# COMP3008: Machine Learning
# Arrhythmia

Ashley J. Robinson

School of Electronics and Computer Science
University of Southampton

May 15, 2013

## 1 Data handling

The data set consists of 452 instances of patient data each containing 279 attributes. This is not consist however as some attributes are not recorded for particular instances. Table 1 contains the offending attributes and the number of instances for which they are not recorded. It is clear that attribute 15 should be be neglected from the data with a missing frequency greater than 80%.

| Attribute number | Attribute name | Missing frequency |
|---|---|---|
| 11 | Vector angles in degrees on front plane of T | 8/452 |
| 12 | Vector angles in degrees on front plane of P | 22/452 |
| 13 | Vector angles in degrees on front plane of QRST | 1/452 |
| 14 | Vector angles in degrees on front plane of J | 376/452 |
| 15 | Heart rate (bpm) | 1/452 |

TABLE 1: Missing data

The reaming data attributes considered have a missing frequency of less than 5% so these can be considered as a positive influence towards classification. As the reaming four attributes are few in comparison to the total number no complex data completion approaches are warranted. The simplest approach is to replace the missing values with the mean of the attribute which is calculated from all instances recorded. Finally all data is normalised using equation (1) before being processed (Prügel-Bennett, 2013). Each element is replaced with a new element where $\mu_i$ is the mean and $\sigma_i$ is the standard deviation of the element attribute across all instances.

$$x_{ki} \rightarrow \frac{x_{ki} - \mu_i}{\sigma_i}, \ \mu_i = \frac{1}{P} \sum_{k=1}^{P} x_{ki}, \ \sigma_i = \sqrt{\frac{1}{P-1} \sum_{k=1}^{P} (x_{ki} - \mu_{ki})^2} \tag{1}$$

## 2 Applicable clustering techniques

The tasks requires classification of the data using unsupervised learning techniques as the only output data is the size of each cluster and the actual classification is unknown. Table 2 contains common

clustering techniques considered for this application (Leskovec and Rajaraman, 2010). K-means clustering is the most sensible choice for a preliminary investigation into this data as some small cluster sizes will prove difficult to identify with multi-class approaches. The algorithm is also simple and can be used as just a binary classifier.

| Technique | Description |
|---|---|
| Hierarchical clustering | A bottom-up approach that starts with each data point belonging to its own cluster and gradually joining pairs of data points to form clusters. |
| K-means clustering | Uses centroids to separate data into a set number (K) of clusters. |
| Fuzzy C-means clustering | Similar to k-means but uses cluster membership functions to create clusters as apposed to hard step functions. |

TABLE 2: Common clustering techniques

# 3    K-means clustering implementation

The first step initialises the algorithm by randomly selecting data points to contribute to each of the K clusters. The cluster centre is manifested as the centroid of all the points then every data point is compared with the cluster centres. Each data point is now assigned to the nearest cluster. The new cluster centres are calculated using their respective assigned data points and the whole assignment process is repeated until there is no change in the cluster vectors. All MATLAB code for this implementation is contained in appendix A.

## 3.1    Low dimension verification

A naive approach to the data set suggests that the condition of a heart can be determined by just the age, weight and heart rate of the patient. This is of course not a useful approach to arrhythmia classification but does produce a three dimensional pattern on which to test the algorithm. It can be seen in figure 1 that the method successfully classifies data; from visual feedback on data clustering. This increases confidence in the algorithms performance in higher dimensions.

## 3.2    Tuning

The aim for this classification is to assign input patterns to the classes that have been identified. The 13 classes for which respective data is contained in the database can be either be individually learnt or 12 of these classes can be grouped as positive for arrhythmia. The approach is to attempt binary classification.

Initial tests yield varied results as the algorithm will converge on different mean values for the two cluster centres. These is because of the random initialisation. This problem can be overcome by iterating over the entire algorithm. The size of the two clusters is know so it can always be assumed that the cluster for which the most data points belongs to is the normal class. A list of the cluster centres can be compiled for the two clusters and then the centroid of all these calculated. This is effectively continuing the algorithm by taking another centroid of data points.
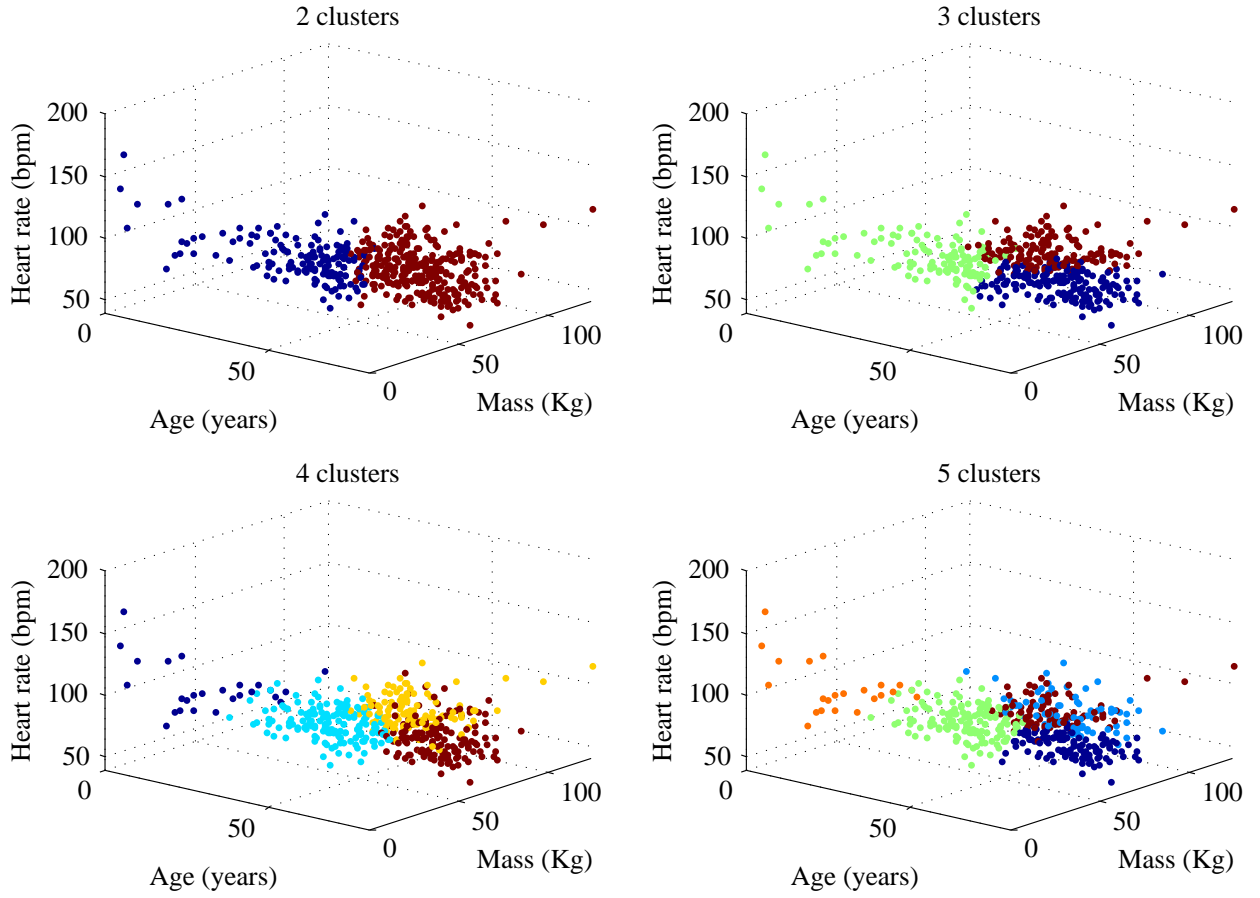
FIGURE 1: K-means clustering on three dimensional data

Occasionally initialisation may lead to cluster sizes that are almost identical in which case it is hard guarantee a correct classification of the centre. This can be overcome by setting a threshold for clusters with similar sizes then simply choosing to disregard the outcome of that iteration.

## 4   Performance

Using the mean output of all the cluster centres the variance can be calculated by comparing the individual runs to the final solution. Equation (2) contains the mean output clusters ($\hat{\mathbf{y}}_i$) from $n$ iterations. The variance is the Euclidean distance between the mean cluster centres ($\hat{\mathbf{y}}_i$) and each of the iteration cluster centres ($\mathbf{y}_{ik}$).

The two histograms in figure 2 display the variance of the algorithm output over 250 iterations. Both exhibit grouping of variance. In fact the groupings represent a convergence on the wrong cluster. Using these incorrectly classified values to calculate the mean, and therefore the variance, throws all values off.

$$V_{ik} = ||\hat{\mathbf{y}}_i - \mathbf{y}_{ik}||, \quad \hat{\mathbf{y}}_i = \frac{1}{n}\sum_{k=1}^{n}\mathbf{y}_{ik}, \quad i \in 1,2 \tag{2}$$
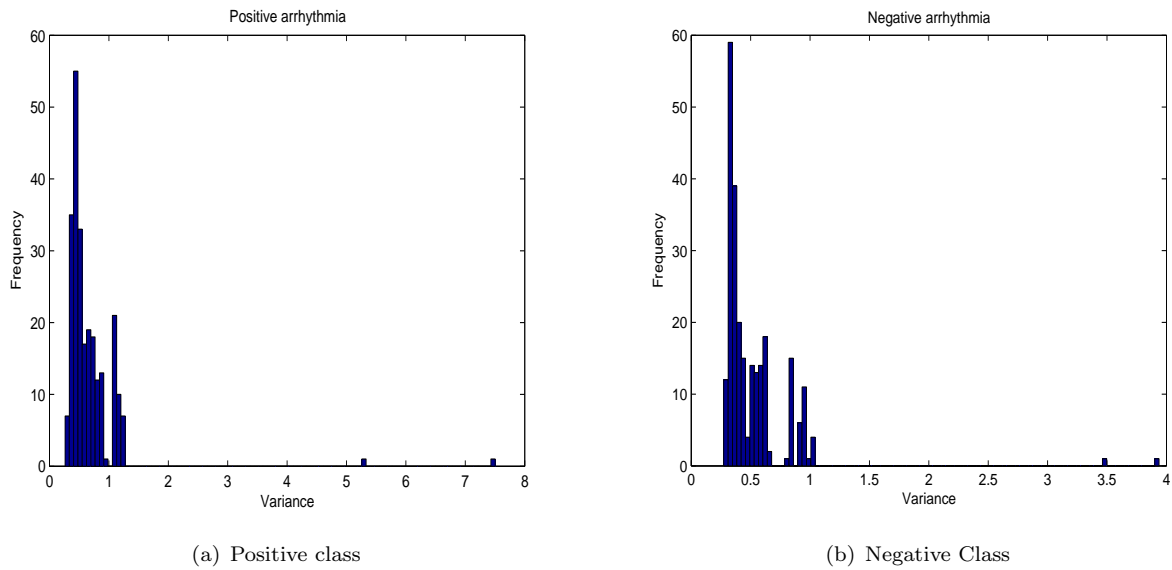
(a) Positive class        (b) Negative Class

FIGURE 2: Histograms of variance through algorithm iterations

The outliers are removed from figure 2 and the variance is recalculated[1] for 248 iterations in figure 3. This also shows a grouping because of the unpopulated variance bins in the middle but this instead is down to the nature of the data. The algorithm will converge in different centres depending on the initialisation.



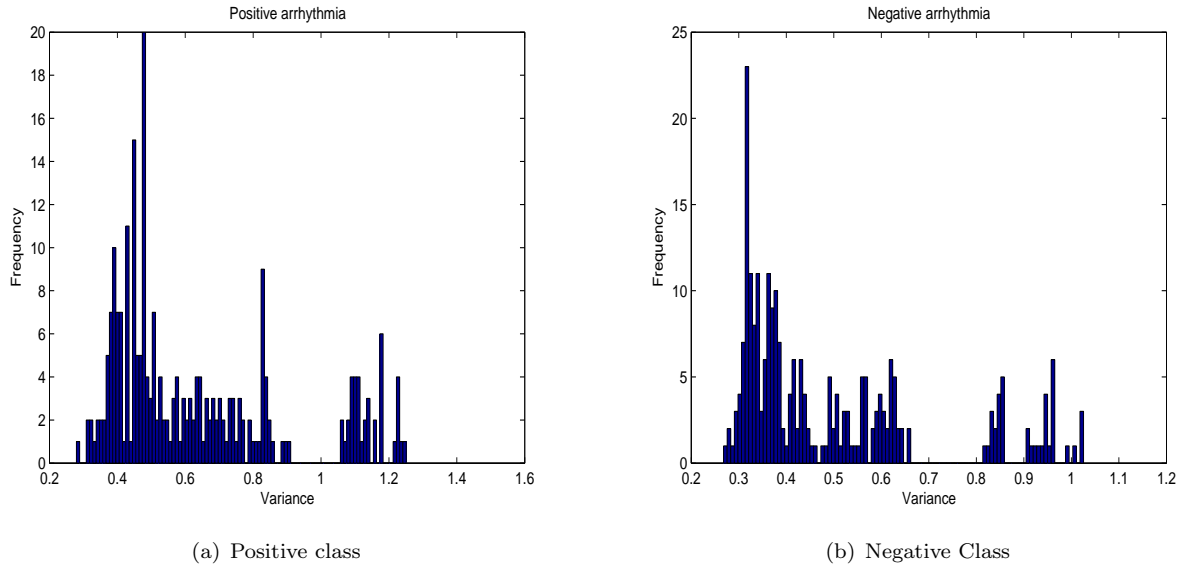(a) Positive class        (b) Negative Class

FIGURE 3: Fixed histograms of variance through algorithm iterations with outliers removed

The mean variance for this data is 0.4259 for the negative clusters and 0.5509 for the positive clusters. This is an estimation on the performance and reveals a need for improvement but shows adequate classification. The positive skew of both histograms in figure 3 verifies the performance.

---

[1]Code for kicking data points and iterating the algorithm is contained in listing 8

# 5   Conclusions and future work

While actual classification of the data in unknown this approach yielded good results after analysis of the variance. The algorithm would need adapting to reject wrongly classified output itself. This would remove the higher values of variance and improve the mean variance of the solution. The quantity threshold for assuming data belongs to a class also needs further investigation.

A much better approach could be produced using supervised learning. Labelled data could be obtained from another data set or through medical analysis of this data set.

# References

Jure Leskovec and Anand Rajaraman. CS345a:Data Mining, 2010. Stanford University, http://www.stanford.edu/class/cs345a/slides/12-clustering.pdf.

Adam Prügel-Bennett. COMP3008 Machine Learning, 2013. Southampton University.

# A  Code listings

```
1   %Ashley Robinson
2   %10/04/12
3   %COMP3008
4   %K_Means_Clustering.m
5   clear
6   % K >= 2
7   K = 2;
8   %Load the data
9   data = LoadData('arrhythmia.data');%Load the data from text file
10  data = Normalise(data);
11  instances = size(data,2);%These are handy to know
12  attributes = size(data,1);
13  %populate all means
14  count = 0;
15  while(min(count) == 0)
16      %initialise
17      assigned_cluster = randi([1,K],1,instances);
18      [old_cluster_mean] = Calculate_Means(assigned_cluster,K,data);
19      assigned_cluster = Assign_New_Cluster(old_cluster_mean,K,data);%Assign data to random
         cluster means
20      [cluster_mean] = Calculate_Means(assigned_cluster,K,data);
21      %Converge
22      i = 0;
23      while(sum(abs(sum(old_cluster_mean - cluster_mean)  ~= zeros(1,K))))%Test for same means
24          old_cluster_mean = cluster_mean;
25          assigned_cluster = Assign_New_Cluster(cluster_mean,K,data);
26          [cluster_mean] = Calculate_Means(assigned_cluster,K,data);
27          i = i + 1;
28          Output(i,K,instances,assigned_cluster);%User output
29      end
30      count = Cluster_Count(K,instances,assigned_cluster);%How many data points to each cluster?
31  end
```

LISTING 1: Main algorithm

```
1   %Ashley Robinson
2   %10/04/12
3   %COMP3008
4   %LoadData.m
5   function[data] = LoadData(Path_to_data)
6       %Open data file in same directory
7       arrhythmia_data = fopen(Path_to_data);
8       %global init
9       instance = 0;
10      while ~feof(arrhythmia_data)%While not at the end of the file
11          instance = instance + 1;
12          line = fgetl(arrhythmia_data);%Read a line from the file
13          start = 1;
14          %disp(line)%Just to see some output
15          attribute = 0;
16          for j=1:length(line)
17              %Look for CSVs
18              if (line(j) == ',')
19                  attribute = attribute + 1;
20                  if (line(start:(j - 1)) == '?')
21                      miss_mask(attribute,instance) = 1;%take a note of missing
22                      data(attribute,instance) = 0;
23                  else
24                      miss_mask(attribute,instance) = 0;%need to keep dimensions up
25                      data(attribute,instance) = str2num(line(start:(j - 1)));
26                  end
27                  start = j + 1;%One in front of the comma
28              end
29          end
30      end
31      fclose(arrhythmia_data);
32      %Replace unkown data with the mean of that attribute
33      for i=1:attribute
```

```
34          total = 0;
35          count = 0;
36          for j=1:instance
37              if(miss_mask(i,j) == 0)
38                  total = total + data(i,j);
39                  count = count + 1;
40              end
41          end
42          total = total./count;%The mean of this attribute
43          miss_mask(i,:) = miss_mask(i,:).*total;
44      end
45      data = data + miss_mask;%Missing values are zero so add miss mask
46      %Remove attribute 14
47      att_point = 1;
48      for i=1:attribute
49          if(i ~= 14)
50              for j=1:instance
51                  temp(att_point,j) = data(i,j);
52              end
53              att_point = att_point + 1;
54          end
55      end
56      data = temp;
57      %All info now in matrix
58      info =  sprintf('Instances: %d \nAttributes: %d ',size(data,2), size(data,1));
59      disp(info)
60  end
```

LISTING 2: Load the data from the raw text file

```
1   %Ashley Robinson
2   %21/04/12
3   %COMP3008
4   %Normalise.m
5   function[new_data] = Normalise(data)
6       instances = size(data,2);%These are handy to know
7       attributes = size(data,1);
8       for i=1:attributes
9           mu = data(i,1);
10          for j=2:instances
11              mu = mu + data(i,j);
12          end
13          mu = mu./instances;
14          sigma = (data(i,1) - mu)^2;
15          for j=2:instances
16              sigma = sigma + (data(i,j) - mu)^2;
17          end
18          sigma = sqrt(sigma./(instances - 1));
19          for j=1:instances
20              if(sigma == 0)
21                  new_data(i,j) = 0;
22              else
23                  new_data(i,j) = (data(i,j) - mu)./sigma;
24          end
25      end
26  end
```

LISTING 3: Normalise the data

```
1   %Ashley Robinson
2   %10/04/12
3   %COMP3008
4   %Calculate_Means.m
5   function[cluster_mean] = Calculate_Means(assigned_cluster,K,data)
6       instances = size(data,2);
7       attributes = size(data,1);
8       for i=1:K
9           cluster_mean(:,i) = zeros(attributes,1);
10          count = 0;
11          for j=1:instances
12              if(assigned_cluster(j) == i)
```

```
13                count = count + 1;%Keep track of the number of instances in the cluster
14                cluster_mean(:,i) = cluster_mean(:,i) + data(:,j);%Total all the vectors
15            end
16        end
17        if(count ~= 0)
18            cluster_mean(:,i) = cluster_mean(:,i)./count;%calculate mean
19        end
20    end
21 end
```

LISTING 4: Calculate the mean of a cluster of data points

```
1  %Ashley Robinson
2  %10/04/12
3  %COMP3008
4  %Assign_New_Cluster.m
5  function[assigned_cluster] = Assign_New_Cluster(cluster_mean,K,data)
6     for i=1:size(data,2)%instances
7         assigned_cluster(i) = 1;%init
8         dist = Calculate_Distance(cluster_mean(:,1),data(:,i));
9         for j=2:K
10            test = Calculate_Distance(cluster_mean(:,j),data(:,i));
11            if(test < dist)
12                assigned_cluster(i) = j;
13                dist = test;
14            end
15        end
16    end
17 end
```

LISTING 5: Assign each data point to the nearest cluster cluster

```
1  %Ashley Robinson
2  %10/04/12
3  %COMP3008
4  %Output.m
5  function[] = Output(iter,K,instances,assigned_cluster)
6     count = Cluster_Count(K,instances,assigned_cluster);%How many data points to each cluster
7     disp(sprintf('Iteration: %d',iter));
8     for i=1:K
9         disp(sprintf('              Cluster %d size: %d ',i,count(i)))
10    end
11 end
```

LISTING 6: View some useful information about what is going on

```
1  %Ashley Robinson
2  %10/04/12
3  %COMP3008
4  %Cluster_Count.m
5  function[count] = Cluster_Count(K,instances,assigned_cluster)
6     count = zeros(K,1);
7     for i=1:instances
8         for j=1:K
9             if(assigned_cluster(i) == j)
10                count(j) = count(j) + 1;
11            end
12        end
13    end
14 end
```

LISTING 7: How many data points are associated with each cluster

```
1  %Ashley Robinson
2  %10/04/12
3  %COMP3008
4  %Iterate_K_Means.m
5  clear
6  %Binary fixed
```

```
 7  K = 2;
 8  iterations = 250;
 9
10  %Load the data
11  data = LoadData('arrhythmia.data');%Load the data from text file
12  data = Normalise(data);
13  instances = size(data,2);%These are handy to know
14  attributes = size(data,1);
15
16
17  for iter=1:iterations
18      info = sprintf('Iteration: %d',iter)
19      disp(info)
20      %populate all means
21      count = 0;
22      while(min(count) == 0)
23          %initialise
24          assigned_cluster = randi([1,K],1,instances);
25          [old_cluster_mean] = Calculate_Means(assigned_cluster,K,data);
26          assigned_cluster = Assign_New_Cluster(old_cluster_mean,K,data);%Assign data to random
         cluster means
27          [cluster_mean] = Calculate_Means(assigned_cluster,K,data);
28          %Converge
29          i = 0;
30          while(sum(abs(sum(old_cluster_mean - cluster_mean)  ~= zeros(1,K))))%Test for same means
31              old_cluster_mean = cluster_mean;
32              assigned_cluster = Assign_New_Cluster(cluster_mean,K,data);
33              [cluster_mean] = Calculate_Means(assigned_cluster,K,data);
34              i = i + 1;
35              %Output(i,K,instances,assigned_cluster);%User output
36          end
37          count = Cluster_Count(K,instances,assigned_cluster)%How many data points to each cluster
        ?
38          if(abs(count(1) - count(2)) < 20)
39              count = 0;
40              disp('close')
41          end
42      end
43      cluster(:,:,iter) = cluster_mean;
44  end
45
46  %First mean holds good patients
47  assigned_cluster = Assign_New_Cluster(cluster(:,:,1),K,data);
48  count = Cluster_Count(K,instances,assigned_cluster);
49  if(count(1) > count(2))
50      mean_cluster(:,1) = cluster(:,1,1);
51      mean_cluster(:,2) = cluster(:,2,1);
52  else
53      mean_cluster(:,2) = cluster(:,1,1);
54      mean_cluster(:,1) = cluster(:,2,1);
55  end
56
57  for iter=2:iterations
58      assigned_cluster = Assign_New_Cluster(cluster(:,:,iter),K,data);
59      count = Cluster_Count(K,instances,assigned_cluster);
60
61      if(count(1) > count(2))
62          mean_cluster(:,1) = mean_cluster(:,1) + cluster(:,1,iter);
63          mean_cluster(:,2) = mean_cluster(:,2) + cluster(:,2,iter);
64      else
65          mean_cluster(:,2) = mean_cluster(:,2) + cluster(:,1,iter);
66          mean_cluster(:,1) = mean_cluster(:,1) + cluster(:,2,iter);
67      end
68  end
69  mean_cluster = mean_cluster./iterations
70  assigned_cluster = Assign_New_Cluster(mean_cluster,K,data);
71  count = Cluster_Count(K,instances,assigned_cluster)
72
73  for iter=1:iterations
74      assigned_cluster = Assign_New_Cluster(cluster(:,:,iter),K,data);
75      count = Cluster_Count(K,instances,assigned_cluster);
```

```
76     if(count(1) > count(2))
77         variance(iter,1) = Calculate_Distance(mean_cluster(:,1),cluster(:,1,iter));
78         variance(iter,2) = Calculate_Distance(mean_cluster(:,2),cluster(:,2,iter));
79     else
80         variance(iter,2) = Calculate_Distance(mean_cluster(:,2),cluster(:,1,iter));
81         variance(iter,1) = Calculate_Distance(mean_cluster(:,1),cluster(:,2,iter));
82     end
83 end
84 figure(1)
85 hist(variance(:,1),100)
86 xlabel('Variance')
87 ylabel('Frequency')
88 title('Cluster 1')
89 figure(2)
90 hist(variance(:,2),100)
91 title('Cluster 2')
92 xlabel('Variance')
93 ylabel('Frequency')
94 variance
95
96
97 kick1 =input('kick 1 (0 to quit):')
98 kick2 =input('kick 2:')
99 while(sum(kick1) ~= 0)
100 %-------------variance fix
101 j = 0;
102 for iter=1:iterations
103     if(sum(variance(iter,1) > kick1) | sum(variance(iter,2) > kick2))
104         disp('out')
105     else
106         j = j + 1;
107         new_cluster(:,:,j) = cluster(:,:,iter);
108     end
109 end
110
111 %First mean holds good patients
112 assigned_cluster = Assign_New_Cluster(new_cluster(:,:,1),K,data);
113 count = Cluster_Count(K,instances,assigned_cluster);
114 if(count(1) > count(2))
115     mean_cluster(:,1) = new_cluster(:,1,1);
116     mean_cluster(:,2) = new_cluster(:,2,1);
117 else
118     mean_cluster(:,2) = new_cluster(:,1,1);
119     mean_cluster(:,1) = new_cluster(:,2,1);
120 end
121
122 for iter=2:j
123     assigned_cluster = Assign_New_Cluster(new_cluster(:,:,iter),K,data);
124     count = Cluster_Count(K,instances,assigned_cluster);
125
126     if(count(1) > count(2))
127         mean_cluster(:,1) = mean_cluster(:,1) + new_cluster(:,1,iter);
128         mean_cluster(:,2) = mean_cluster(:,2) + new_cluster(:,2,iter);
129     else
130         mean_cluster(:,2) = mean_cluster(:,2) + new_cluster(:,1,iter);
131         mean_cluster(:,1) = mean_cluster(:,1) + new_cluster(:,2,iter);
132     end
133 end
134 mean_cluster = mean_cluster./j
135 assigned_cluster = Assign_New_Cluster(mean_cluster,K,data);
136 count = Cluster_Count(K,instances,assigned_cluster)
137
138 for iter=1:j
139     assigned_cluster = Assign_New_Cluster(new_cluster(:,:,iter),K,data);
140     count = Cluster_Count(K,instances,assigned_cluster);
141     if(count(1) > count(2))
142         new_variance(iter,1) = Calculate_Distance(mean_cluster(:,1),new_cluster(:,1,iter));
143         new_variance(iter,2) = Calculate_Distance(mean_cluster(:,2),new_cluster(:,2,iter));
144     else
145         new_variance(iter,2) = Calculate_Distance(mean_cluster(:,2),new_cluster(:,1,iter));
146         new_variance(iter,1) = Calculate_Distance(mean_cluster(:,1),new_cluster(:,2,iter));
```

```
147       end
148  end
149  figure(3)
150  hist(new_variance(:,1),100)
151  xlabel('Variance')
152  ylabel('Frequency')
153  title('Cluster 1')
154  figure(4)
155  hist(new_variance(:,2),100)
156  title('Cluster 2')
157  xlabel('Variance')
158  ylabel('Frequency')
159  new_variance
160  kick1 =input('kick 1 (0 to quit):')
161  kick2 =input('kick 2:')
162  end
```

LISTING 8: Used to iterate over the algorithm and kick output from the variance calculation