

# COMPM012 - Coursework 3

Esben A. Sørig

November 27, 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

```

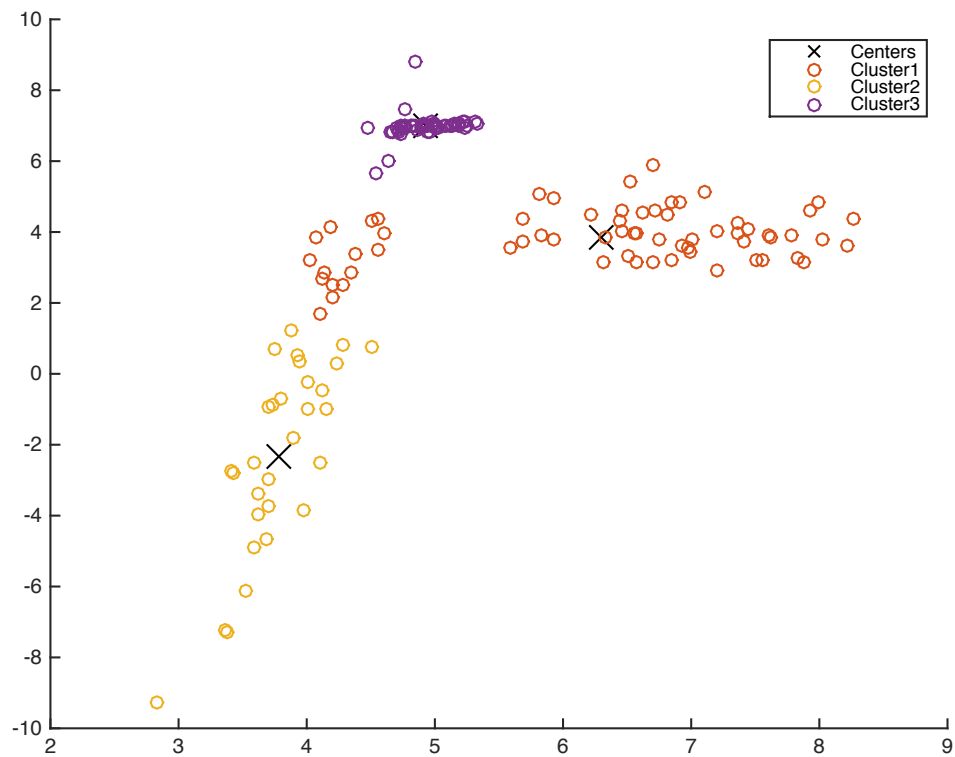
1 figure

```

```

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

```

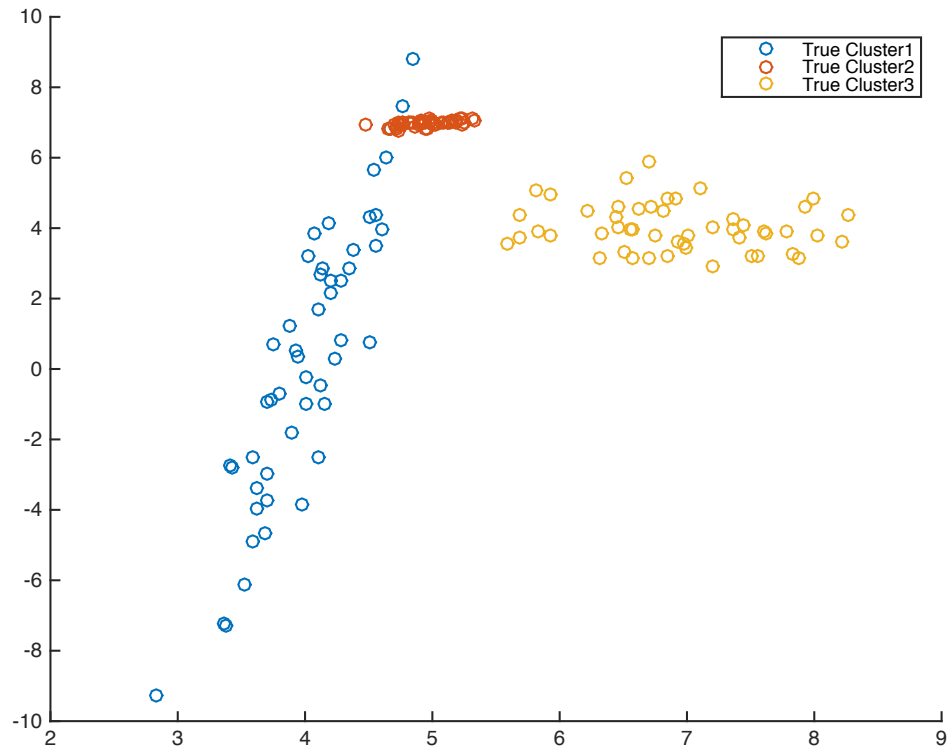


Comparing this to the true classifications

```

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

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

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

## 3 Extension

### 3.1 k-means segmentation

### 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         r(i, :) = [repmat(0, 1, cluster-1) 1 repmat(0, 1, k-
26 cluster)];
27     end
28 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);

```