

COMPM012 - Coursework 3

Esben A. Sørig

December 1, 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

```

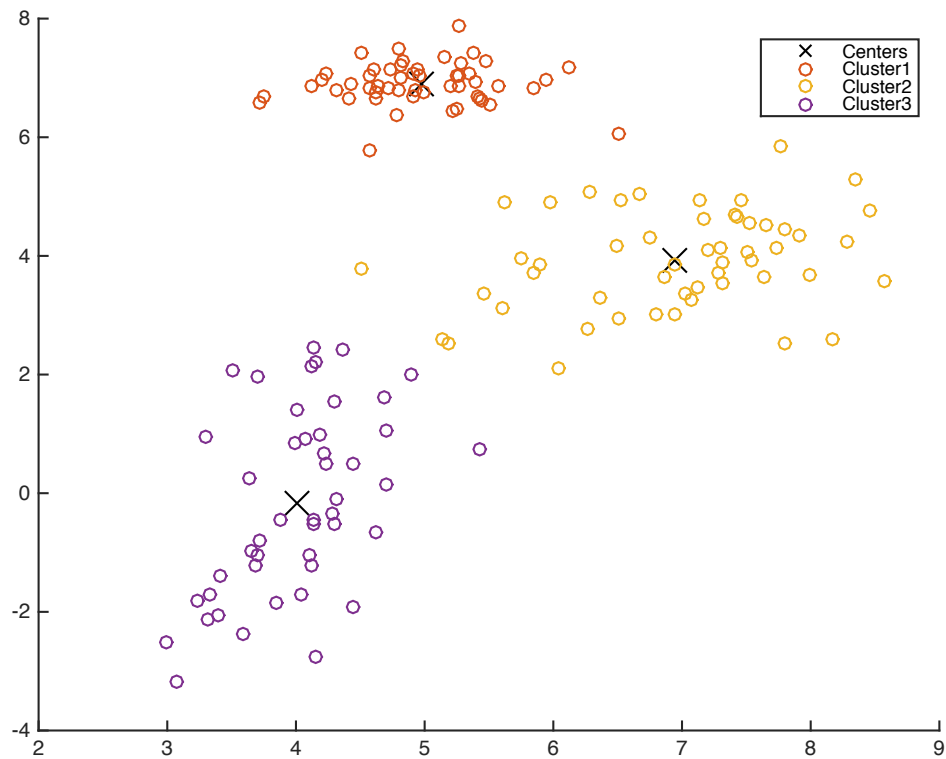
1 figure

```

```

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

```

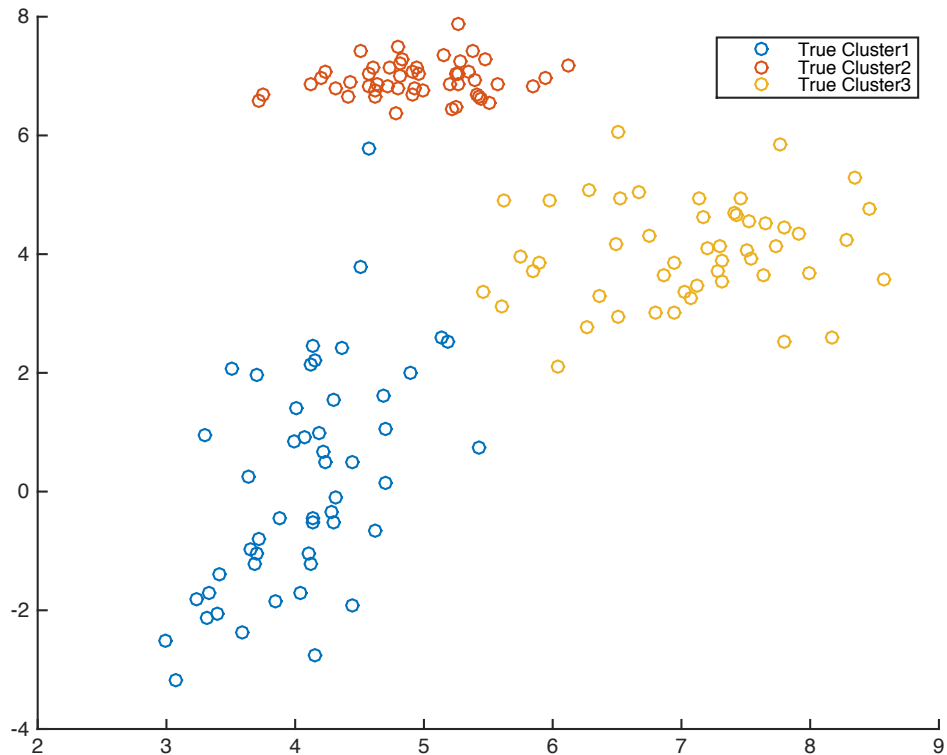


Comparing this to the true classifications

```

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     0.0333
3 stddeviationerror =
4     2.0922e-17

```

1.3 Iris data

Performing a similar experiment and error measure on the iris dataset is done in a similar way

```

1 iris = readtable('iris.dat');
2 params = table2array(iris(:, 1:4));
3 classifications = table2array(iris(:, 5));
4 classes = unique(classifications);
5 class1 = find(strcmp(classifications, classes(1)));
6 class2 = find(strcmp(classifications, classes(2)));
7 class3 = find(strcmp(classifications, classes(3)));
8
9 errors = [];
10 for n = 1:100
11     clusters = kmeans(params, 3);
12     clusters1 = clusters(class1);
13     clusters2 = clusters(class2);
14     clusters3 = clusters(class3);

```

```

15
16     label1 = mode(clusters1);
17     label2 = mode(clusters2);
18     label3 = mode(clusters3);
19
20     errcount = 0;
21     for i = 1:50
22         if clusters1(i) ~= label1
23             errcount = errcount + 1;
24         end
25         if clusters2(i) ~= label2
26             errcount = errcount + 1;
27         end
28         if clusters3(i) ~= label3
29             errcount = errcount + 1;
30         end
31     end
32
33     e = errcount/size(params, 1);
34     errors = [errors e];
35 end
36
37 mean = mean(errors)
38 stddeviation = std(errors)

```

Which returns

```

1 mean =
2     0.1096
3 stddeviation =
4     0.0117

```

2 Questions

2.1 Dataset with local minima

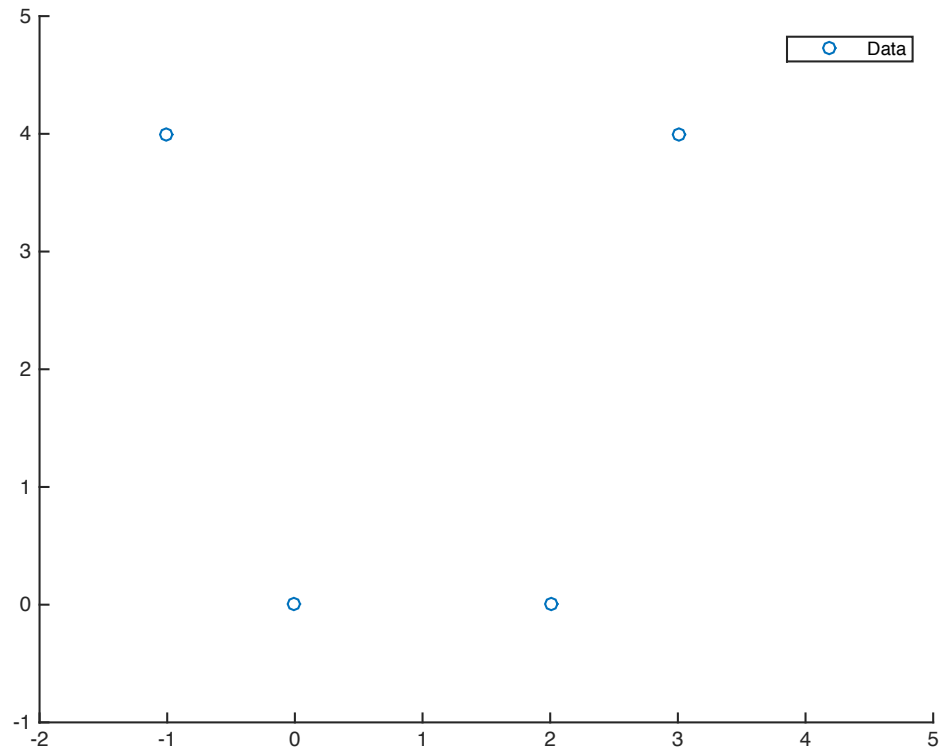
We can easily construct a dataset on which our k-means algorithm has local minima. An example of such a dataset is the set $\{(0, 0), (2, 0), (-1, 4), (3, 4)\}$. This is clear if we plot the data.

```

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:



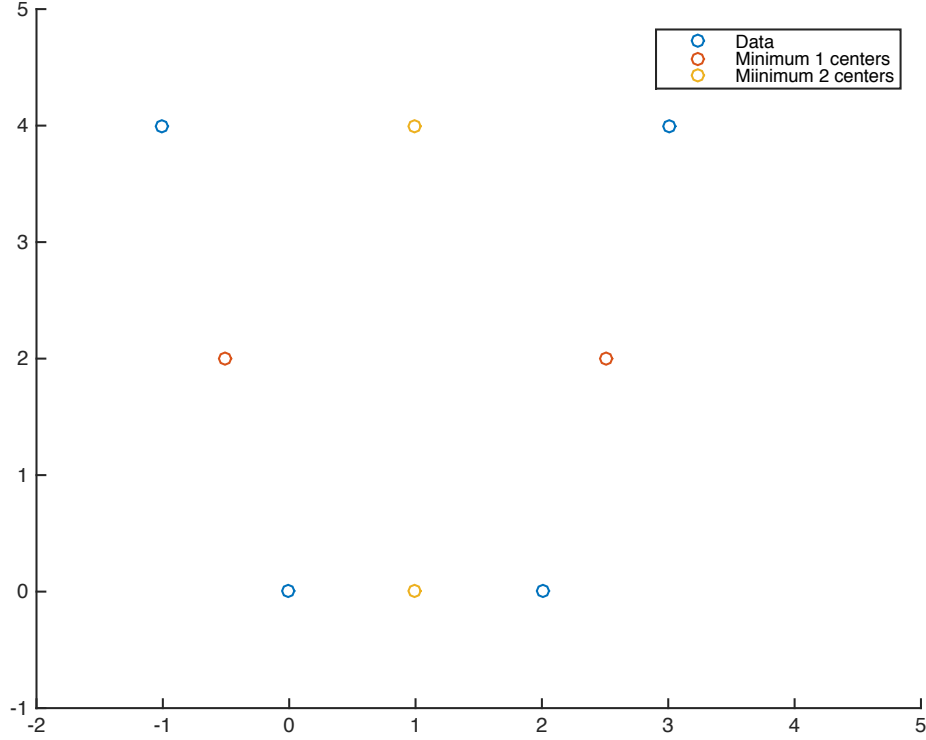
Running k-means on this dataset will converge in one of two local minima. We run the algorithm twice and plot the centers of the clusters: (in order to get two different results, this code may have to be run several times)

```

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 designing the segmentation algorithm. Let $\mathbf{i}_{n,k}^*$ be the set of segment delimiters $\{i_1^*, \dots, i_{k-1}^*\}$ that minimises the error of segmenting the points x_1, \dots, x_n into k segments. The key observation for construction of a polynomial time algorithm is that $\mathbf{i}_{n,k}^*$ can be described as

$$\mathbf{i}_{n,k}^* = \mathbf{i}_{j,k-1}^* \cup \{j\} \text{ where} \quad (4)$$

$$j = \underset{1 \leq j < n}{\operatorname{argmin}}(E(\mathbf{i}_{j,k-1}^* \cup \{j\}))$$

Here E denotes the error function of the segmentation produced a set of delimiters. From this observation we can construct a dynamic programming algorithm that builds a matrix of the optimal segmentations and a matrix of their errors. When the algorithm has reached the optimal segmentation $\mathbf{i}_{n,k}^*$ we have found the solution. The Matlab implementation of the algorithm is listed here:

```

1 function delimiters = segmentation(X, k)
2     % E(i, k) represents the minimum error of segmenting
3     % X(1:i, :) into k segments.
4     E = [];
5
6     % D(i, j) stores the jth delimiter that segments X(1:i, :)
7     % into the optimal j+1 segments.
8     D = [];
9
10    % We initialise E's first column (error of a single
11    % segment over the points X(1:i, :))
12    for i = 1:size(X, 1)
13        E(i,1) = error([], X(1:i, :));
14    end

```

```

15
16 % We then proceed to filling out each entry in E until
17 % we have reached our objective (the minimum k
18 % segmentation of all points in X).
19 for i = 1:size(X, 1)
20     for l = 2:min(i, k)
21         % Find min error of segmenting x_1, ..., x_i
22         % into l segments
23         minerror = Inf;
24         for j = 1:(i-1)
25             if j >= l
26                 e = E(j, l-1) + error([], X((j+1):i, :));
27                 if e < minerror
28                     minerror = e;
29                     D(i, l) = j; % Store delimiter
30                 end
31             end
32         end
33         E(i, l) = minerror;
34     end
35 end
36
37 % We retrieve the optimal delimiters from D
38 delimiters = [];
39 delimiter = size(X, 1);
40 for i = k:-1:2
41     delimiter = D(delimiter, i);
42     delimiters = [delimiters delimiter];
43 end
44 delimiters = sort(delimiters)
45 end

```

Informal argument of time complexity: The running time is polynomial in l , n , and k since the algorithm simply consists of nested loops running over l and k . In the inner most loop, the error function is called, which is similarly polynomial in complexity in k and l but also n (since we sum over the points). The complexity of the algorithm can therefore be bounded above by a constant multiple of some polynomial in n , l , and k .

3.2 (p, k)-means

We derive the minimiser of the objective for each C_i in the same way as with the squared distance. We let the derivative with respect to \mathbf{c}_k be equal to zero and solve for \mathbf{c}_k . We start by finding the derivative with respect to \mathbf{c}_k .

$$\frac{\delta}{\delta \mathbf{c}_k} \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d_p(\mathbf{c}_i, \mathbf{x}) = \sum_{\mathbf{x} \in C_k} \frac{\delta}{\delta \mathbf{c}_k} d_p(\mathbf{c}_k, \mathbf{x}) \quad (5)$$

We assume that the components of the inputs are nonnegative and obtain

$$\begin{aligned} \frac{\delta}{\delta \mathbf{c}_k} d_p(\mathbf{c}_k, \mathbf{x}) &= \frac{\delta}{\delta \mathbf{c}_k} \left(\sum_{i=1}^n \mathbf{c}_{k,i}^p - \mathbf{x}_i^p - p(\mathbf{c}_{k,i} - \mathbf{x}_i) \mathbf{x}_i^{p-1} \right) \\ &= p \mathbf{c}_k^{p-1} - p \mathbf{x}^{p-1} \end{aligned} \quad (6)$$

Where \mathbf{x}^n denotes that each component of \mathbf{x} is raised to the n th power. Inserting back in (5) and setting this equal to zero, we can now solve for \mathbf{c}_k

$$\begin{aligned} \sum_{\mathbf{x} \in C_k} (p \mathbf{c}_k^{p-1} - p \mathbf{x}^{p-1}) &= 0 \\ \implies \sum_{\mathbf{x} \in C_k} \mathbf{c}_k^{p-1} &= \sum_{\mathbf{x} \in C_k} \mathbf{x}^{p-1} \\ \implies \mathbf{c}_k &= \sqrt[p-1]{\frac{\sum_{\mathbf{x} \in C_k} \mathbf{x}^{p-1}}{|C_k|}} \end{aligned} \quad (7)$$

Observe that this result is consistent with our result for the euclidean distance. We can now design the (p, k) -means algorithm.

1. **for** $i = 1, \dots, k$ **do**

$$C_i = \{\mathbf{x} \in x \mid i = \operatorname{argmin}_{1 \leq j \leq k} d_p(\mathbf{c}_j, \mathbf{x})\}$$

2. **for** $j = 1, \dots, k$ **do**

$$\mathbf{c}_j = \sqrt[p-1]{\frac{\sum_{i=1}^l r_{ij} \mathbf{x}_i^{p-1}}{\sum_{i=1}^l r_{ij}}}$$

3. **While** not converged **go to** step 1.

Again, \mathbf{x}^n denotes that each component of \mathbf{x} is raised to the n th power and $\sqrt[n]{\mathbf{x}}$ denotes that we take the n th root of each component of \mathbf{x} . The argument of convergence for this algorithm is similar to the argument for the standard k -means algorithm. That is, the objective decreases in step 1 and 2 and the objective is bounded below. Hence the algorithm converges.

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         r(i, :) = [repmat(0, 1, cluster-1) 1 repmat(0, 1, k-
26 cluster)];
27     end
28
29 %MOVIE GENERATION
30 % Split the data into the three clusters
31 cluster1 = [];
32 cluster2 = [];
33 cluster3 = [];
34
35 for i = 1:size(X, 1)
36     if r(i,1) == 1
37         cluster1 = [cluster1; X(i, :)];
38     end
39     if r(i,2) == 1
40         cluster2 = [cluster2; X(i, :)];
41     end
42     if r(i,3) == 1
43         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);

```