# NEURAL NETWORKS

# ASSIGNMENT
# UNSUPERVISED LEARNING
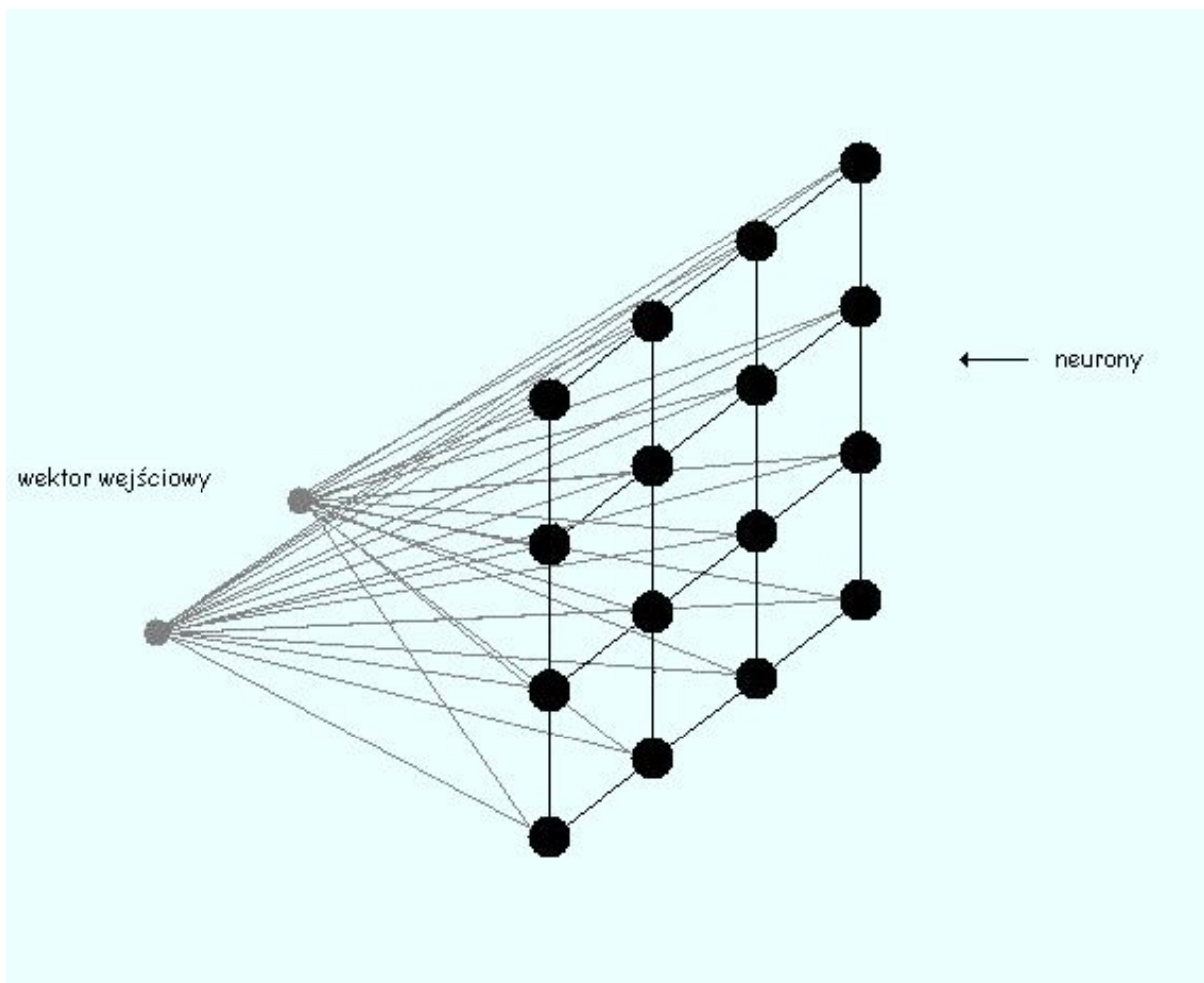# Non-Parametric CLustering

**Keya Desai**

**201501012**

# Self organizing Maps

## Concept:

- SOM is an Unsupervised Neural Network learning.
- A self-organizing map consists of components called nodes or neurons. Associated with each node are a weight vector of the same dimension as the input data vectors, and a position in the map space.
- The neighboring cells in the neural network compete in their activities by means of mutual lateral interactions.

Functioning of self-organizing neural network is divided into three stages:

- construction
- learning
- identification



**2-D map of neurons**

## Steps:

Training occurs in several steps and over many iterations:

1. Each node's weights are initialized.
2. A vector is chosen at random from the set of training data and presented to the lattice.
3. Every node is examined to calculate which one's weights are most like the input vector. The winning node is commonly known as the Best Matching Unit (BMU).
4. The radius of the neighbourhood of the BMU is now calculated. This is a value that starts large, typically set to the 'radius' of the lattice, but diminishes each time-step. Any nodes found within this radius are deemed to be inside the BMU's neighbourhood.
5. Each neighbouring node's (the nodes found in step 4) weights are adjusted to make them more like the input vector. The closer a node is to the BMU, the more its weights get altered.
6. Repeat step 2 for N iterations.

# IRIS Dataset:

For Iris dataset, the number of neurons in output layer is 3 and in input layer(i.e number of features) is 4.

Code:

```matlab
%self organization map implementation
clc;
load iris_dat.dat;
  X = iris_dat;

[vectors_num, dim] = size(X);

k=3;                              %number of cllusters
max_iterations=100;
weights=initweights(X,dim,k);    %initizializing the weights
alpha=0.6;
M=zeros(vectors_num,1);

radius=5;
disp('radius=');
disp(radius);
for i=1:max_iterations

dnodes=zeros(k,k);        %Matrix for distance between output nodes
for w=1:k-1              %calculating the distance between nodes
    for j=w+1:k
```

```matlab
23
24 -            for n=1:dim-1
25 -                dnodes(w,j)=dnodes(w,j)+ ((weights(n,w)-weights(n,j))^2);
26 -            end
27 -            dnodes(w,j)=dnodes(w,j)^(1/2);
28 -            dnodes(j,w)=dnodes(w,j);
29 -        end
30 -    end
31
32 -    for j=1:vectors_num
33 -        [distance,p]=getdistancemetrix(X,dim,k,j,weights);
34 -        weights=updateweights(X,dim,j,alpha,weights,p);  %updating the weights of PMU
35 -        for y=1:k                                        %updating weights of neighbouring units
36 -            if (dnodes(p,y)<radius) && (dnodes(p,y)>0)
37 -                weights=updateweights(X,dim,j,alpha,weights,y);
38 -            end
39 -        end
40 -    end
41 -    alpha=alpha/2;
42 -    radius=radius/2;
43
44 -    end

46 -    for j=1:vectors_num
47 -        [distance,p]=getdistancemetrix(X,dim,k,j,weights);
48 -        M(j,1)=p;
49 -    end
50 -    C = confusionmat(M,X(:,5));
51 -    C1 = C(1,1)/50;
52 -    C2 = C(2,2)/50;
53 -    C3 = C(3,3)/50;
54 -    disp(' Individual class efficiency for');
55 -    disp('Class1:');
56 -    disp(C1);
57 -    disp('Class 2:')
58 -    disp(C2);
59 -    disp('Class 3');
60 -    disp(C3);
61 -    SUM = C1 + C2 + C3;
62 -    SUM1 = SUM/3;
63 -    disp('Average Efficiency=');
64 -    disp(SUM1);
65 -    SUM2=(C(1,1)+C(2,2)+C(3,3))/150;
66 -    disp('Overall Efficieny=');
67 -    disp(SUM2);
68 -    disp('Confustion Matrix:');
69 -    disp(C);
```
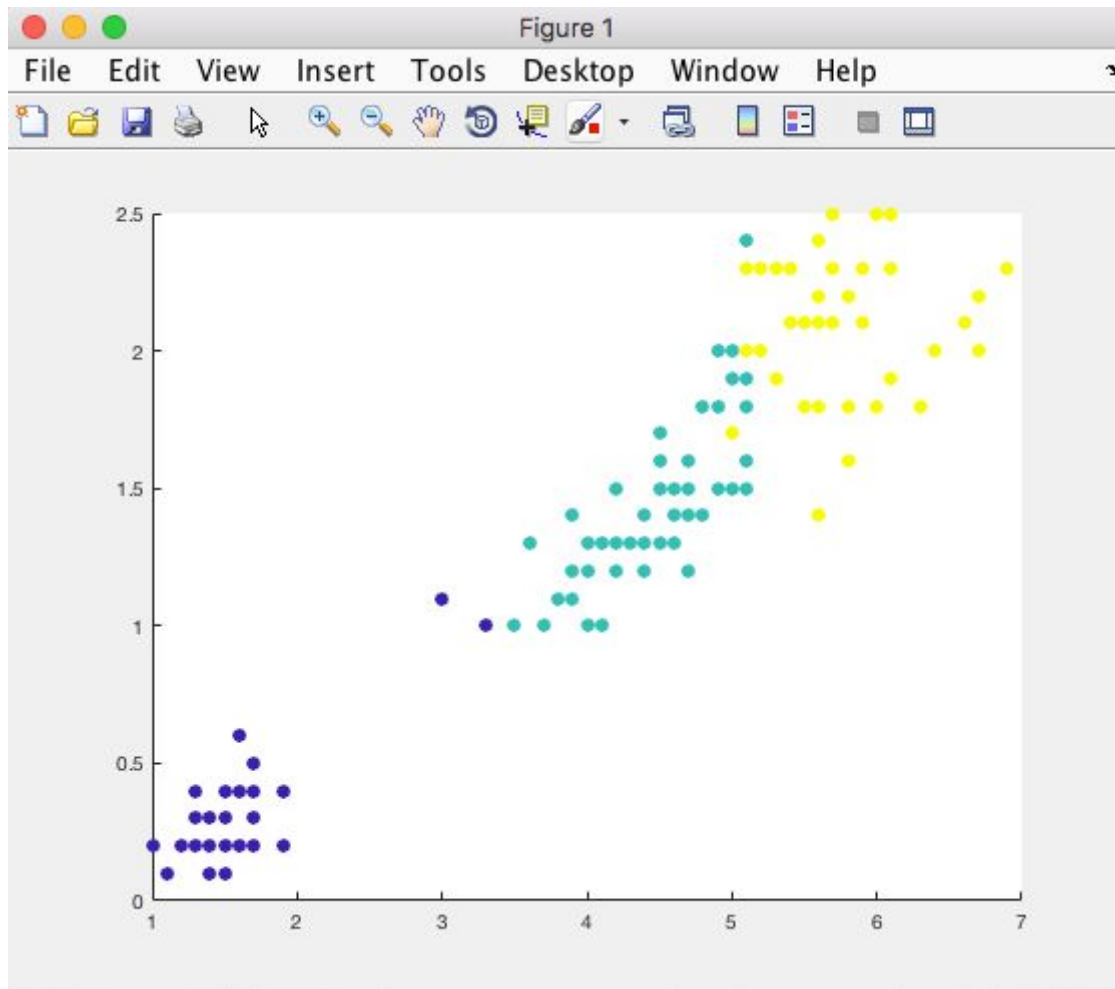
```matlab
70 -      scatter(X(:,3),X(:,4),50,M,'filled');
71
72        function weights = initweights(X,dim,k)
73 -          weights=zeros(dim-1,k);
74 -      for i=1:dim-1
75 -      maxi = max(X(:,i));
76 -      mini = min(X(:,i));
77 -      weights(i,:) = mini + (maxi-mini).*rand (1,k);
78 -      end
79 -      end
80
81        function [distance,p]=getdistancemetrix(X,dim,k,j,weights)
82
83 -          distance=zeros(k,1);
84 -           for m=1:k
85
86 -               for n=1:dim-1
87 -                   distance(m,1)=distance(m,1)+ ((X(j,n)-weights(n,m))^2);
88 -               end
89 -               distance(m,1)=distance(m,1)^(1/2);
90 -           end
91 -           [~,I] = min(distance);
92 -             p=I(1,1);
93 -      end


94
95        function weights=updateweights(X,dim,j,alpha,weights,p)
96
97
98 -           for i=1:dim-1
99 -               temp=weights(i,p);
100 -              weights(i,p)=temp+(alpha*(X(j,i)-temp) );
101 -           end
102
103 -      end
104
105
```
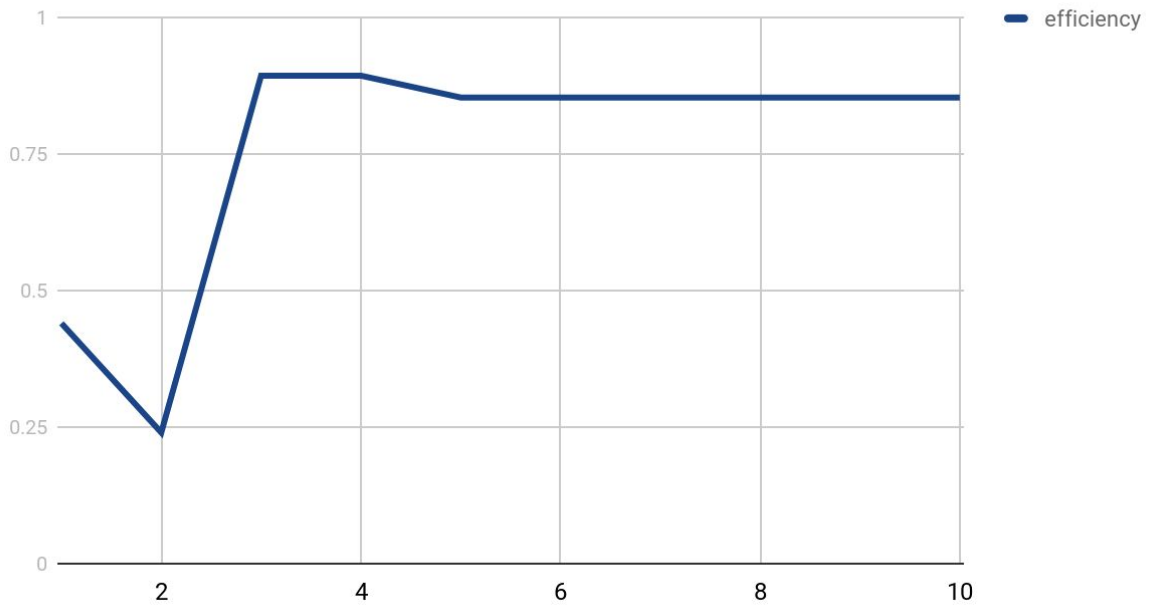
**Output:**

# Radius Analysis:

## Concept

### Determining the Best Matching Unit's Local Neighbourhood

- After for each iteration,the BMU(Best Matching Unit) has been determined, the next step is to calculate which of the other nodes are within the BMU's neighbourhood. All these nodes have their weight vectors altered in the next step.
- To get the neighbours, the distance between each neuron is calculated and if it below a specific threshold(given by variable radius) the weights of that neuron is also updated.
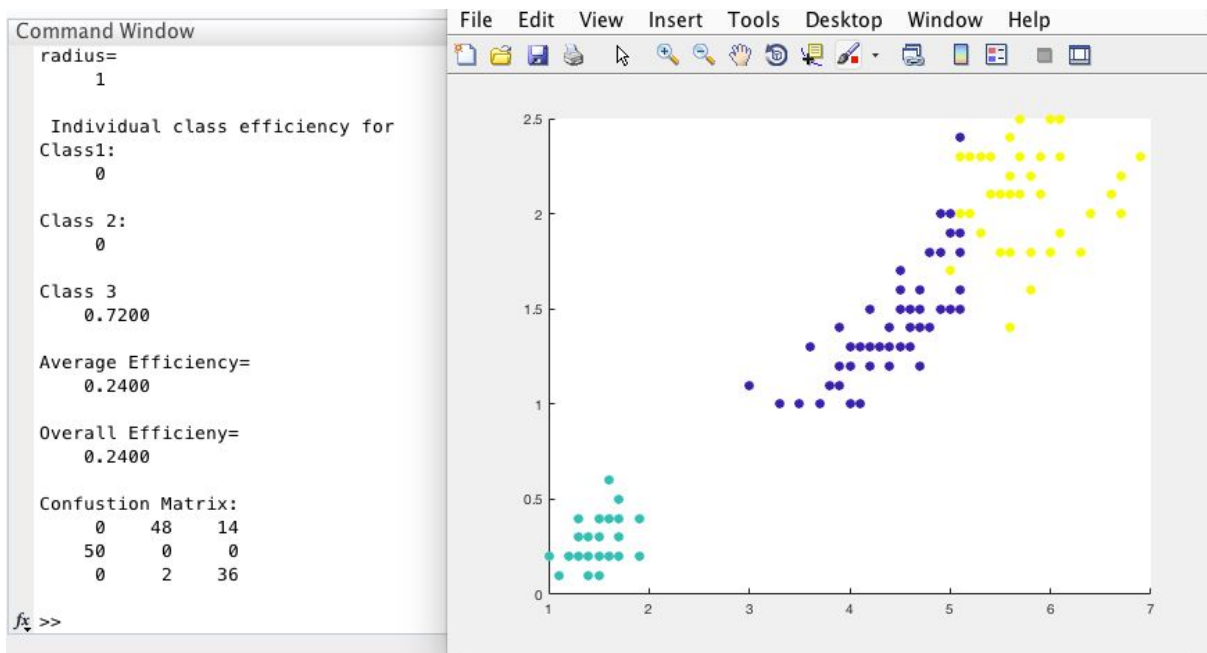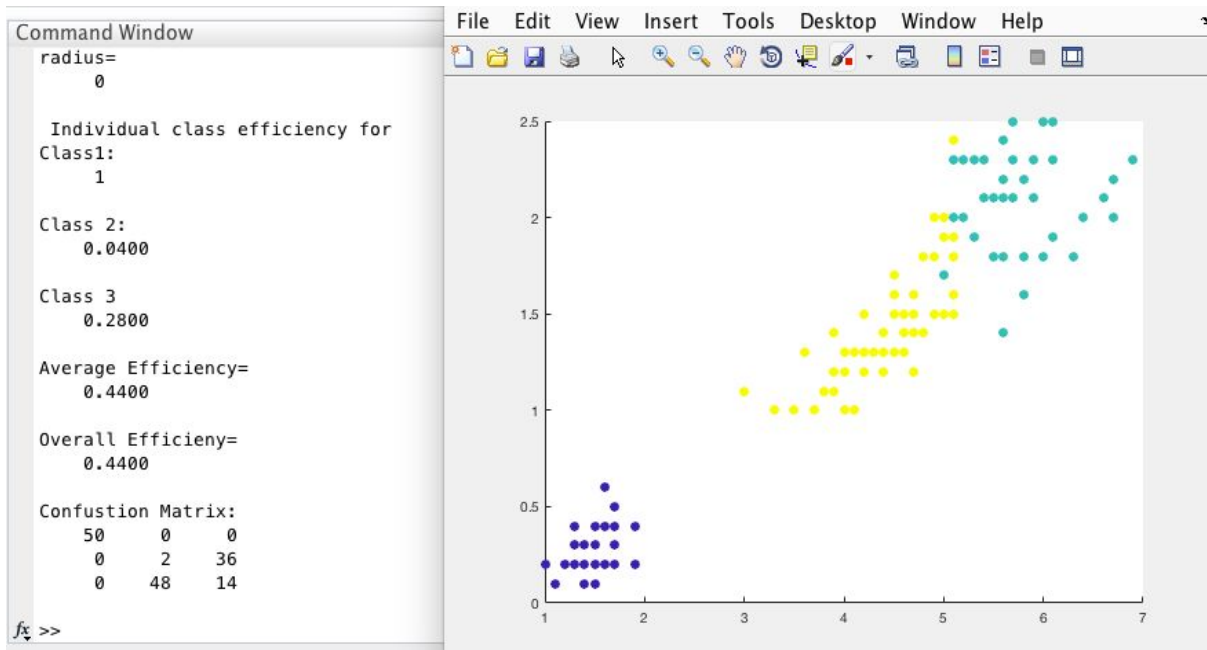
## Observation

- For different values of radius,we calculated the average efficiency.We found that by initializing the radius with 3 and 4 ,the efficiency was maximum.As the efficiency was increased the efficiency decreased by 0.01 and remained constant thereafter.
- For radius=0 i.e not using the concept of neighbourhood neurons, we got the least efficiency.
- Also, by keeping radius to be a finite value, we got a consistent result on running the code multiple times which is not the case when radius is kept 0.
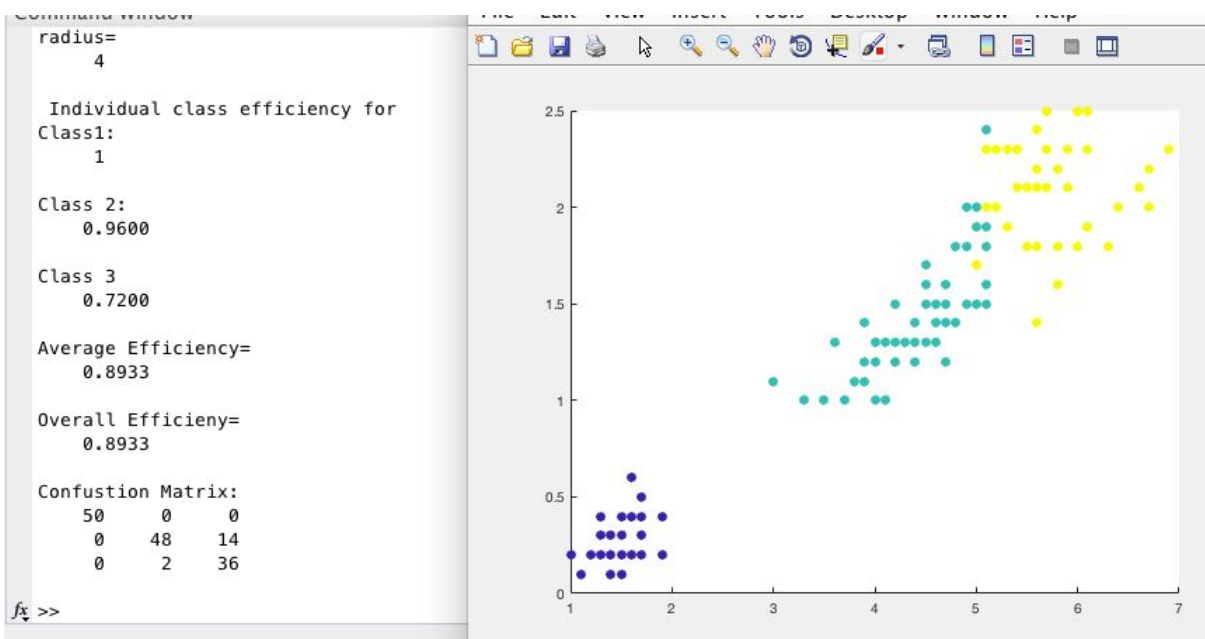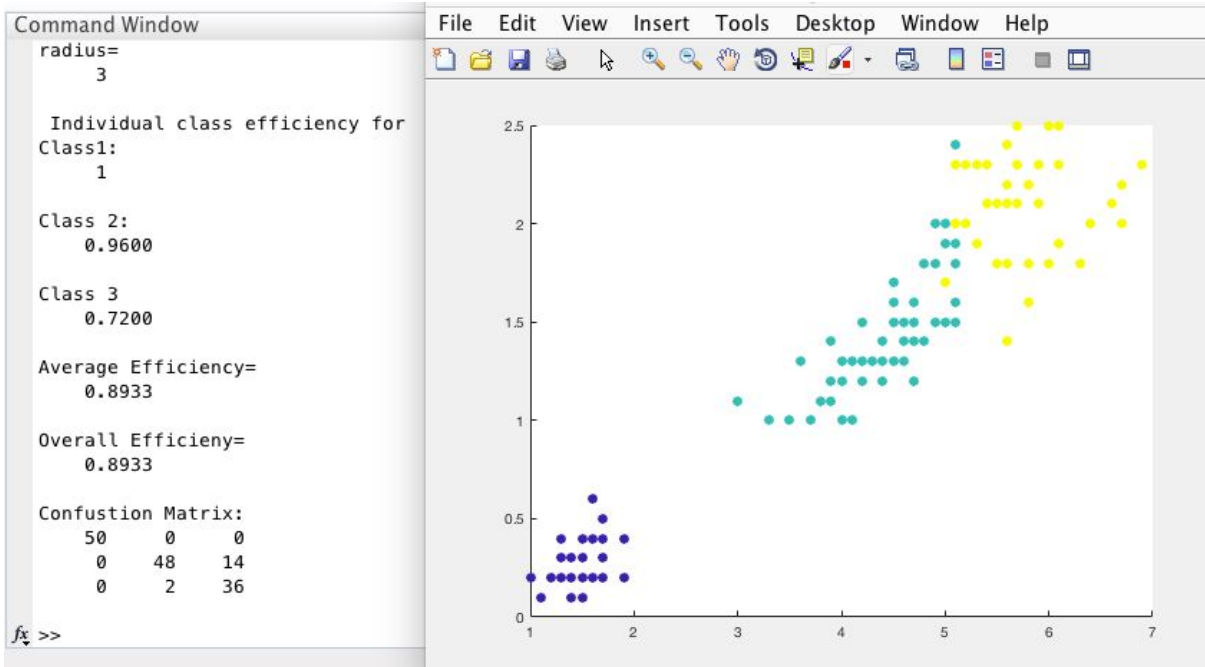
## Radius vs efficiency plot



- The following output was taken by keeping the number of iterations constant at 10.

| Radius | Average Efficiency |
|--------|--------------------|
| 0      | 0.4400             |
| 1      | 0.2400             |
| 2      | 0.2400             |
| 3      | 0.8933             |
| 4      | 0.8933             |
| 5      | 0.8533             |
| 10     | 0.8533             |

**Command Window (top):**

```
radius=
     0

 Individual class efficiency for
Class1:
     1

Class 2:
    0.0400

Class 3
    0.2800

Average Efficiency=
    0.4400

Overall Efficieny=
    0.4400

Confustion Matrix:
    50     0     0
     0     2    36
     0    48    14

fx >>
```



**Command Window (bottom):**

```
radius=
     1

 Individual class efficiency for
Class1:
     0

Class 2:
     0

Class 3
    0.7200

Average Efficiency=
    0.2400

Overall Efficieny=
    0.2400

Confustion Matrix:
     0    48    14
    50     0     0
     0     2    36

fx >>
```

Command Window
```
radius=
     5

 Individual class efficiency for
Class1:
     1

Class 2:
    0.8400

Class 3
    0.7200

Average Efficiency=
    0.8533

Overall Efficieny=
    0.8533

Confustion Matrix:
    50     4     0
     0    42    14
     0     4    36
```
fx >>



Command Window
```
radius=
    10

 Individual class efficiency for
Class1:
     1

Class 2:
    0.8400

Class 3
    0.7200

Average Efficiency=
    0.8533

Overall Efficieny=
    0.8533

Confustion Matrix:
    50     4     0
     0    42    14
     0     4    36
```
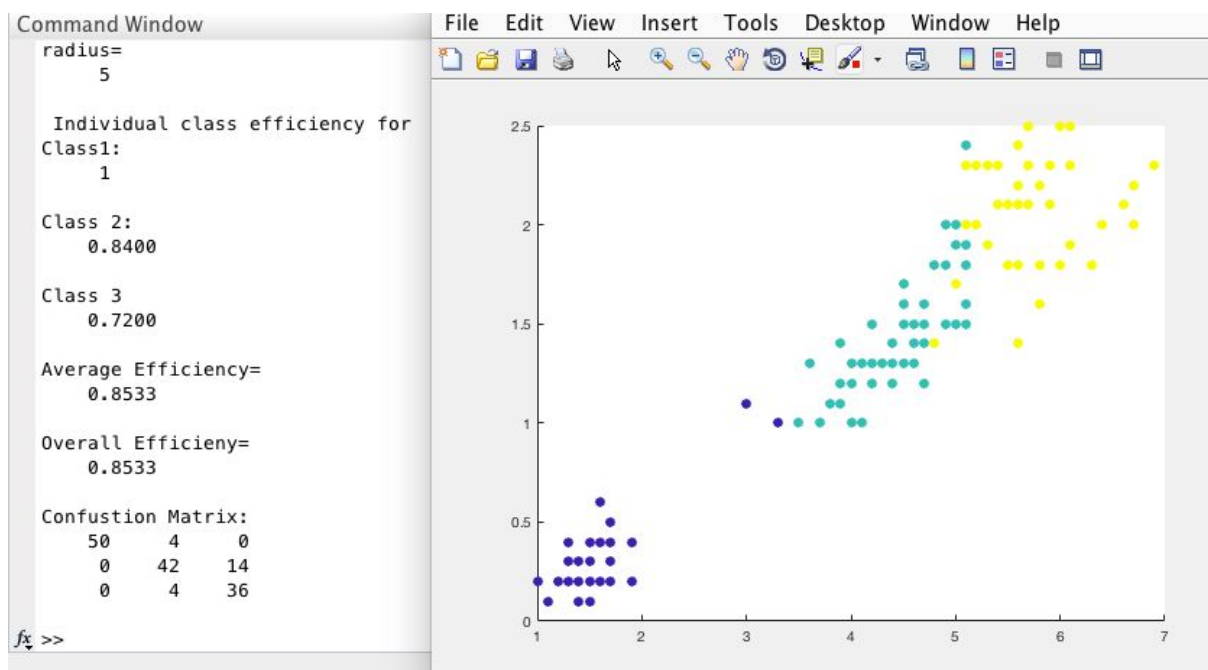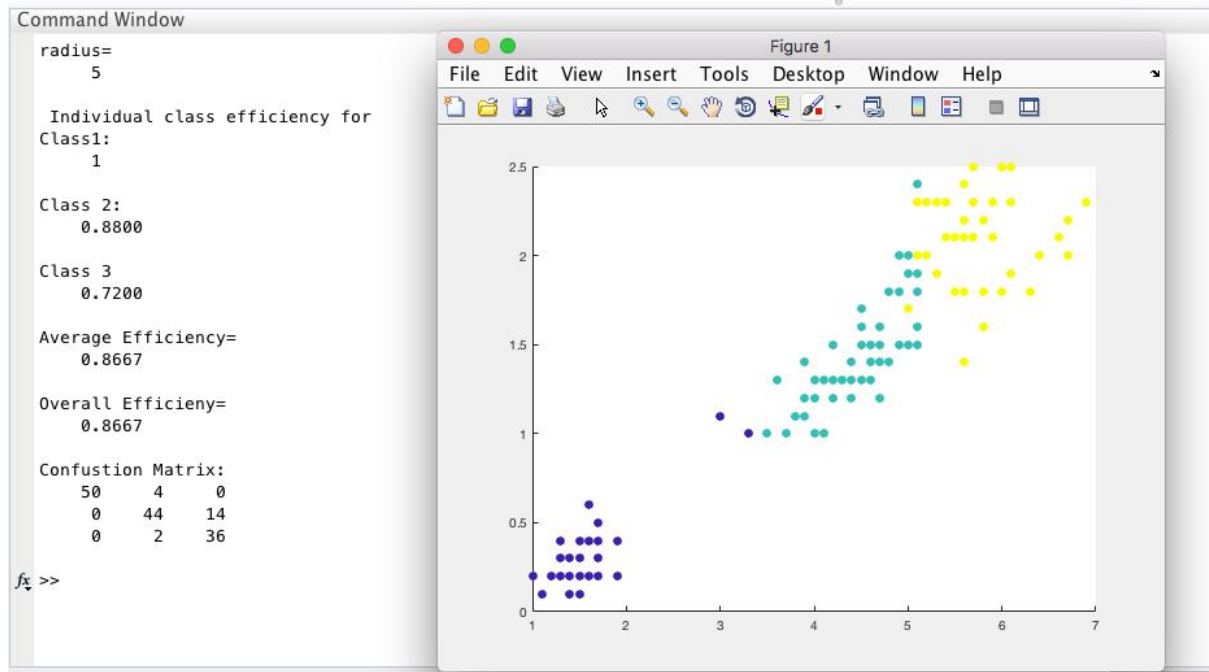fx >>

# Number Of Iterations

## Observation:

- We observed that by keeping the radius constant at 5 we get the maximum efficiency by iterating 100 times.On increasing the number of iterations further the efficiency remains the same.
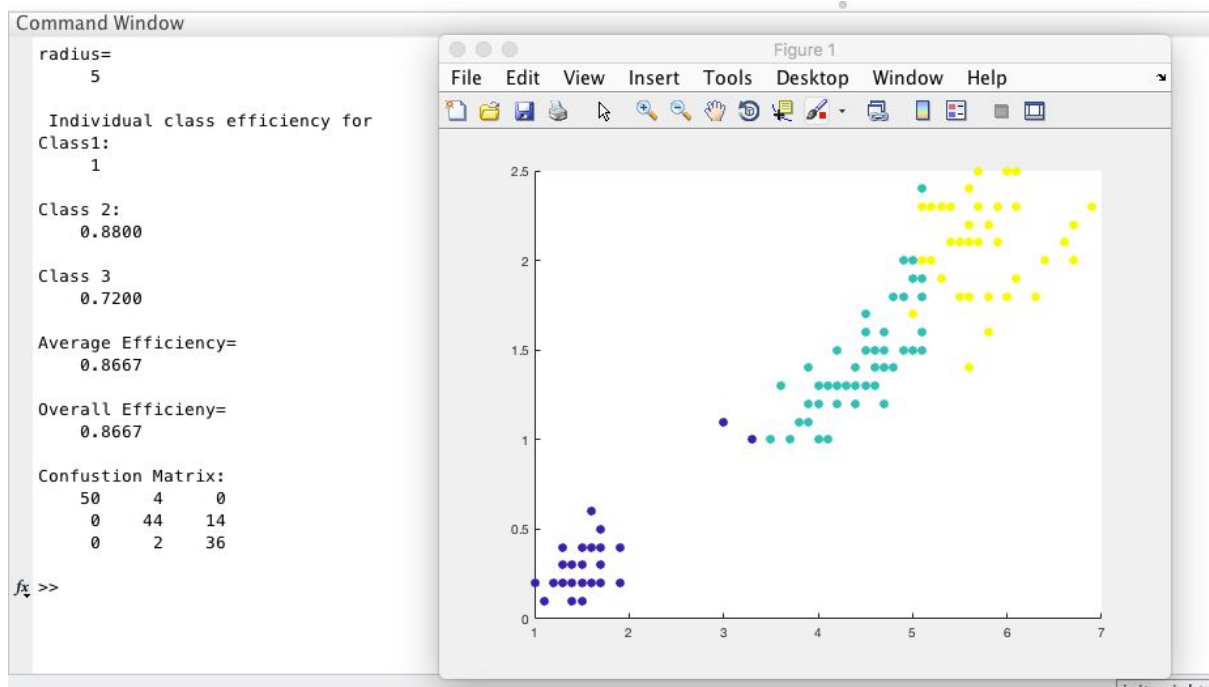
Number of iterations=10

Number of iterations=100
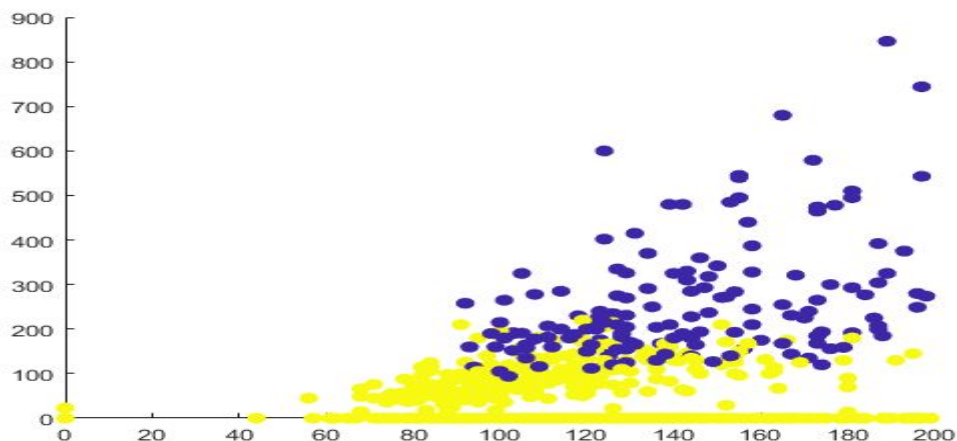


Number of iterations=500

# PIMA Diabetes data set

## Code:

```matlab
1 -     clc;
2 -     clear all;
3 -     close all;
4
5
6 -     data = load('diabetes.txt');
7
8       %number of nodes is 3
9 -     global nDim nNeurons
10 -    nNeurons =2;
11 -    nDim =8;
12 -    learning_rate =0.1;
13 -    n = size(data(:,1))
14 -    disp(n);
15 -    learning_rate =0.2;
16 -    weights=ones(nDim,nNeurons);
17 -    for i=1:nDim
18 -    maxi = max(data(:,i));
19 -    mini = min(data(:,i));
20 -    weights(i,:) = mini + (maxi-mini).*rand (1,nNeurons);
21 -    end
```

```matlab
40 -    nIterations =100;
41
42 -    for k=1:nIterations
43 -        for c=1:n
44
45 -            min_dist = 100000;
46 -            dist = zeros(nNeurons,1);
47 -            for i=1:nNeurons
48 -                for j=1:nDim
49 -                    dist(i) = dist(i)+((data(c,j) - weights(j,i))^ 2);
50 -                end
51 -                dist(i) = sqrt(dist(i));
52
53
54 -                if(min_dist>dist(i))
55 -                    indices(c,1)=i;
56 -                    min_dist=dist(i);
57 -                end
58
59 -            end
60
61 -            for j =1:nDim
62 -                weights(j,indices(c,1)) = weights(j,indices(c,1)) + learning_rate*(data(c,j)-weights(j,indices(c,1)));
66               % weights = updateweights(weights,nNeurons,nDim,data,c,learning_rate);
67 -            end
68 -        end
69 -    figure();
70 -    disp(weights);
71 -    scatter(data(:,2),data(:,5),50,indices,'filled');
72 -    ans = confusionmat(indices,data(:,9));
73 -    disp(ans);
```

**Confusion matrix**

| 426 | 198 |
|-----|-----|
| 74  | 70  |

# Mean-Shift Clustering

## Concept

- Mean Shift is a powerful and versatile non parametric iterative algorithm that can be used for lot of purposes like finding modes, clustering etc.
- Mean shift considers feature space as an empirical probability density function. If the input is a set of points, then Mean shift considers them as sampled from the underlying probability density function. If dense regions (or clusters) are present in the feature space, then they correspond to the mode (or local maxima) of the probability density function. We can also identify clusters associated with the given mode using Mean Shift.
- To estimate the density function of a random variable, we will be using the kernel density function. We have used the Gaussian kernel function here.

## Code:

- First we create a function for Mean shift clustering called the HGMeanShiftClustering.m whose parameters are the data matrix, the bandwidth size ,the kernel and the plotflag.
- The data matrix the input data that we have to cluster.
- The bandwidth size is the bandwidth we need to set for the mean shift clustering.
- There are several kinds of kernel that we can give as argument like the 'gaussian' and the 'flag' etc.
- The plot flag is the flag for whether we want to get a plot of the output or not.For our example the plot flag is 0.

```matlab
1     function [clustCent,data2cluster,cluster2dataCell] = HGMeanShiftCluster(dataPts,bandWidth,kernel,plotFlag);
2
3         %**** Initialize stuff ***
4  -      [numDim,numPts] = size(dataPts);
5  -      numClust = 0;
6  -      bandSq = bandWidth^2;
7  -      initPtInds = 1:numPts;
8  -      maxPos = max(dataPts,[],2); % biggest size in each dimension
9  -      minPos = min(dataPts,[],2); % smallest size in each dimension
10 -      boundBox = maxPos-minPos; % bounding box size
11 -      sizeSpace = norm(boundBox); % indicator of size of data space
12 -      stopThresh = 1e-3*bandWidth; % when mean has converged
13 -      clustCent = []; % center of clust
14 -      beenVisited= false(1,numPts); % track if a points been seen already
15 -      numInitPts = numPts; % number of points to posibaly use as initilization points
16 -      clusterVotes = zeros(1,numPts,'uint16'); % used to resolve conflicts on cluster membership
17 -      clustMembsCell = [];
18
19         %*** mean function with the chosen kernel ****
20 -      switch kernel
21 -      case 'flat' % flat kernel
22 -          kmean = @(x,dis) mean(x,2);
23 -      case 'gaussian' % approximated gaussian kernel
24 -          kmean = @(x,d) gaussfun(x,d,bandWidth);
25 -      otherwise
26 -          error('unknown kernel type');
27 -      end
28
29         while numInitPts
30 -          tempInd = ceil( (numInitPts-1e-6)*rand); % pick a random seed point
31 -          stInd = initPtInds(tempInd); % use this point as start of mean
32 -          myMean = dataPts(:,stInd);   % intilize mean to this points location
33 -          myMembers = []; % points that will get added to this cluster
34 -          thisClusterVotes = zeros(1,numPts,'uint16'); % used to resolve conflicts on cluster membership
35
36 -          while true %loop untill convergence
37 -              sqDistToAll = sum(bsxfun(@minus,myMean,dataPts).^2); % dist squared from mean to all points still active
38
39 -              inInds = find(sqDistToAll < bandSq); % points within bandWidth
40 -              thisClusterVotes(inInds) = thisClusterVotes(inInds)+1; % add a vote for all the in points belonging to this cl
41
42 -              myOldMean = myMean; % save the old mean
43 -              myMean = kmean(dataPts(:,inInds),sqrt(sqDistToAll(inInds))); % compute the new mean
44 -              myMembers = [myMembers inInds]; % add any point within bandWidth to the cluster
45 -              beenVisited(myMembers) = true; % mark that these points have been visited
46
47             %*** plot stuff ****
48 -              if plotFlag
```

```matlab
49 -                figure(12345),clf,hold on
50 -                if numDim == 2
51 -                    plot(dataPts(1,:),dataPts(2,:),'.')
52 -                    plot(dataPts(1,myMembers),dataPts(2,myMembers),'ys')
53 -                    plot(myMean(1),myMean(2),'go')
54 -                    plot(myOldMean(1),myOldMean(2),'rd')
55 -                    pause(0.1);
56 -                end
57 -            end

59             %**** if mean doesn't move much stop this cluster ***
60 -            if norm(myMean-myOldMean) < stopThresh
61                 %check for merge posibilities
62 -                mergeWith = 0;
63 -                for cN = 1:numClust
64 -                    distToOther = norm(myMean-clustCent(:,cN)); % distance to old clust max
65 -                    if distToOther < bandWidth/2 % if its within bandwidth/2 merge new and old
66 -                        mergeWith = cN;
67 -                        break;
68 -                    end
69 -                end

71 -                if mergeWith > 0 % something to merge
72 -                    nc = numel(myMembers); % current cluster's member number
```

```matlab
72 -            nc = numel(myMembers); % current cluster's member number
73 -            no = numel(clustMembsCell{mergeWith}); % old cluster's member number
74 -            nw = [nc;no]/(nc+no); % weights for merging mean
75 -            clustMembsCell{mergeWith} = unique([clustMembsCell{mergeWith},myMembers]);   %record which points inside
76 -            clustCent(:,mergeWith) = myMean*nw(1) + myOldMean*nw(2);
77 -            clusterVotes(mergeWith,:) = clusterVotes(mergeWith,:) + thisClusterVotes;    %add these votes to the merg
78 -        else % it's a new cluster
79 -            numClust = numClust+1; %increment clusters
80 -            clustCent(:,numClust) = myMean; %record the mean
81 -            clustMembsCell{numClust} = myMembers; %store my members
82 -            clusterVotes(numClust,:) = thisClusterVotes; % creates a new vote
83 -        end
```

CODE FOR THE KERNEL-

For calculation purpose we use Gaussian function.

```
10 -        ns = 1000; % resolution of guassian approximation
11 -        xs = linspace(0,bandWidth,ns+1); % approximate ticks
12 -        kfun = exp(-(xs.^2)/(2*bandWidth^2));
13 -        w = kfun(round(d/bandWidth*ns)+1);
14 -        w = w/sum(w); % normalise
15 -        out = sum( bsxfun(@times, x, w ), 2 );
16 -    └ end
```

# IRIS DATA

```
1 -    clc;
2 -    clear all;
3 -    close all;
4
5 -    load iris_dat.dat;
6 -    data = iris_dat;
7 -    data = transpose(data);
8 -    bandWidth =2;
9 -    kernel = 'gaussian';
10 -    [clustCent,data2cluster,cluster2dataCell] = HGMeanShiftCluster(data,bandWidth,kernel,0);
11 -    scatter(data(3,:),data(4,:),50,data2cluster,'filled');
12 -    hist(data2cluster);
13 -    |
```

# PIMA DIABETES DATASET

```
1 -    clc;
2 -    clear all;
3 -    close all;
4
5 -    load diabetes.txt;
6 -    data = diabetes;
7 -    data = transpose(data);
8 -    bandWidth =73;
9 -    kernel = 'gaussian';
10 -    [clustCent,data2cluster,cluster2dataCell] = HGMeanShiftCluster(data,bandWidth,kernel,0);
11 -    scatter(data(2,:),data(5,:),50,data2cluster,'filled');
12
13 -    disp(confusionmat(data2cluster,data(9,:)));
```

# Output

## 1. For iris data:

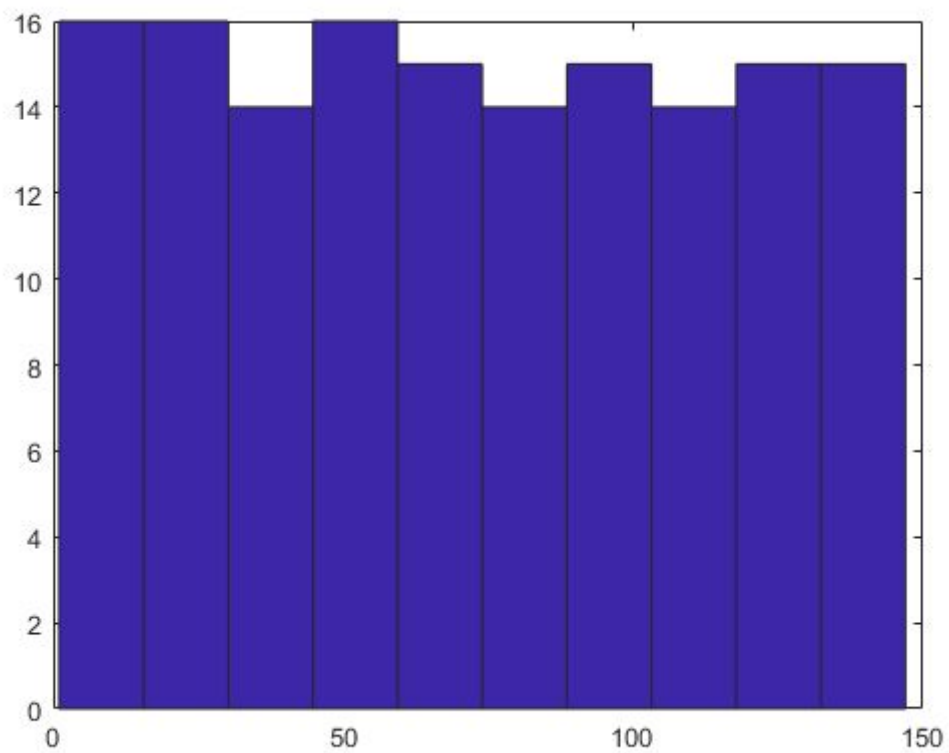### a. Bandwidth( h parameter=1)

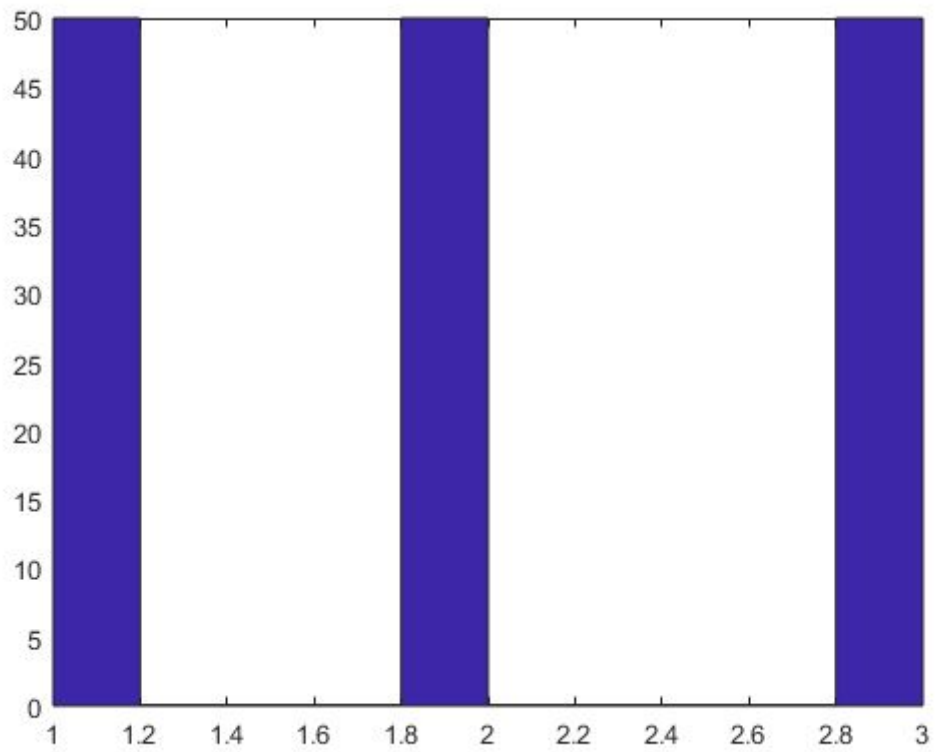**Histograms for Iris data:**

X- axis:num of clusters

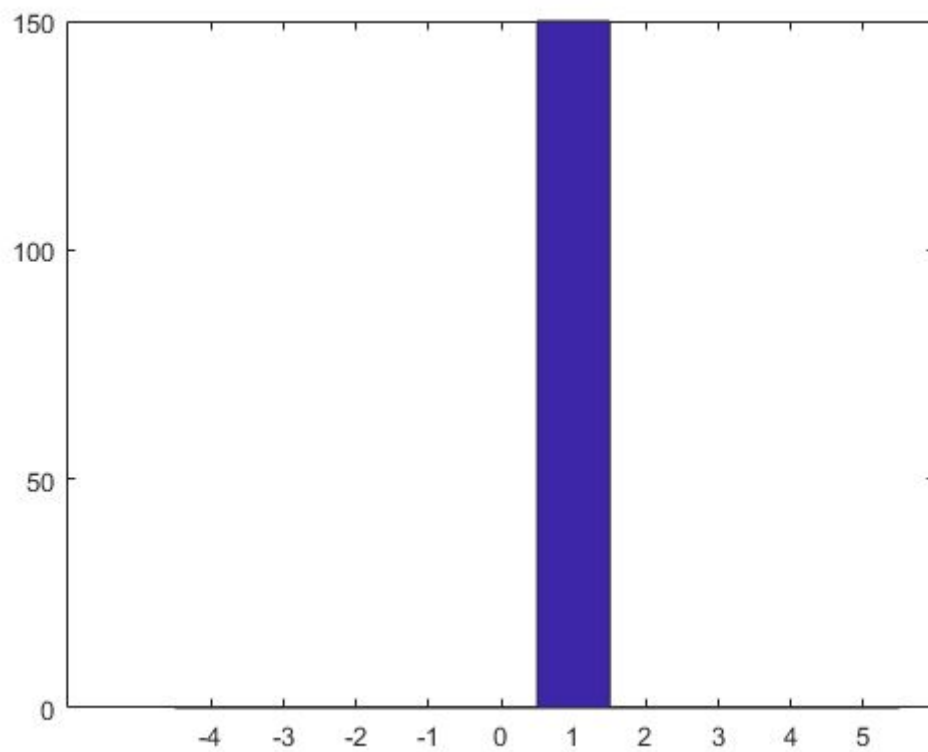Y-axis: num of data-sets per cluster

h= Bandwidth Parameter

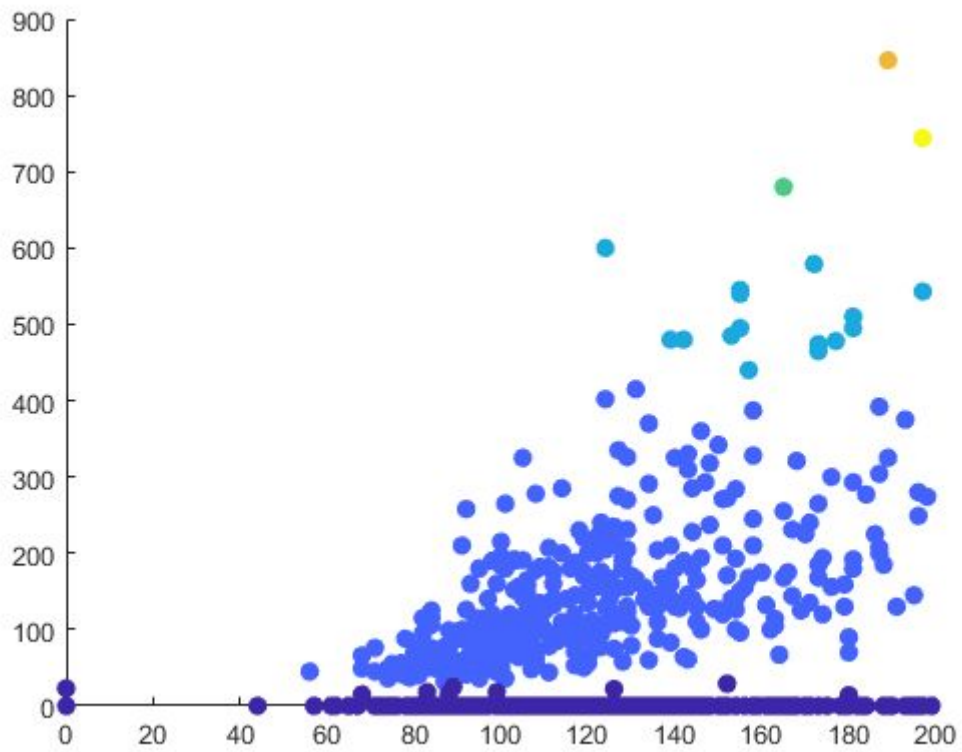# 1. Iris_overclustering h=0.1
# (No of Clusters=150)

## Normal Cluster -Iris Plot h=1
## (No of Clusters= 3)


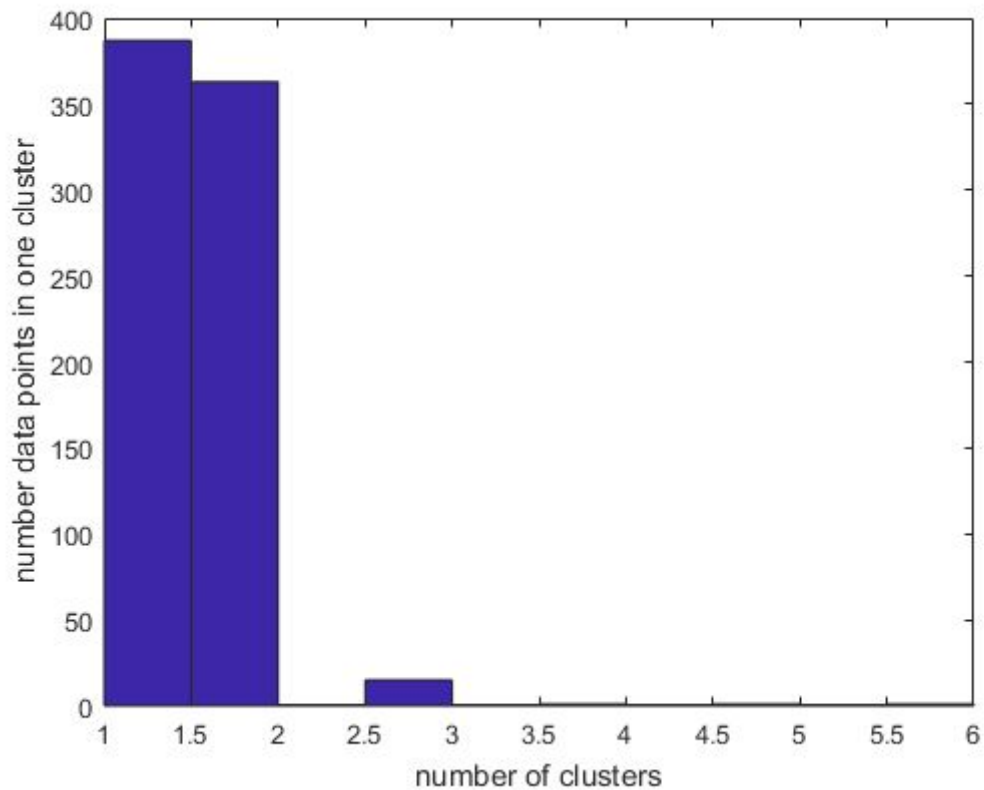
## Under Clustering h = 5
## (No of clusters=1)
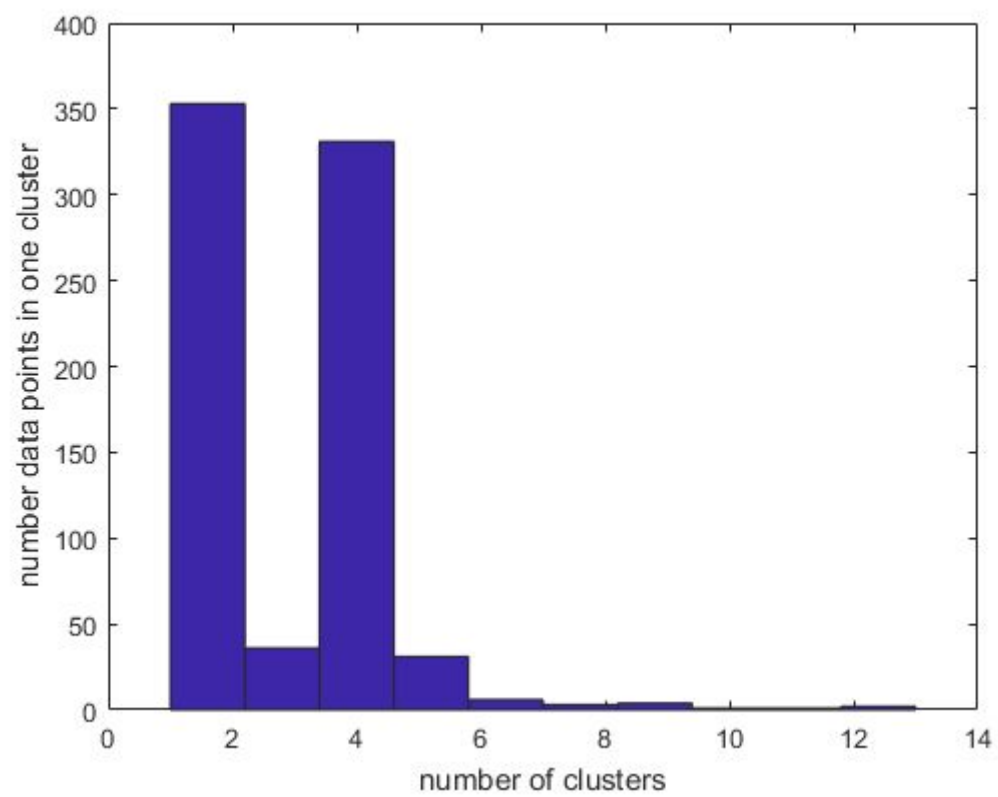
2. Diabetes Data-set
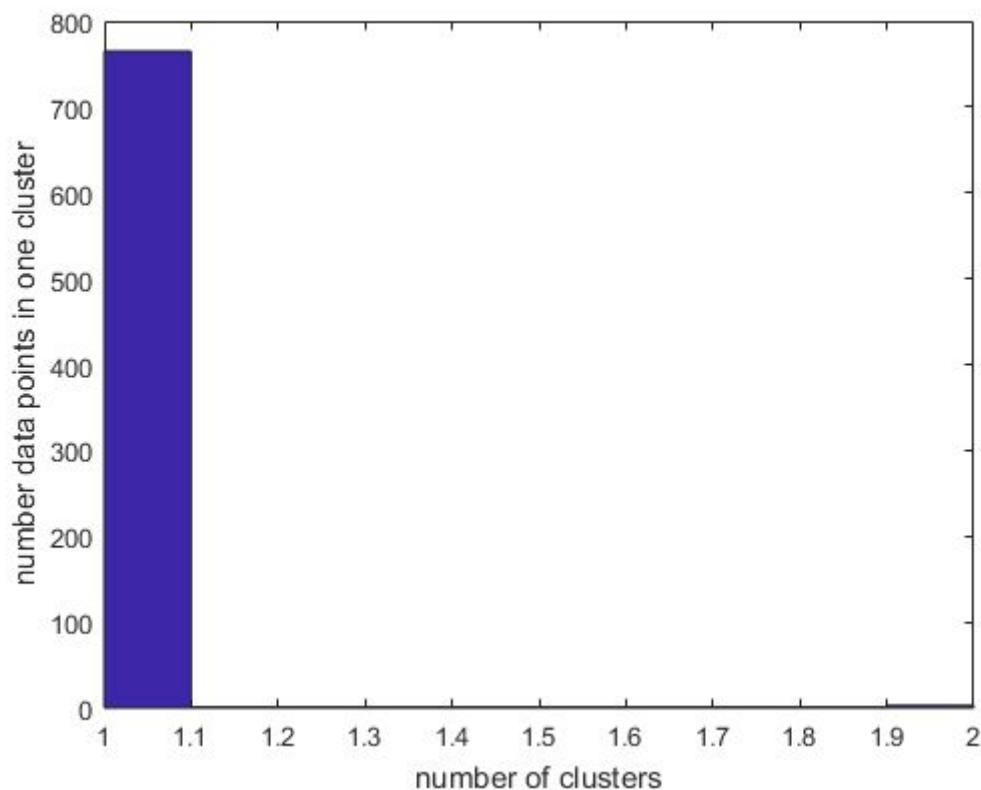
# Histogram for Diabetes Data-set:

## Normal Clustering h=72
## (No of Clusters=3)



## Over-clustering h=40
## (No of Clusters=12)

Under-clustering h=100
(No of Clusters=1)

## Observations:

- By tuning proper bandwidth parameter, we get:
- 3 clusters for Iris Data-set for bandwidth(h)=1
- 2 clusters( +1 misclassified) for Diabetes data-set for bandwidth(h)=72.

Changing the bandwidth leads to under-clustering and over-clustering.


## Analysis:

- We could see that the mean-shift algorithm is time intensive.The time complexity of mean-shift algorithm is O(T*n^2) where

    T= number of iterations

    n= number of data points in the data set

- As we increase the data points, time complexity increases.

As we saw in Iris-data set which has 150 data points and Diabetes data-set which has 768 data points, time taken to reach convergence is more in Diabetes data-set.

**Solution:** We searched some methods to make the mean shift algorithm to converge faster.

1.    The adaptive Mean Shift method is where you let the bandwidth parameter vary for each data point. The h parameter is calculated using Knn (k-nearest neighbour method) algorithm. Different distance metric methods are used in the Knn algorithm.

2.    An alternate way to speed up convergence is to alter the data points during the course of Mean Shift.

II.    Even though mean shift is a non-parametric algorithm, it does require the bandwidth parameter h to be tuned. We can use kNN to find out the bandwidth. The choice of bandwidth in influences convergence rate and the number of clusters. The ideal bandwidth which results in perfect number of clusters in case of Iris Data is 1 and in case of Diabetes data is 72.

A large h results in incorrect clustering and merges distinct clusters. A very small h results in too many clusters.

- **Dunn Index**

The Dunn Index assesses the goodness of a clustering, by measuring the maximal diameter of clusters and relating it to the minimal distance between clusters. This measure is quite conservative and prone to outliers, since it bases its calculation only on minimal and maximal distances.

$$D = \frac{\min\limits_{c_i \neq c_j \in C}\{d(c_i, c_j)\}}{\max\limits_{c_k \in C}\{d'(c_k)\}}$$

- **Davies–Bouldin index**

The **Davies–Bouldin index (DBI)** is a metric for evaluating clustering algorithms. This is an internal evaluation scheme, where the validation of how well the clustering

has been done is made using quantities and features inherent to the dataset. This has a drawback that a good value reported by this method does not imply the best information retrieval.

**Observation:**
- The very notion of "good clustering" is relative. The same applies for the quality measures. It depends on what you want to evaluate.
- For instance, the Davies-Bouldin Index evaluates intra-cluster similarity and inter-cluster differences. If you consider these to be good criteria, go for the Davies-Bouldin.

- Dunn index does not give precise output if there are many points.
- When there are many distinct or different groups output of DB index is better.
- DB index gives better output when there is noisy data
- Dunn index takes more time

## Applications of Mean-shift Algorithm:

i. The most important application is using Mean Shift for clustering. The fact that Mean Shift does not make assumptions about the number of clusters or the shape of the cluster makes it ideal for handling clusters of arbitrary shape and number.

ii. Mean Shift is used in multiple tasks in Computer Vision like segmentation, tracking, discontinuity preserving smoothing etc.

iii. Although, Mean Shift is primarily a mode finding algorithm, we can find clusters using it. The stationary points obtained via gradient ascent represent the modes of the density function. All points associated with the same stationary point belong to the same cluster.

## Comparison with K-means algorithm:

We studied K-means algorithm in the parametric clustering and based on this algorithm is ISO-DATA algorithm for non-parametric clustering.
We thus give a basic comparison between K-means algorithm and Mean-Shift
 One of the most important difference is that K-means makes two broad assumptions – the number of clusters is already known and the clusters are shaped spherically (or elliptically). Mean shift, being a non-parametric algorithm, does not assume anything about number of clusters. The number of modes give the number of clusters. Also, since it is based on density estimation, it can handle arbitrarily shaped clusters. K-means is very sensitive to initializations. A wrong initialization can delay convergence or sometimes even result in wrong clusters. Mean shift is fairly robust to initializations. Typically, mean shift is run for each point or sometimes points are selected uniformly from the feature space. Similarly, K-means is sensitive to outliers but Mean Shift is not very sensitive.  K-means is fast and has a time complexity $O(k*n*T)$, where k is the number of clusters, n is the number of points and T is the number of iterations. Classic mean shift is computationally expensive with a time complexity $O(T*n*n)$. Mean shift is sensitive to the selection of bandwidth, h. A small h can slow down the convergence. A large h can speed up convergence but might merge two modes. But still, there are many techniques to determine h reasonably well.