

# 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

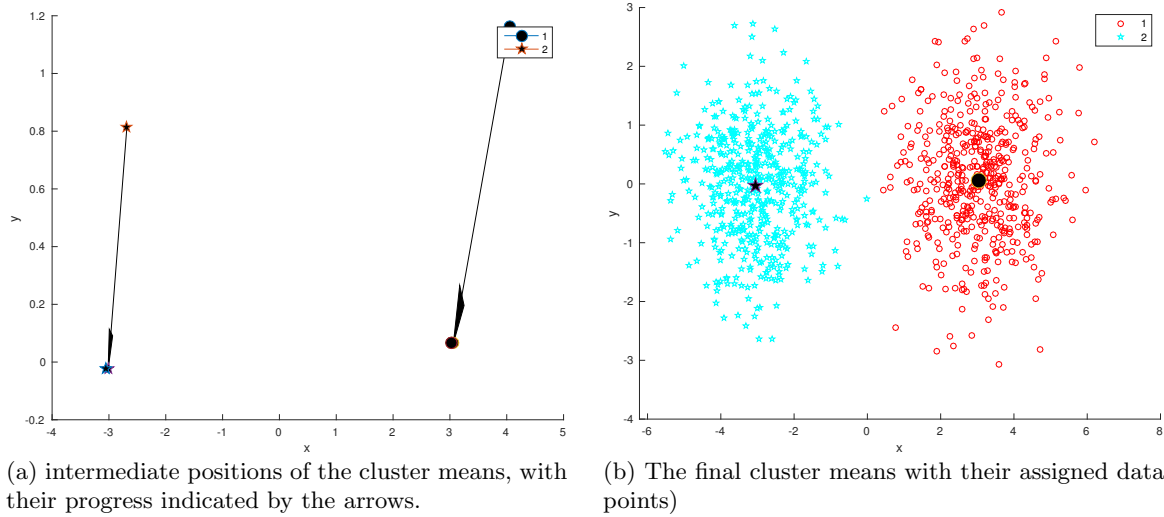


Figure 2: Results for  $k=4$

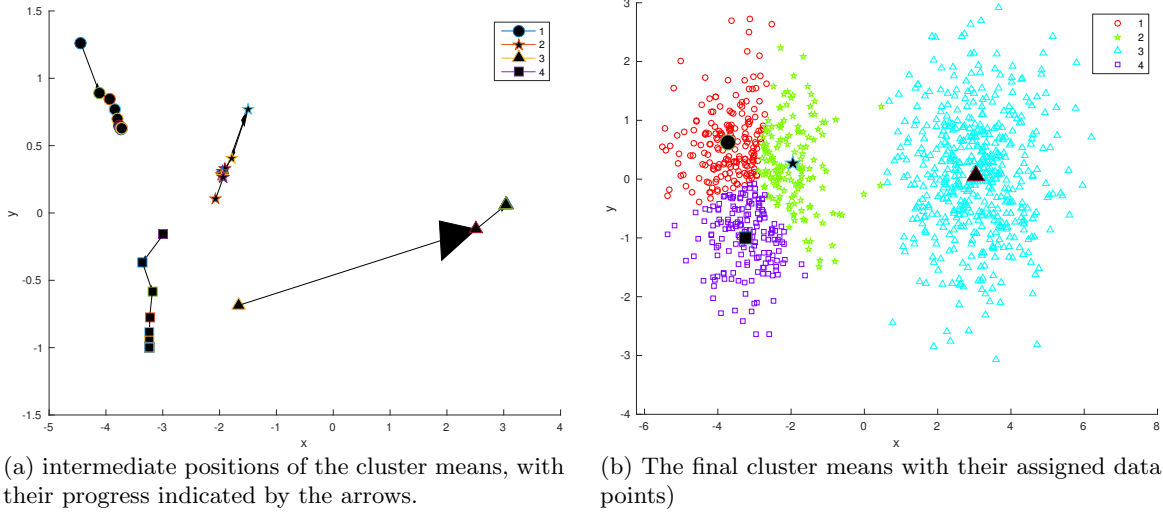
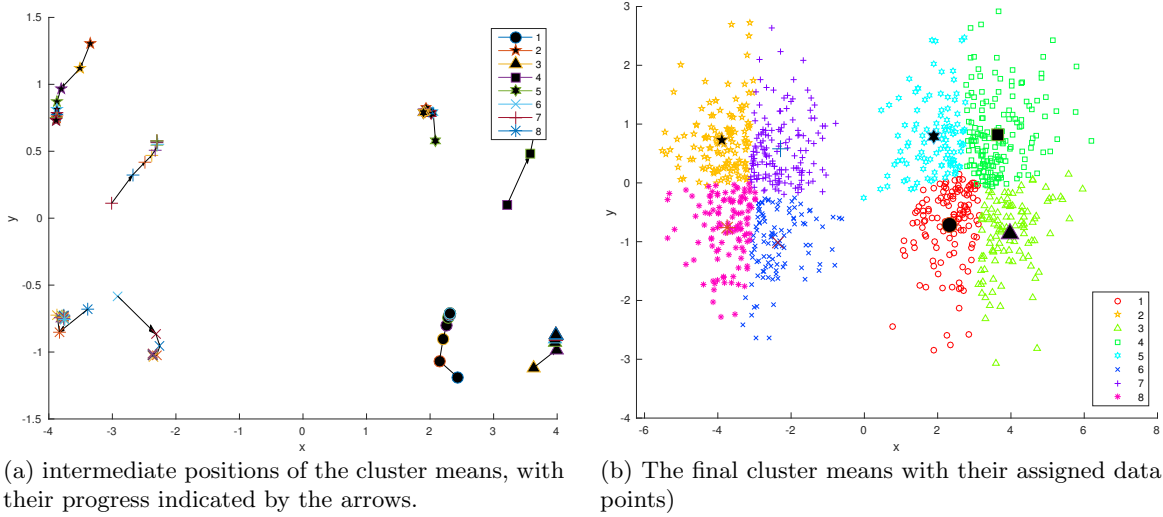


Figure 3: Results for  $k=8$



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  $k_{opt}$

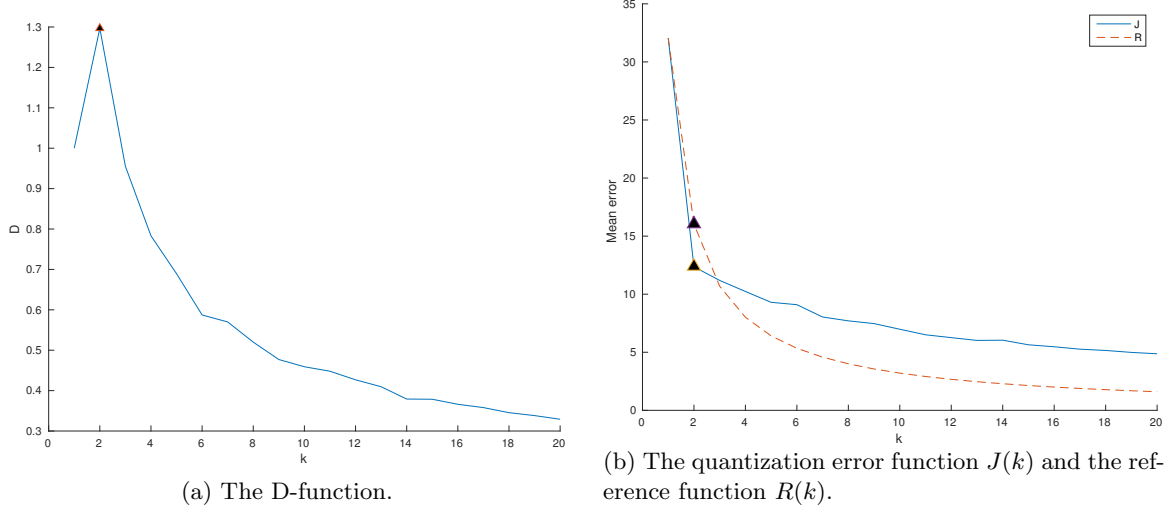


Figure 4 shows the  $D(k)$ ,  $J(k)$  and  $R(k)$  for a  $kmax$  of 20. Here  $k_{opt} = \operatorname{argmax}_k D(k)$  is the  $k$  for which the difference between the error function and the reference function is the largest. Figure 4 shows that  $k_{opt}$  can be found at a  $k$  of 2, which was 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.

## 3

### a

See the part in `kmeans.m` which is commented “Kmeans plusplus initialization”.

### b

We computed the following means and standard deviations given in 1.

Table 1: Mean and standard deviation of quantization error for  $k = 100$ .

Algorithm	mean q-error	sd q-error
<i>kmeans</i> ++	10.0310	0.3502
<i>kmeans</i>	10.6001	0.4268

### c

We calculated the following p-value using a an unpaired one-tailed two-sample t-test (see `runKMeans.m` at the bottom):  $2.3816e^{-5}$ . This p-value is significant ( $< 0.05$ ), which means that we have enough certainty to reject the null hypothesis that there is no difference in performance between the two algorithmes. Hence, we assume a difference between the performance of the two algorithmes (*kmeans* ++ performs better).

## Assignment 2 Batch Neural gas vs k-means

### Appendix

../Code/kmeans.m

```
1 function [qError] = kmeans(dat, k, plusplus, writeOutput)
2 % K-means clustering algorithm
3 close all;
4 shapes = 'op^shx+*dv<>.';
5
6 if plusplus == 0
7     % Init the prototypes to a random point
8     prototypes = zeros(k, ndims(dat));
9     for i = 1:k
10         newPoint = dat(randi(length(dat)), 1:2);
11         while (sum(pdist2(prototypes, newPoint) == 0) ~= 0)
12             newPoint = dat(randi(length(dat)), 1:2);
13         end
14         prototypes(i, :) = newPoint;
15     end
16 else
17     % Kmeans plusplus initialization
18     prototypes = zeros(k, ndims(dat));
19     newPoint = dat(randi(length(dat)), 1:2);
20     distances = pdist2(newPoint, dat(:, 1:2));
21     prototypes(1, :) = newPoint;
22
23     for i = 2:k
24         probabilities = distances.^2;
25         % Choose a random number in the range of sum(probabilities)
26         newPointIndex = find(cumsum(probabilities) > rand*sum(probabilities), 1);
27         newPoint = dat(newPointIndex, :);
28         prototypes(i, :) = newPoint;
29
30         % Calculate the new distances (select min value)
31         distances = min(distances, pdist2(newPoint, dat(:, 1:2)));
32     end
33 end
34
35 % Init the first figure
36
37 % figure(1)
38 % hold on;
39 % xlabel('x');
40 % ylabel('y');
41 % for i = 1 : size(prototypes, 1)
42 %     plot(prototypes(i, 1), prototypes(i, 2), 'Marker', shapes(i), 'MarkerSize', 10, '
43 %         MarkerFaceColor', 'black')
44 % end
45
46 % Perform k-means
47 % loopCounter = 0;
48 loop = 1;
49 while(loop == 1)
50     loop = 0;
51     % Uncomment to show loop counter :D
52     % loopCounter = loopCounter + 1
53
54     for point = 1 : length(dat)
55         dat(point, 3) = find(pdist2(dat(point, 1:2), prototypes) == min(pdist2(dat(point, 1:2),
56             prototypes)), 1);
57     end
58 end
```

```

56
57     for prototype = 1 : size(prototypes, 1)
58         newMean = mean(dat(dat(:,3) == prototype,1:2));
59         if newMean ~= prototypes(prototype,:)
60             loop = 1;
61         end
62
63         plot_arrow( prototypes(prototype,1), prototypes(prototype,2), newMean(:,1), newMean
64                     (:, 2));
65         prototypes(prototype,:) = newMean;
66         plot(newMean(1),newMean(2),'Marker', shapes(prototype), 'MarkerSize', 10, '
67             MarkerFaceColor', 'black')
68
69     end
70
71     % Calculate the quantization error
72     qError = 0;
73     for i = 1 : size(prototypes, 1)
74         qError = qError + sum(pdist2(prototypes(i,:), dat(dat(:,3) == i,1:2)));
75     end
76
77     % More figure stuff
78     % legend(num2str(1:k))
79     % if writeOutput == 1
80     %     print(sprintf(' ../Report/Fig1_k%d', k), '-depsc');
81     % end
82     % figure(2)
83     % hold on;
84     % gscatter(dat(:,1), dat(:,2), dat(:,3), [], shapes, 5)
85     %
86
87     %
88     % for i = 1 : size(prototypes, 1)
89     %     plot(prototypes(i,1), prototypes(i,2), 'Marker', shapes(i), 'MarkerSize', 13, '
90         %     MarkerFaceColor', 'black')
91     % end
92     %
93     % xlabel('x');
94     % ylabel('y');
95     if writeOutput == 1
96         print(sprintf(' ../Report/Fig2_k%d', k), '-depsc');
97     end

```

../Code/runKMeans.m

```

1  load('kmeans1.mat', 'kmeans1');
2
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, 2, 1, 0);
15     end
16     J(k) = mean(error)/10;
17     R(k) = J(1) * k^(-2/ndims(kmeans1));
18 end

```

```

19 D = R ./ J;
20
21
22 % Plot D
23 [maxVal maxInd] = max(D);
24 figure(3)
25 hold on;
26 plot(D);
27
28 plot(maxInd, maxVal, 'Marker', '^', 'MarkerSize', 6, 'MarkerFaceColor', 'black')
29 xlabel('k');
30 ylabel('D');
31 print(sprintf('.. / Report / Fig3 '), '-depsc');
32
33 % Plot J and R
34 figure(4)
35 hold on;
36 plot(J);
37 plot(R, '—');
38 plot(maxInd, J(maxInd), 'Marker', '^', 'MarkerSize', 10, 'MarkerFaceColor', 'black')
39 plot(maxInd, R(maxInd), 'Marker', '^', 'MarkerSize', 10, 'MarkerFaceColor', 'black')
40 xlabel('k');
41 ylabel('Mean error');
42 legend('J', 'R');
43 print(sprintf('.. / Report / Fig4 '), '-depsc');
44
45 % Perform the kmeans++ test
46 k = 100;
47 nRuns = 20;
48 error_without = zeros(1, nRuns);
49 error_with = zeros(1, nRuns);
50 for i = 1:nRuns
51     datestr(now)
52     i
53     error_without(i) = kmeans(checkerboard, k, 0, 0);
54     error_with(i) = kmeans(checkerboard, k, 1, 0);
55 end
56 error_without = error_without/nRuns;
57 error_with = error_with/nRuns;
58
59 mean(error_without)
60 std(error_without)
61 mean(error_with)
62 std(error_with)
63
64 % using an unpaired one-tailed two-sample t-test
65 ttest2(error_without, error_with)

```