

# COMPM012 - Coursework 3

Esben A. Sørig

November 29, 2014

## 1 Practical

### 1.1 k-means implementation

```
1 function [clusterings, centers] = mykmeans(X, k)
2     % Randomly initialize centers of clusters
3     c = datasample(X, k, 'Replace', false);
4     r = repmat(0, size(X, 1), k);
5     oldr = 1; % something that is not equal to r initially
6
7     dist = @(x, y) norm(x-y);
8
9     % Loop as long as the clustering is changing
10    while ~isequal(r, oldr)
11        oldr = r;
12
13        % Assign points to clusters
14        for i = 1:size(X,1)
15            cluster = 1;
16            for j = 1:k
17                if dist(X(i, :), c(j, :)) < dist(X(i, :), c(
18                    cluster, :))
19                    cluster = j;
20                end
21            end
22            r(i, :) = [repmat(0, 1, cluster-1) 1 repmat(0, 1, k-
23                cluster)];
24        end
25
26        % Update center positions
27        for i = 1:k
28            npoints = 0;
29            c(i, :) = 0;
```

```

30         c(i, :) = c(i, :) + r(j, i)*X(j, :);
31         npoints = npoints + r(j, i);
32     end
33     c(i, :) = c(i, :)/npoints;
34 end
35 end
36
37 centers = c;
38 % Output formatting (vector with cluster index for each row
39 % in X)
39 clustering = repmat(0, size(X,1), 1);
40 for i = 1:size(X,1)
41     for j = 1:k
42         if r(i, j) == 1
43             clustering(i) = j;
44             break;
45         end
46     end
47 end
48
49 clusterings = clustering;
50 end

```

## 1.2 k-means test on data generated from three gaussians

Firstly, we generate the data and run the k-means algorithm on it.

```

1 % Generate data
2 data = genData2;
3
4 % Put data for each cluster in a separate list
5 cluster1 = [];
6 cluster2 = [];
7 cluster3 = [];
8 [clusterings, centers] = mykmeans(data, 3);
9
10 for i = 1:size(data, 1)
11     if clusterings(i) == 1
12         cluster1 = [cluster1; data(i, :)];
13     end
14     if clusterings(i) == 2
15         cluster2 = [cluster2; data(i, :)];
16     end
17     if clusterings(i) == 3
18         cluster3 = [cluster3; data(i, :)];
19     end
20 end

```

We can now plot the clusters the algorithm has found

```

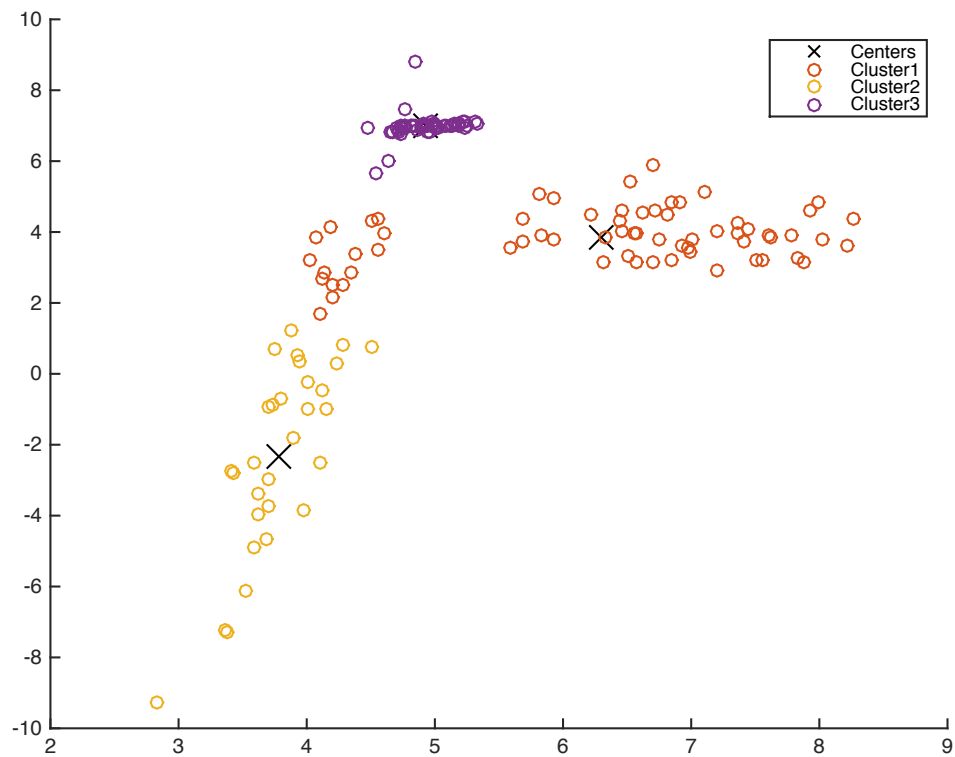
1 figure

```

```

2 hold on
3
4 scatter(centers(:, 1), centers(:, 2), 200, 'xblack');
5 scatter(cluster1(:, 1), cluster1(:, 2));
6 scatter(cluster2(:, 1), cluster2(:, 2));
7 scatter(cluster3(:, 1), cluster3(:, 2));
8 legend('Centers', 'Cluster1', 'Cluster2', 'Cluster3');
9 hold off

```

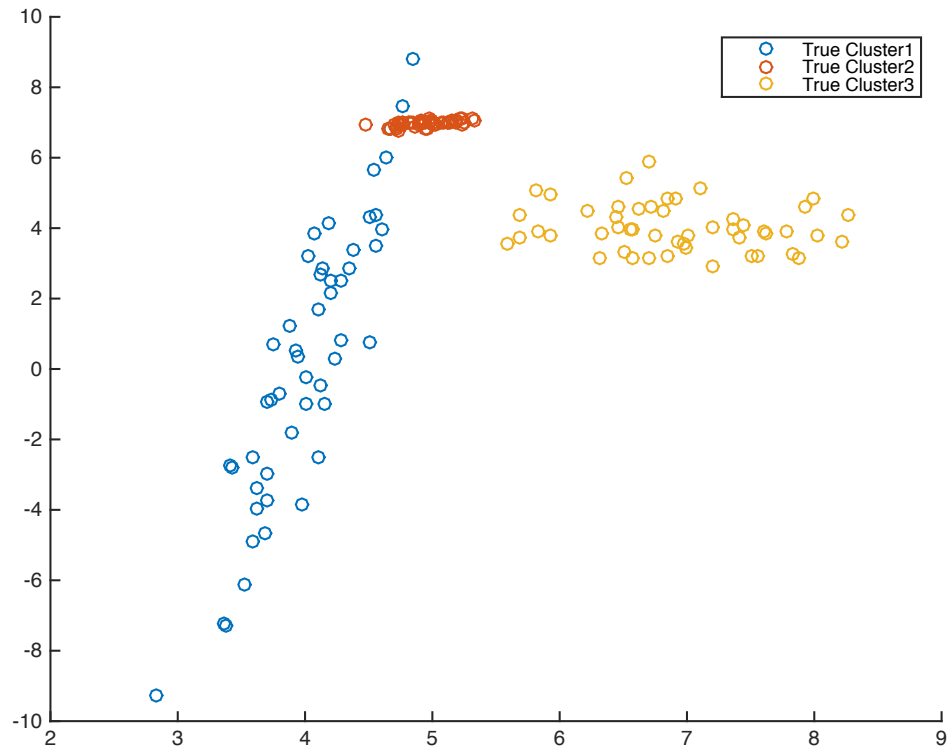


Comparing this to the true classifications

```

1 figure
2 hold on
3 scatter(data(1:50, 1), data(1:50, 2));
4 scatter(data(51:100, 1), data(51:100, 2));
5 scatter(data(101:150, 1), data(101:150, 2));
6 legend('True Cluster1', 'True Cluster2', 'True Cluster3');
7 hold off

```



We clearly get an intuition of how the k-means algorithm is performing on the data. Please see appendix 4.1 for the code that creates the animation of the convergence to a solution for this dataset.

We measure the mean error and standard deviation of the error of the clustering over a 100 runs of the algorithm.

```

1 %% Error measurements
2 errors = [];
3 for j = 1:100
4     [clusterings, centers] = mykmeans(data, 3);
5
6     % keep track of the classifications of each true cluster
7     firstcluster = [0,0,0];
8     secondcluster = [0,0,0];
9     thirdcluster = [0,0,0];
10    for i = 1:50
11        firstcluster(clusterings(i)) = 1 + firstcluster(
clusterings(i));
12    end
13    for i = 51:100

```

```

14         secondcluster(clusterings(i)) = 1 + secondcluster(
clusterings(i));
15     end
16     for i = 101:150
17         thirdcluster(clusterings(i)) = 1 + thirdcluster(
clusterings(i));
18     end
19     % We assume that the mode of the classifications of a true
cluster is the class of the true cluster. Therefore, the
error of a cluster is the frequency of other classifications
appearing for that cluster.
20     firstcluster = sort(firstcluster);
21     secondcluster = sort(secondcluster);
22     thirdcluster = sort(thirdcluster);
23     misclassifications = firstcluster(1) + firstcluster(2) +
secondcluster(1) + secondcluster(2) + thirdcluster(1) +
thirdcluster(2);
24     errors = [errors misclassifications/150];
25 end
26
27 meanerror = mean(errors)
28 stddeviationerror = std(errors)

```

Which returns:

```

1 meanerror =
2     0.1438
3 stddeviationerror =
4     0.0155

```

## 2 Questions

### 2.1 Dataset with local minima

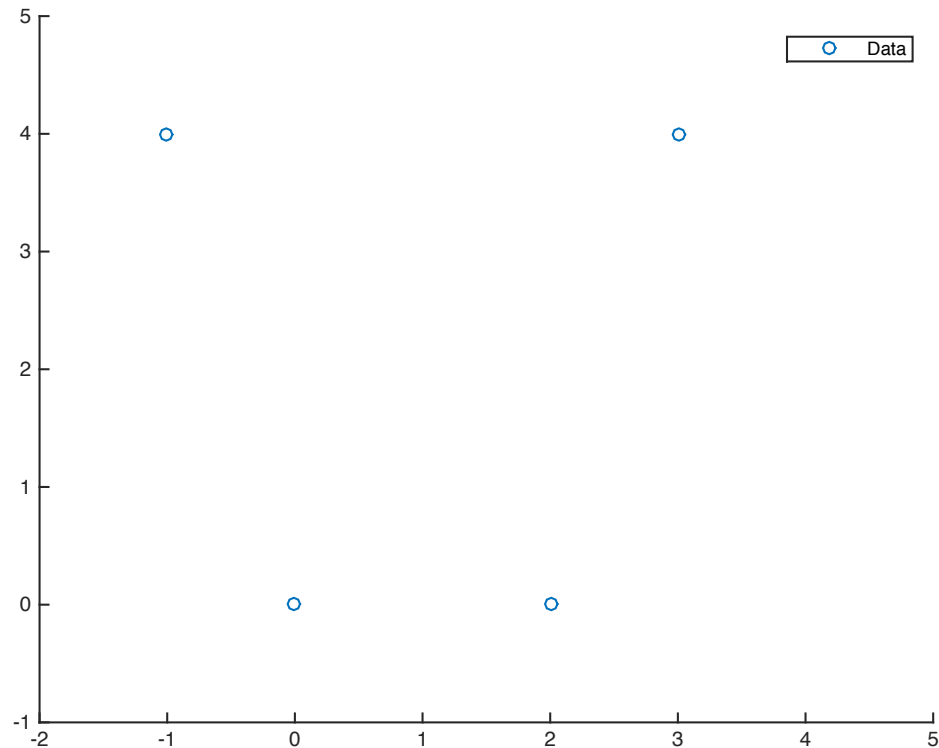
We can easily construct a dataset on which our k-means algorithm has local minima. An example of such a dataset is the set  $\{(0, 0), (2, 0), (-1, 4), (3, 4)\}$ . This is clear if we plot the data.

```

1 X = [0, 0; 2, 0; -1, 4; 3, 4];
2
3 hold on
4 axis([-2 5 -1 5]);
5 scatter(X(:, 1), X(:, 2));

```

Which gives the plot:



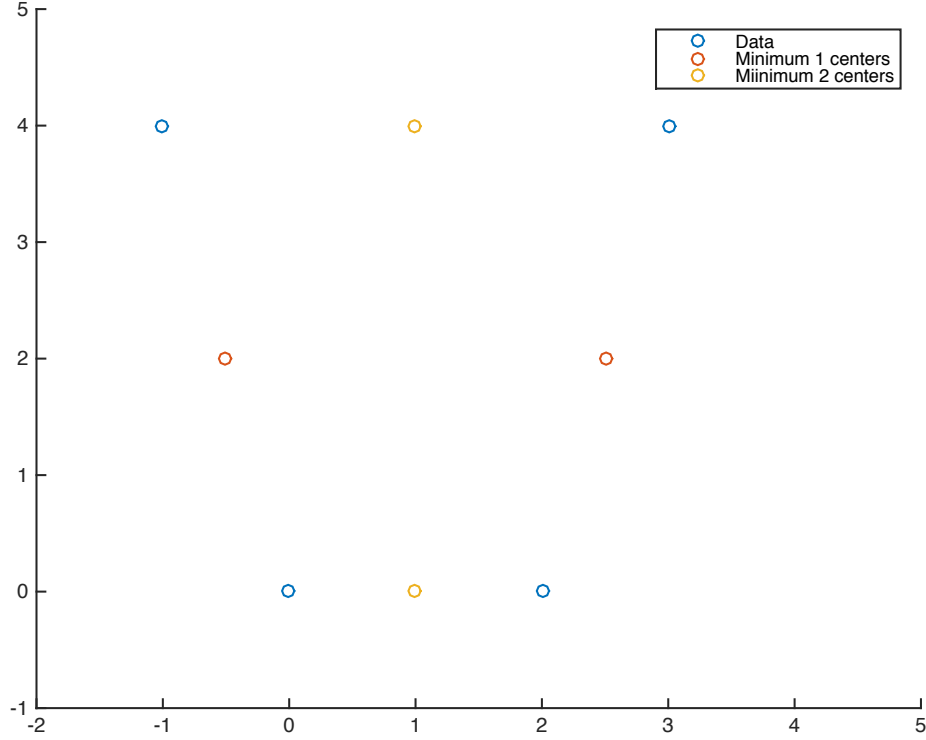
Running k-means on this dataset will converge in one of two local minima. We run the algorithm twice and plot the centers of the clusters: (in order to get two different results, this code may have to be run several times)

```

1 [clustering1, centers1] = mykmeans(X, 2);
2 [clustering2, centers2] = mykmeans(X, 2);
3 legend('Data');
4
5 scatter(centers1(:, 1), centers1(:, 2));
6 scatter(centers2(:, 1), centers2(:, 2));
7
8 legend('Data', 'Minimum 1 centers', 'Minimum 2 centers');
9 hold off

```

Which produces:



Here we clearly see the two local minima of the convergence on the four data points.

## 2.2 Argument that the centroid is the minimizer of the sum of squared distances

We can minimize the summed squared error

$$SSE = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \mathbf{c}_i\|^2 = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \mathbf{c}_i)^2 \quad (1)$$

for the  $k^{th}$  centroid by setting the derivative with respect to the  $k^{th}$  centroid equal to zero.

$$\begin{aligned}
\frac{\delta}{\delta \mathbf{c}_k} \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \mathbf{c}_i)^2 &= \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \frac{\delta}{\delta \mathbf{c}_k} (\mathbf{x} - \mathbf{c}_i)^2 \\
&= \sum_{\mathbf{x} \in C_k} 2(\mathbf{x} - \mathbf{c}_k) = 0 \\
\Rightarrow \sum_{\mathbf{x} \in C_k} \mathbf{c}_k &= \sum_{\mathbf{x} \in C_k} \mathbf{x} \\
\Rightarrow \mathbf{c}_k &= \frac{1}{\sum_{\mathbf{x} \in C_k} 1} \sum_{\mathbf{x} \in C_k} \mathbf{x}
\end{aligned} \tag{2}$$

We see that the minimizer of  $\mathbf{c}_k$  is equivalent to the centroid of the cluster.

### 2.3 Proof of convergence in finite amount of steps

We know that the k-means algorithm converges but we do not know if it does so in finitely or infinitely many steps. Here we show that the k-means algorithm does indeed converge in finitely many steps.

**Proof:** The k-means algorithm stops once there is no longer any improvement in the clustering. It will therefore never reach the same clustering twice. Each possible k-clustering of the input consists of k subsets. The number of subsets of the input is a finite number and the number of k subsets of the input is bounded above by this. Consequently, there is a finite number of possible clusterings. Hence the k-means algorithm will terminate after a finite number of iterations.

## 3 Extension

### 3.1 k-means segmentation

The objective of the algorithm is to find the segmentation  $\{i_1, \dots, i_{k-1}\}$  of a sequence of points  $\{\mathbf{x}_1, \dots, \mathbf{x}_l\} \in \mathbb{R}^n$  that is the global minimum of the following optimisation problem:

$$\operatorname{argmin}_{i_1, \dots, i_{k-1}; \mathbf{c}_1, \dots, \mathbf{c}_k} \sum_{j=1}^k \sum_{p=i_{j-1}+1}^{i_j} \|\mathbf{x}_p - \mathbf{c}_j\|^2 \tag{3}$$



The algorithm must run in polynomial time in  $l$ ,  $n$ , and  $k$ . We start by defining the error function as it will be useful when searching for minima of in the error.

```

1 function e = error(delimiters, X)
2     e = 0;
3     bounddelimiters = [0 delimiters size(X, 1)];
4     for j = 2:size(bounddelimiters, 2)
5         indices = (1+bounddelimiters(j-1)):bounddelimiters(j);
6         segment = X(indices, :);
7         centroid = sum(segment)/size(segment, 1);
8         for k = 1:size(segment, 1)
9             e = e + norm(segment(k, :)-centroid)^2;
10        end
11    end
12 end

```

We can now proceed to design the segmentation algorithm. Consider the set of global minimum delimiters  $\mathbf{i}^* = \{i_1^*, \dots, i_{k-1}^*\}$ . Observe that each  $i_j^*$  is the delimiter that given  $\mathbf{i}^* \setminus \{i_j^*\}$  minimises the error of  $\mathbf{i}^*$ . That is, each  $i_j^*$  is the delimiter that minimises the error when we add it to the set of the other delimiters. We can use this observation to construct a polynomial time algorithm.

The idea is to start with an empty set of optimal delimiters, i.e.  $\mathbf{i}_0^* = \emptyset$  where  $\mathbf{i}_n^*$  denotes the set of  $n$  optimal delimiters (which produces  $n + 1$  segments). From our observation it follows that  $\mathbf{i}_n^* = \mathbf{i}_{n-1}^* \cup \{i_n^*\}$  where  $i_n^*$  is the delimiter that minimises the error of  $\mathbf{i}_{n-1}^* \cup \{i_n\}$ . This easily translates to a recursive algorithm:

```

1 function delimiters = segmentation(X, k)
2     delimitercount = k-1;
3     if delimitercount <= 0
4         delimiters = [];
5         return;
6     end
7
8     delimiters = segmentation(X, k-1);
9     mindelimiters = [];
10    for i = 1:(size(X,1) - 1)
11        newdelimiters = sort([delimiters i]);
12        if error(newdelimiters, X) < error(mindelimiters, X)
13            mindelimiters = newdelimiters;
14        end
15    end
16    delimiters = mindelimiters;
17 end

```

**Informal argument of time complexity:** The running time is polynomial in  $l$ ,  $n$ , and  $k$ . The recursive call is performed  $k - 1$  times. At each recursive call, a loop iterates  $l$  times, at each iteration calling the error function. The error function has two nested loops iterating over the  $k$  segments delimiters and each of the, on average,  $l/k$  points in each segment. It is assumed that the summation function used in the error function has a complexity of  $O(l \times n)$ . Without going in further detail it seems that the running time is around  $O(k \times n \times l^2)$  (may be different under certain conditions such as  $k \times \log k > n \times l$  in which case sorting the delimiters has higher complexity than calculating the error) and definitely not larger than polynomial in  $k$ ,  $n$ , and  $l$ .

### 3.2 (p, k)-means

## 4 Appendix

### 4.1 Movie generation code

```

1 %% Movie
2 X = data;
3 k = 3;
4 % Randomly initialize centers of clusters
5 c = datasample(X, k, 'Replace', false);
6 r = repmat(0, size(X, 1), k);
7 oldr = 1; % something that is not equal to r initially
8 movie = [];
9
10 dist = @(x, y) norm(x-y);
11
12 % Loop as long as the clustering is changing
13 while ~isequal(r, oldr)
14     oldr = r;
15
16     % Assign points to clusters
17     for i = 1:size(X,1)
18         cluster = 1;
19         for j = 1:k
20             if dist(X(i, :), c(j, :)) < dist(X(i, :), c(cluster
21 , :))
22                 cluster = j;
23             end
24         end
25
26         r(i, :) = [repmat(0, 1, cluster-1) 1 repmat(0, 1, k-
27 cluster)];

```

```

26     end
27
28     %MOVIE GENERATION
29     % Split the data into the three clusters
30     cluster1 = [];
31     cluster2 = [];
32     cluster3 = [];
33
34     for i = 1:size(X, 1)
35         if r(i,1) == 1
36             cluster1 = [cluster1; X(i, :)];
37         end
38         if r(i,2) == 1
39             cluster2 = [cluster2; X(i, :)];
40         end
41         if r(i,3) == 1
42             cluster3 = [cluster3; X(i, :)];
43         end
44     end
45
46     % Plot the centers and the three clusters
47     hold on
48     scatter(c(:, 1), c(:, 2), 200, 'xblack');
49     if size(cluster1, 1) ~= 0
50         scatter(cluster1(:, 1), cluster1(:, 2));
51     end
52     if size(cluster2, 1) ~= 0
53         scatter(cluster2(:, 1), cluster2(:, 2));
54     end
55     if size(cluster3, 1) ~= 0
56         scatter(cluster3(:, 1), cluster3(:, 2));
57     end
58     legend('Centers', 'Cluster1', 'Cluster2', 'Cluster3');
59     hold off
60     movie = [movie getframe];
61     clf;
62     % MOVIE GENERATION OVER
63
64     % Update center positions
65     for i = 1:k
66         npoints = 0;
67         c(i, :) = 0;
68         for j = 1:size(r, 1)
69             c(i, :) = c(i, :) + r(j, i)*X(j, :);
70             npoints = npoints + r(j, i);
71         end
72         c(i, :) = c(i, :)/npoints;
73     end
74 end

```

```
75 movie2avi(movie, 'k-means.avi', 'fps', 1);
```