (s1867733)    Steven Bosch (s1861948)

October 15, 2015

## Assignment 1    k-means clustering, quantization error, gap statistic

### 1

Using the code given in the Appendix(kmeans.m and runKMeans.m), we created the plots shown in figures 1, 2 and 3.
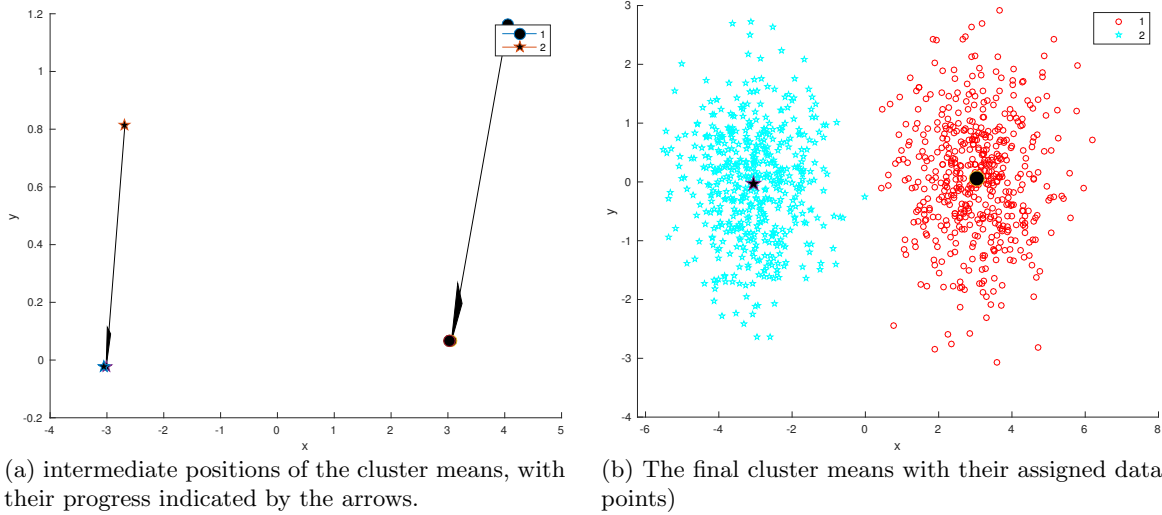
Figure 1: Results for k=2



(a) intermediate positions of the cluster means, with their progress indicated by the arrows.

(b) The final cluster means with their assigned data points)

Figure 2: Results for k=4



(a) intermediate positions of the cluster means, with their progress indicated by the arrows.

(b) The final cluster means with their assigned data points)

Figure 3: Results for k=8



(a) intermediate positions of the cluster means, with their progress indicated by the arrows.
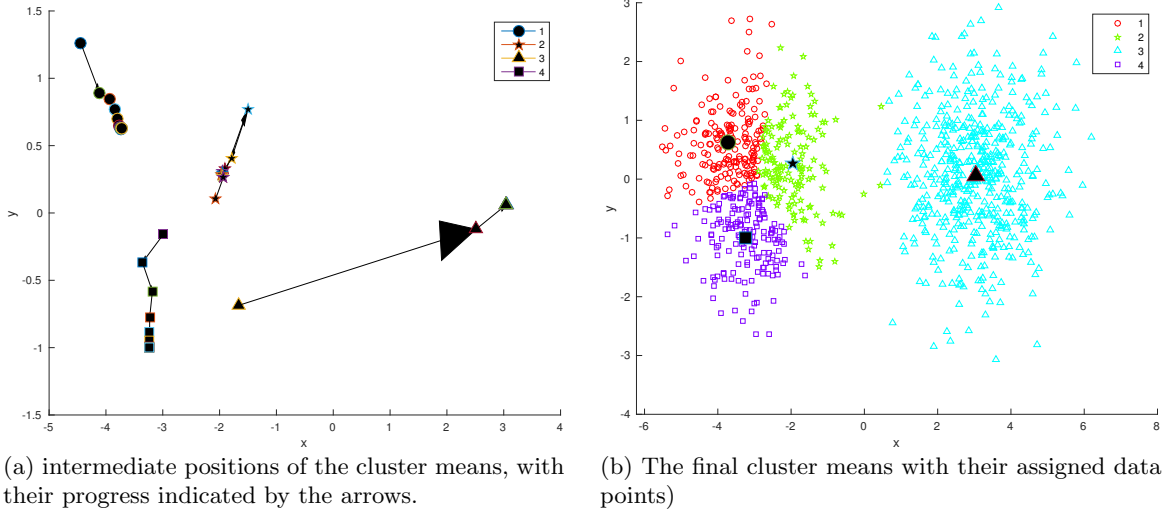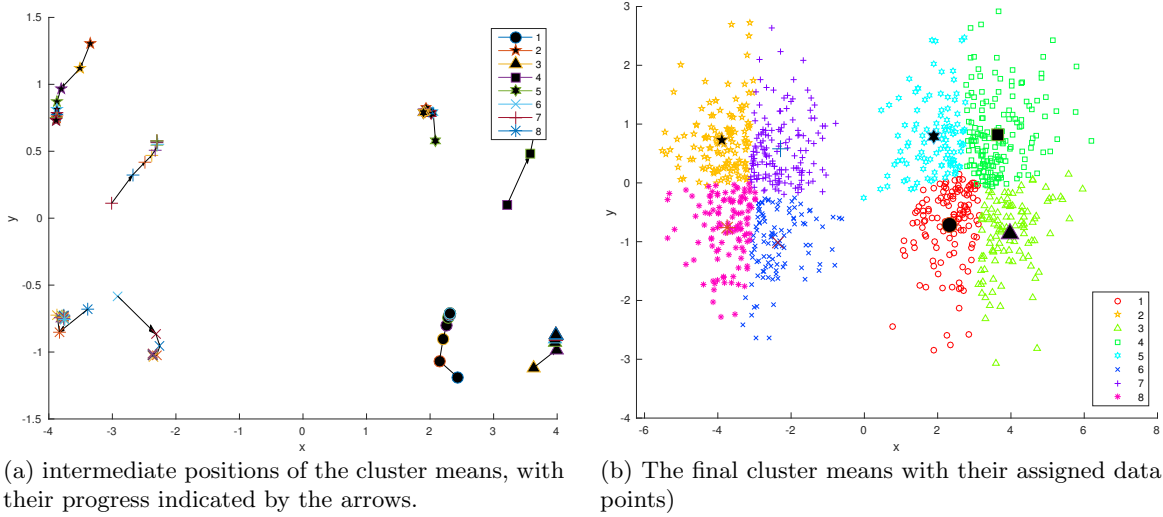
(b) The final cluster means with their assigned data points)

We we can clearly see that the data form two clusters. Therefore figure 1a shows the quickest convergence to the final cluster centers. Usually it takes about two epochs for the cluster centers to converge, as is shown in the figure. Figure 1b shows that these centers form in the places which the human eye observes to be the correct centers. When we choose k as 4, as shown in figure 2, we can see that, dependent on the initialization, sometimes one main cluster gets divided into three subclusters and the other remains one cluster, and sometimes the two clusters get separated into two clusters each. The number of epochs it takes to reach convergence is high compared to a run using $k = 2$. This can be explained by the fact that the data are not naturally separated into four clusters but in two, so the distances between the data points within a main cluster are small. This causes the algorithm to take longer to find a convergence, since the cluster centers switch often during the clustering process. Finally figure 3 shows the clustering for $k = 8$, which takes the longest amount of epochs to converge, because of the same reasoning. It separates both of the clusters into four subclusters.

## 2

Using the code given in the appendix (kmeans.m and runKMeans.m) we computed the quantization errors and D-function given in figure 4.

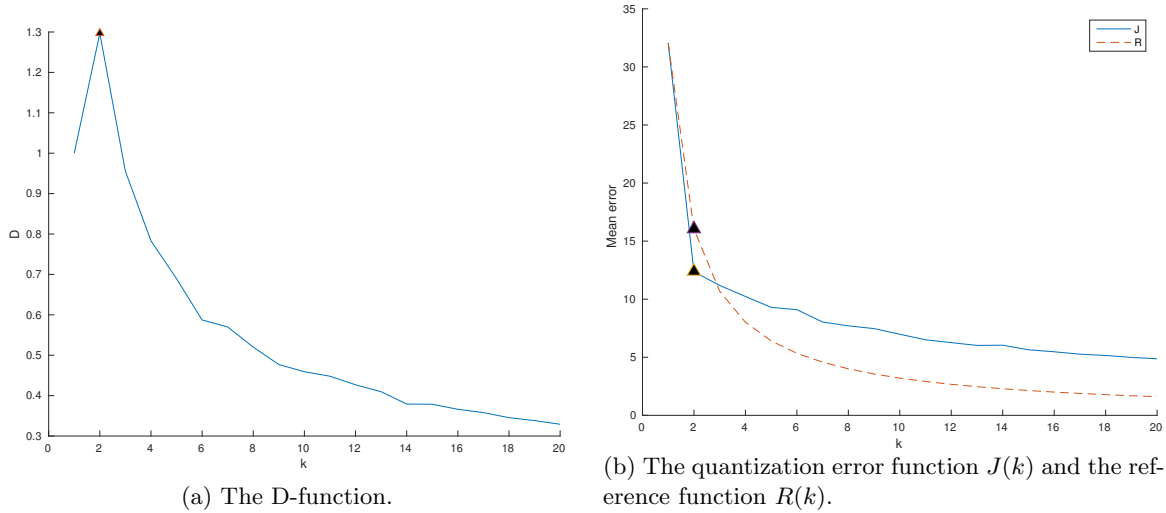Figure 4: Results for $kmax = 20$. The triangles give $k_{opt}$



(a) The D-function.



(b) The quantization error function $J(k)$ and the reference function $R(k)$.

Figure 4 shows the $D(k)$, $J(k)$ and $R(k)$ for a $kmax$ of 20. Here $k_{opt} =_k D(k)$ is the $k$ for which the difference between the error function and the reference function is the largest. Figure 4 shows that $k_{opt}$ can be found at a $k$ of 2, which was expected because it could clearly be seen by a human that the data are divided into two main clusters. $D(k)$ decreases as $k$ increases, clusters that have no natural division need to be divided by even more clusters, which causes them to shift a lot before finally reaching convergence.

# Assignment 2     Batch Neural gas vs k-means

# Appendix

../Code/kmeans.m

```
1  function [qError] = kmeans(dat, k, writeOutput)
2  % K-means clustering algorithm
3  close all;
4  shapes = 'op^shx+*dv<>.';
5
6  % Init the prototypes to a random point
7  prototypes = zeros(k,ndims(dat));
8  for i = 1:k
9      newPoint = dat(randi(length(dat)),1:2);
10     while (sum(pdist2(prototypes, newPoint) == 0) ~= 0)
11         newPoint = dat(randi(length(dat)),1:2);
12     end
13     prototypes(i,:) = newPoint;
14 end
15
16 %Init the first figure
17 figure(1)
18 hold on;
```

```matlab
19  xlabel('x');
20  ylabel('y');
21
22  for i = 1 : size(prototypes, 1)
23      plot(prototypes(i,1),prototypes(i,2),'Marker', shapes(i), 'MarkerSize', 10, '
            MarkerFaceColor', 'black')
24  end
25
26
27  % Perform k-means
28  loop = 1;
29  while(loop == 1)
30      loop = 0;
31
32      for point = 1 : length(dat)
33          dat(point,3) = find(pdist2(dat(point,1:2), prototypes) == min(pdist2(dat(point,1:2),
                prototypes)),1);
34      end
35
36      for prototype = 1 : size(prototypes, 1)
37          newMean = mean(dat(dat(:,3) == prototype,1:2));
38          if newMean ~= prototypes(prototype,:)
39              loop = 1;
40          end
41          plot_arrow( prototypes(prototype,1),   prototypes(prototype,2), newMean(:,1), newMean
                (:, 2));
42          prototypes(prototype,:) = newMean;
43          plot(newMean(1),newMean(2),'Marker', shapes(prototype), 'MarkerSize', 10, '
                MarkerFaceColor', 'black')
44      end
45
46
47  end
48
49  % Calculate the quantization error
50  qError = 0;
51  for i = 1 : size(prototypes, 1)
52      qError = qError + sum(pdist2(prototypes(i,:), dat(dat(:,3) == i,1:2)));
53  end
54
55  % More figure stuff
56  legend(strtrim(cellstr(num2str((1:k)'))'));
57  if writeOutput == 1
58      print(sprintf('../Report/Fig1_k%d', k), '-depsc');
59  end
60  figure(2)
61  hold on;
62  gscatter(dat(:,1),dat(:,2),dat(:,3),[],shapes, 5)
63
64
65  for i = 1 : size(prototypes, 1)
66      plot(prototypes(i,1),prototypes(i,2),'Marker', shapes(i), 'MarkerSize', 13, '
            MarkerFaceColor', 'black')
67  end
68
69  xlabel('x');
70  ylabel('y');
71  if writeOutput == 1
72      print(sprintf('../Report/Fig2_k%d', k), '-depsc');
73  end
```

../Code/runKMeans.m

```matlab
1  load('kmeans1.mat', 'kmeans1');
2
```

```matlab
 3  error= zeros(1,10);
 4  kmax = 20;
 5  J = zeros(1, kmax);
 6  R = zeros(1,kmax);
 7
 8  % Run for 1 to kmax clusters
 9  for k = 1 : kmax
10      k
11      % Run it 10 times for every cluster and calculate the mean error and
12      % reference
13      for i = 1:10
14          error(i) = kmeans(kmeans1,k, 0);
15      end
16      J(k) = mean(error)/10;
17      R(k) = J(1) * k^(-2/ndims(kmeans1));
18  end
19
20  D = R ./ J;
21
22  % Plot D
23  [maxVal maxInd] = max(D);
24  figure(3)
25  hold on;
26  plot(D);
27  plot(maxInd, maxVal,'Marker', '^', 'MarkerSize', 6, 'MarkerFaceColor', 'black')
28  xlabel('k');
29  ylabel('D');
30  print(sprintf('../Report/Fig3'), '-depsc');
31
32  % Plot J and R
33  figure (4)
34  hold on ;
35  plot(J);
36  plot(R, '--');
37  plot(maxInd, J(maxInd),'Marker', '^', 'MarkerSize', 10, 'MarkerFaceColor', 'black')
38  plot(maxInd, R(maxInd),'Marker', '^', 'MarkerSize', 10, 'MarkerFaceColor', 'black')
39  xlabel('k');
40  ylabel('Mean error');
41  legend('J', 'R');
42  print(sprintf('../Report/Fig4'), '-depsc');
```