

DM - Assignment 04

Problem 1

1)

Formulating non-linear soft-margin SVM:

$$\min_{w,b,\xi} \frac{1}{2} \|w^2\| + C \sum_{i=1}^N \xi_i$$

SVM with $\phi(x)$:

$$\begin{aligned} \text{s.t. } \forall i \in [1, N], \quad y_i(\langle w, \phi(x_i) \rangle + b) &\geq 1 - \xi_i \\ \forall i \in [1, N], \quad \xi_i &\geq 0. \end{aligned}$$

↵

Quadratic Programming & Dual Function:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \sum_{i=1}^N \alpha_i$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0, \text{ where } \alpha_i \in [0, C], i = 1, 2, \dots, N.$$

↵

$$f(x) = \text{sign} \left(\sum_{i=1}^N \alpha_i y_i \phi(x_i) \phi(x) + b \right) = \text{sign} \left(\sum_{i=1}^N \alpha_i y_i K(x_i, x) + b \right).$$

2)

1. Given $K_{i,j} = K(x_i, x_j)$, we only need to find a mapping $\phi(x): X \rightarrow H$, s.t. for any x, z belongs to X , $k(x, z) = \phi(x) \cdot \phi(z)$; saying one kernel matrix, we can find multiple such mappings out of it, and feature vectors w and bias b can thus be calculated with that.

2. No. Time and space complexity increases by calculating feature vectors comparing to merely calculate the kernel matrix.

Problem 2

1) Complexity of k-means & k-medoids (PAM)

Time requires compute the distance between two data points is linear to their dimensions, that is a complexity of $O(n)$. Saying the number of iterations of either algorithm is i .

K-means steps are like below:

1. Initialization: choose k random centroids from X . $\rightarrow O(k)$
2. Repeat the following steps:
 - Assignment: perform 1-NN classification to (re)assign each $x \in X$ to one of the cluster centroids. $\rightarrow O(Nn)$
 - Update: recompute each c_j to be the centroid of all the points assigned to its cluster. $\rightarrow O(Nn)$

Until the centroids/clusters stabilize (i.e., do not change)

So the total complexity of k-means is $O(iknN)$.

For each iteration, the PAM steps are like below:

1. Initialize: select k of the N data points as the medoids. $\rightarrow O(k)$
2. Associate each data point to the closest medoid. $\rightarrow O(Nkn)$
3. While the cost of the configuration decreases:
 - For each medoid m , for each non-medoid data point o : \rightarrow multiply by $k(N-k)$
 - Swap m and o , recompute the cost (sum of distances of points to their medoid) $\rightarrow O((N-k)n)$ cuz medoids need not to be computed
 - If the total cost of the configuration increased in the previous step, undo the swap $\rightarrow O(C)$

So the total complexity of PAM is $i * (O(k) + O(Nkn) + k(N-k) * O((N-k)n)) = O(ikn(N-k)^2)$.

2) Shake & Bake Disimilarity

1. When $D(x,y) = L^{-1} \sum_{j=1}^L D_j(x,y) = 0$, because $D_j \in \{0, 0.5, 1\}$, so for every j , D_j has to be 0, in that condition $x=y$.
2. $D(x,y)$ is clearly symmetrical regarding x and y , $D(x,y) = D(y,x)$ stands.
3. Pick some $z \in X$, if z equals to x or y , then $D(x,y) = D(x,z) + D(z,y)$; otherwise, for each j , we have the

following conditions:

$C(x)=C(y)$	$C(y)=C(z)$	$C(x)=C(z)$	$D(x,y)$	$D(y,z)$	$D(x,z)$
T	T	T	0.5	0.5	0.5
T	F	F	0.5	1	1
F	T	F	1	0.5	1
F	F	T	1	1	0.5
F	F	F	1	1	1

In all above cases, $D(x,y) \leq D(x,z) + D(y,z)$, so triangle rule stands.

Thus, by 1-3 we have proved shake & bake dissimilarity is a distance metric.

3) Centroid by Attribute-wise Mode

The sum of absolute error (SAE) of a cluster w.r.t. mismatch distance $d(x,y) = \#\{x[j] \neq y[j]\}$ ($j \in [1,n]$ where n is the dimensionality of points) can be expressed as:

$$\sum_j \sum_i \text{int}(x(i,j) \neq c(j))$$

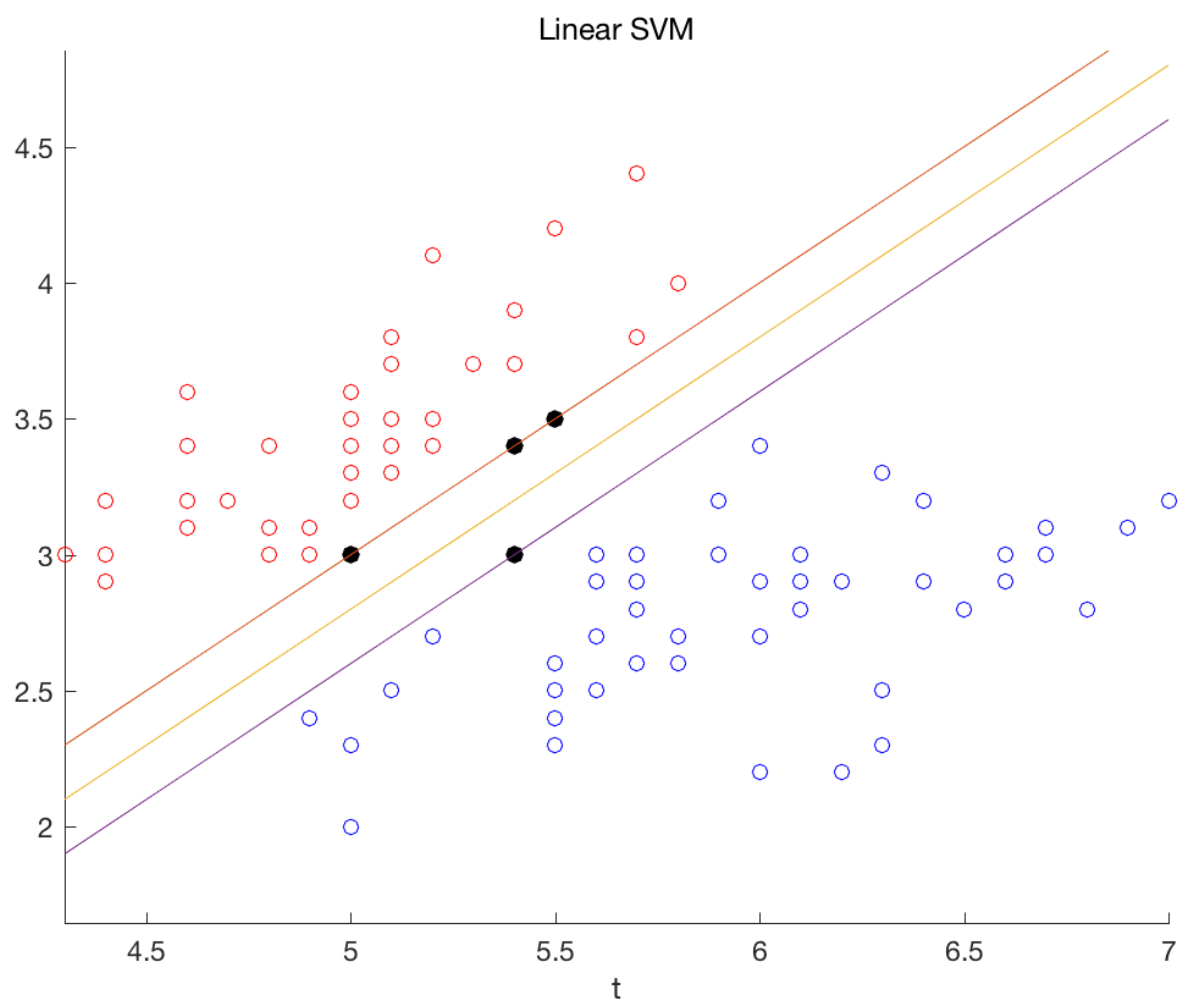
where $i \in [1,N]$, N is number of points in this cluster, and C is the centroid ($n \times 1$). It is clear that each attribute (j) here is independent, making it j different minimization problem, each one formulated as:

$$\sum_i \text{int}(x(i,j) \neq c(j))$$

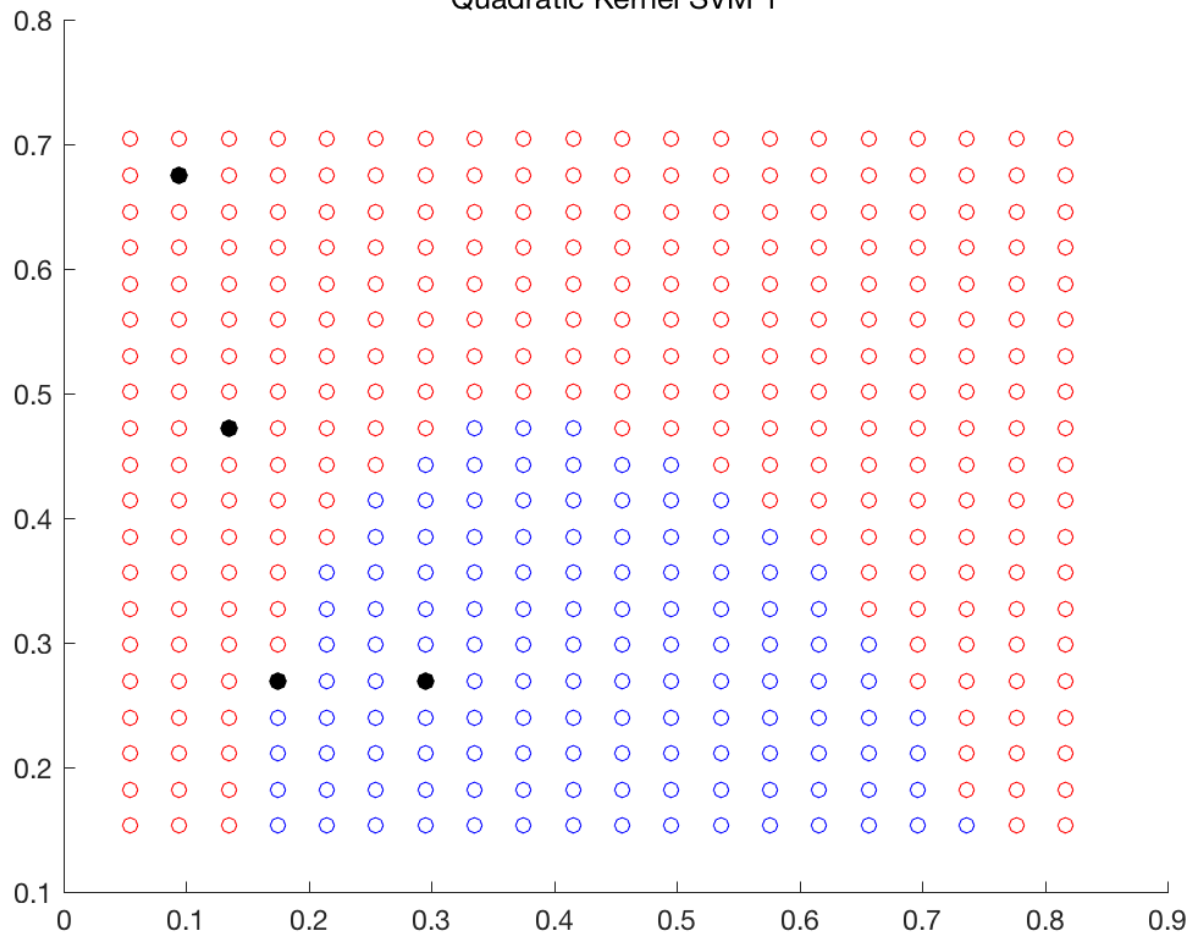
To make it smaller in each j , we should set as many $X(i,j)$ equals to $C(j)$, which means setting $C(j) = \text{mode}(X(:,j))$.

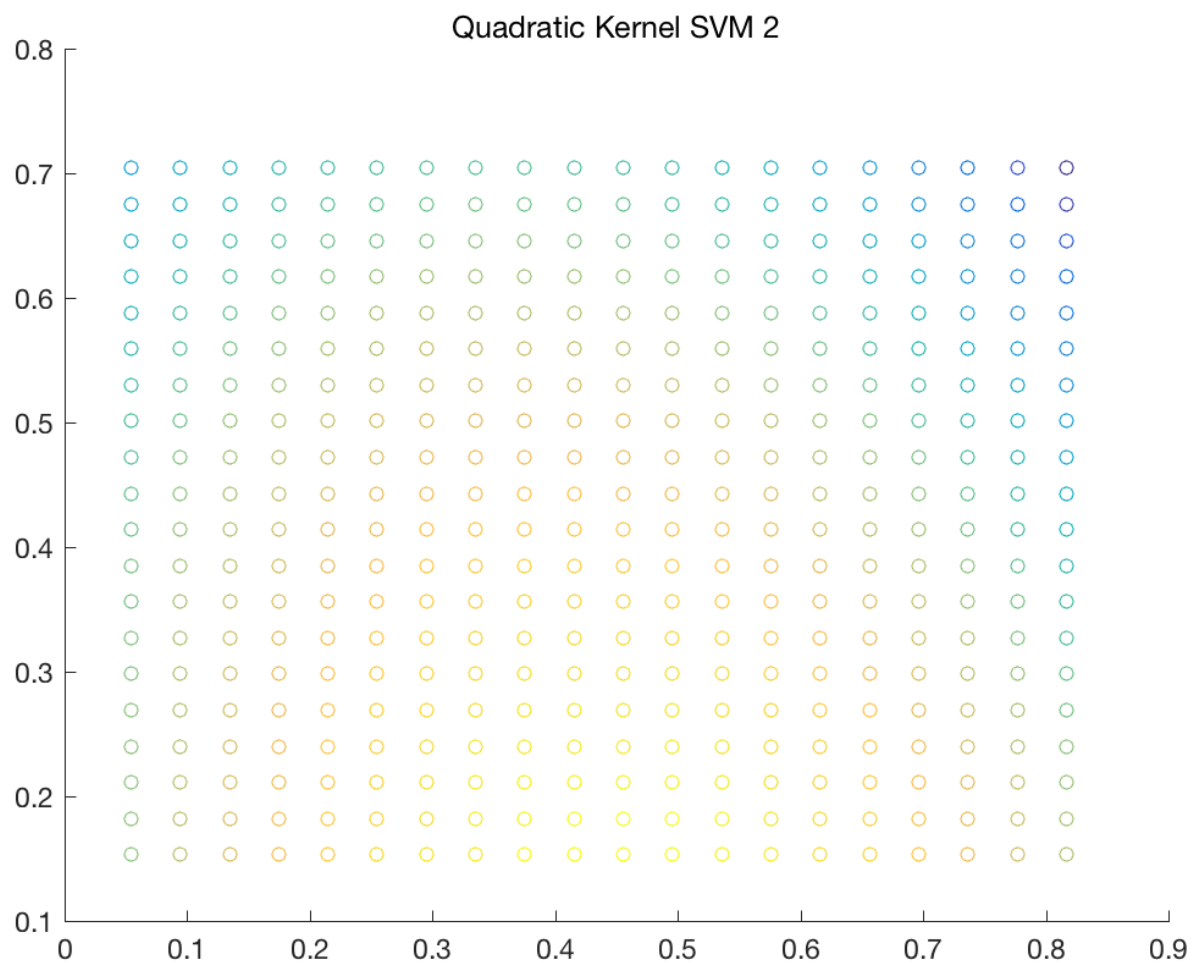
Thus the attribute-wise mode is indeed the centroid w.r.t the mismatch dissimilarity in terms of minimizing SAE in each cluster.

Problem 3

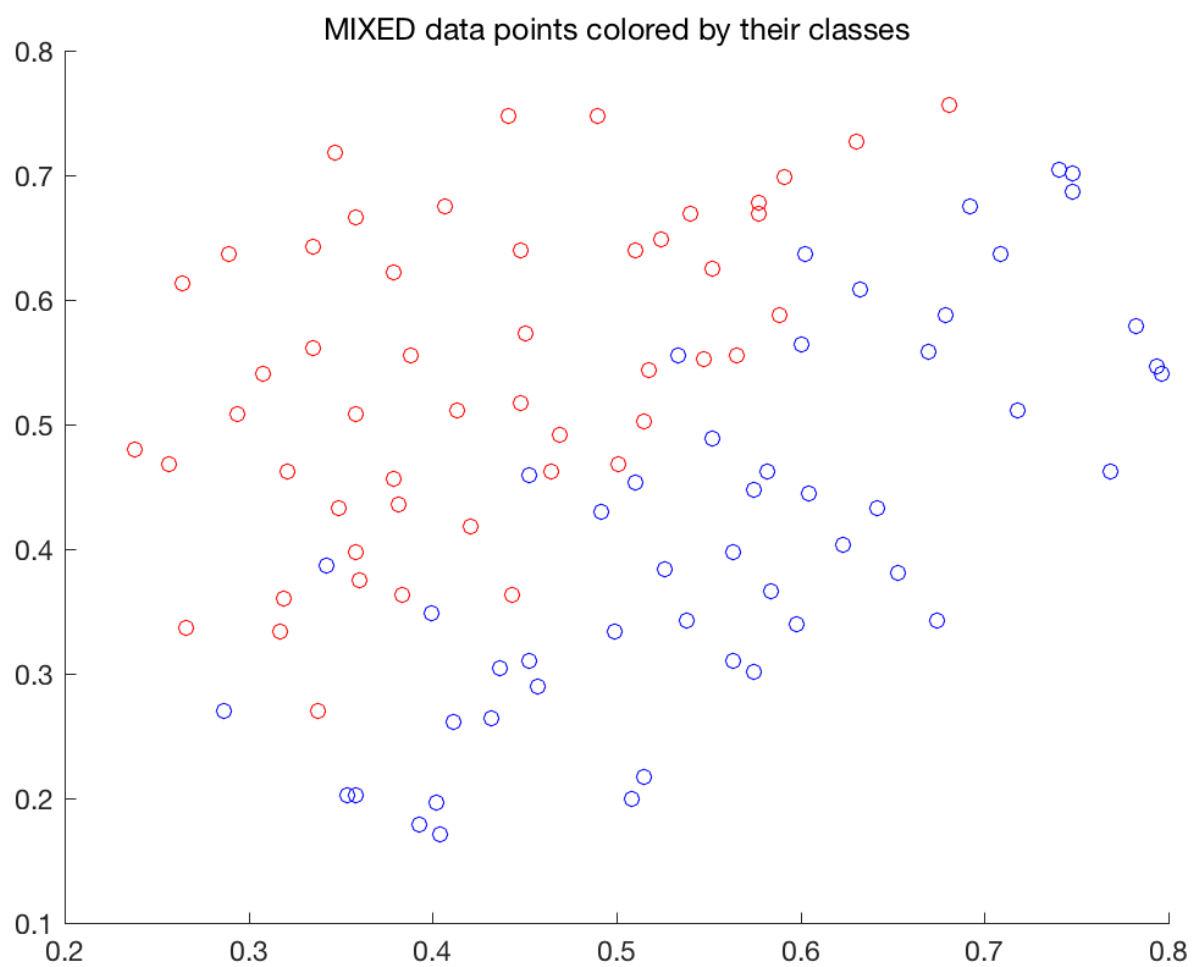


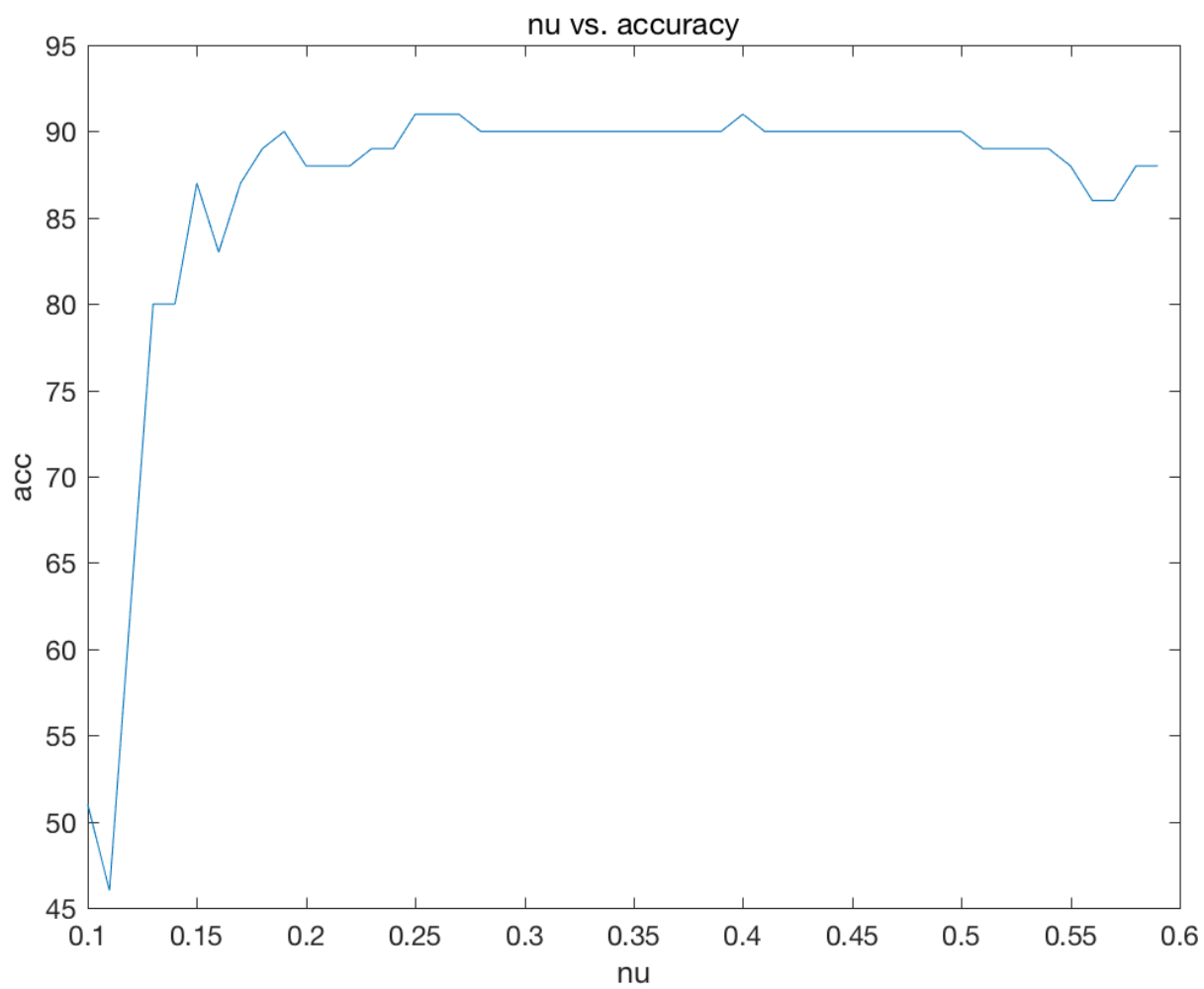
Quadratic Kernel SVM 1

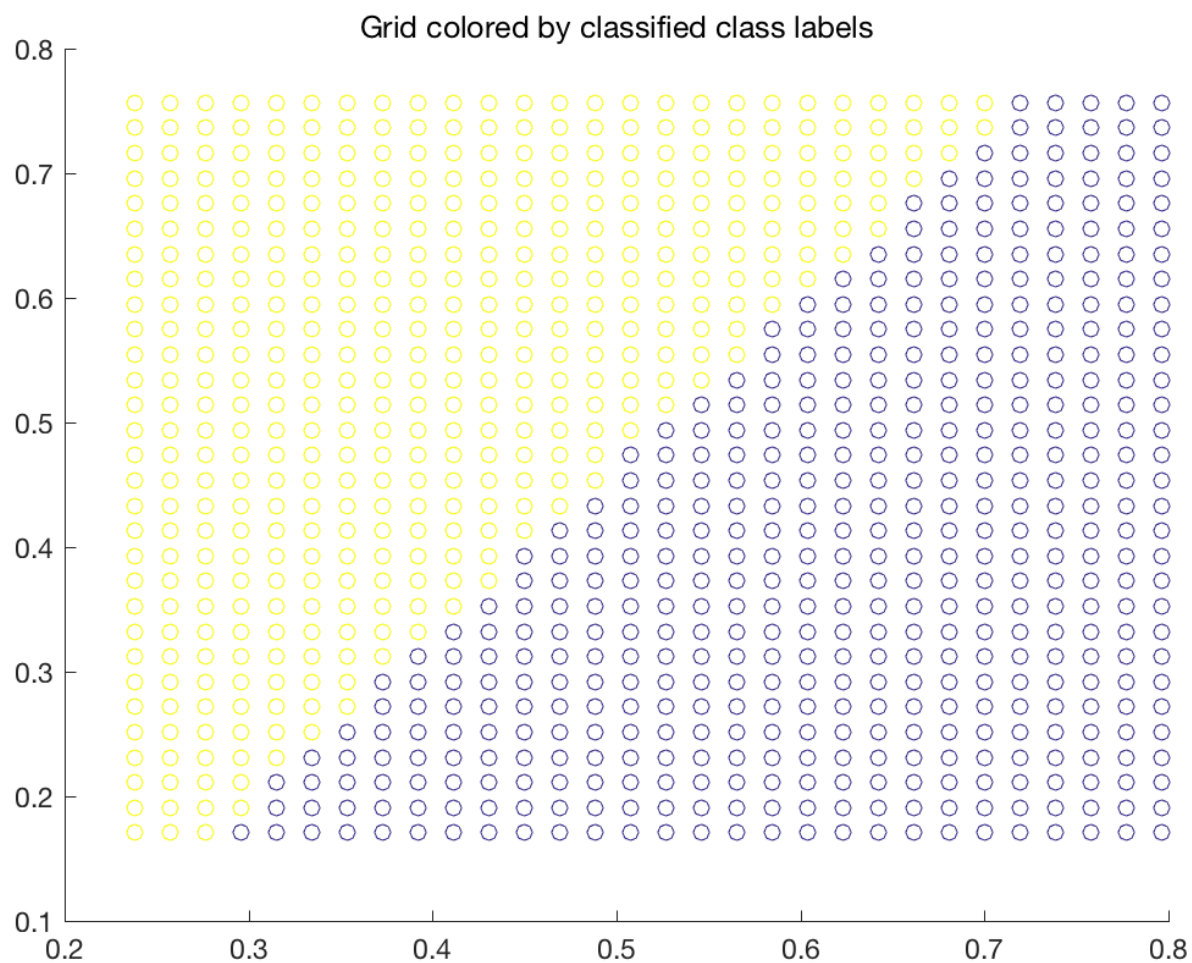


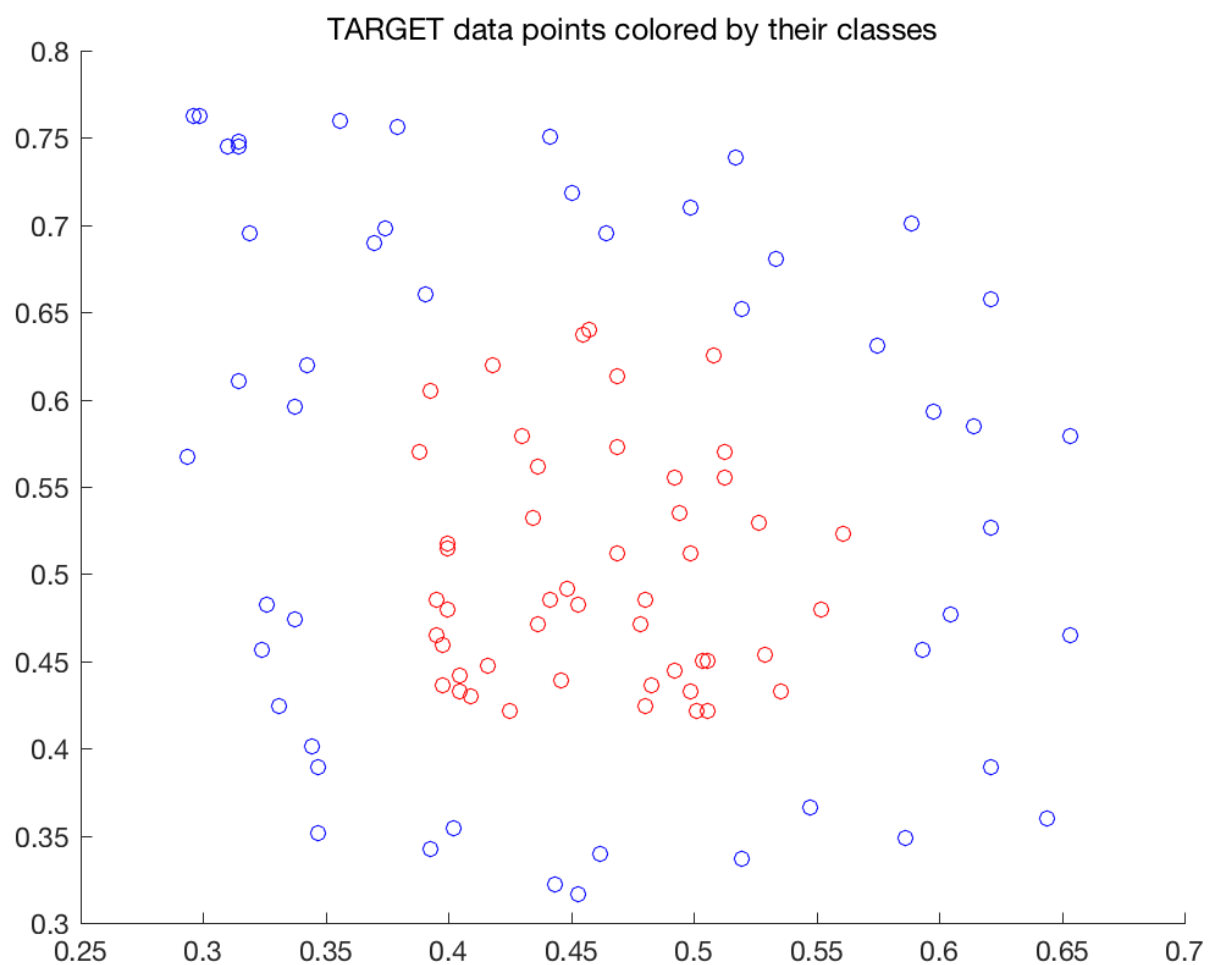


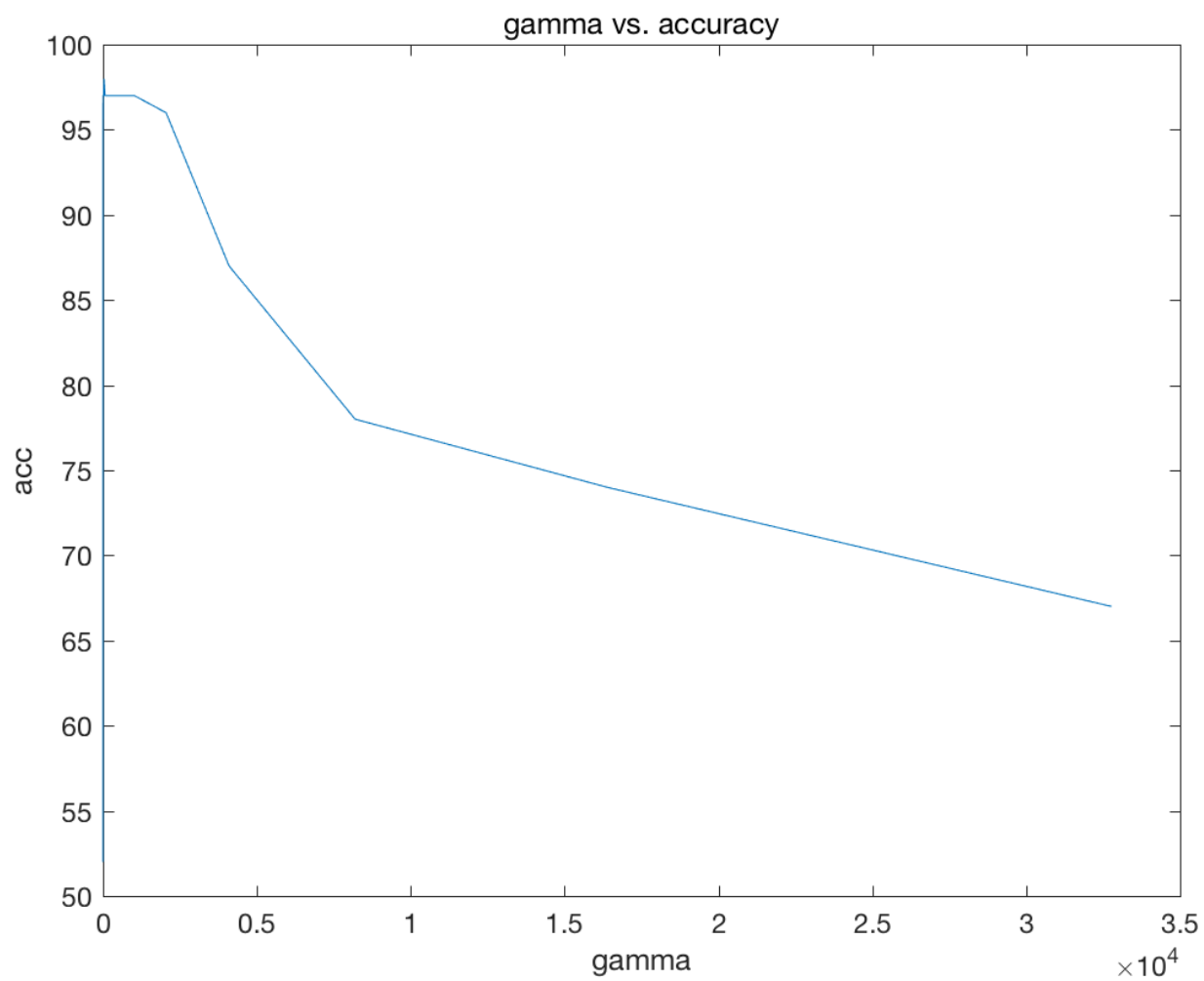
Problem 4







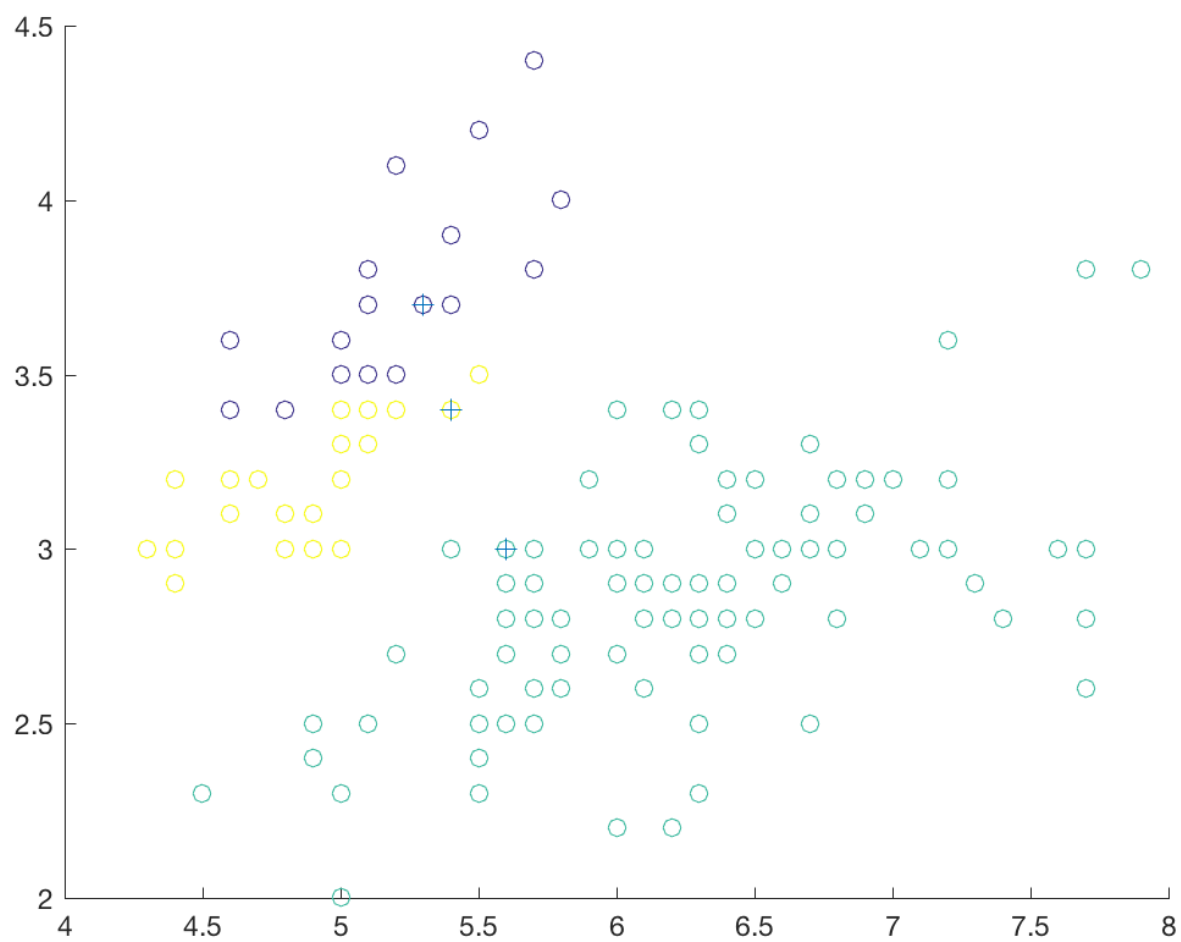


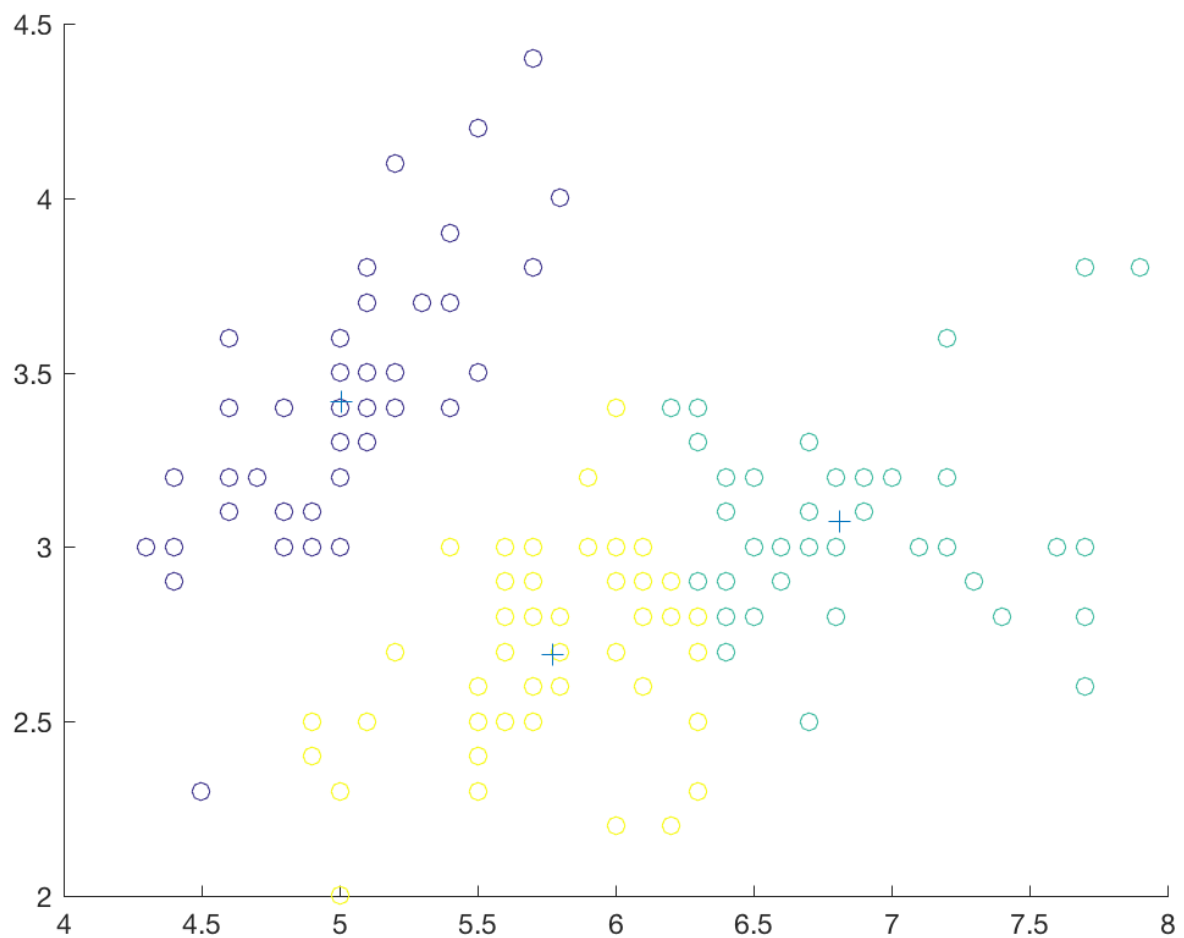




Problem 5

Illustrate clustering on a two-dimensional dataset

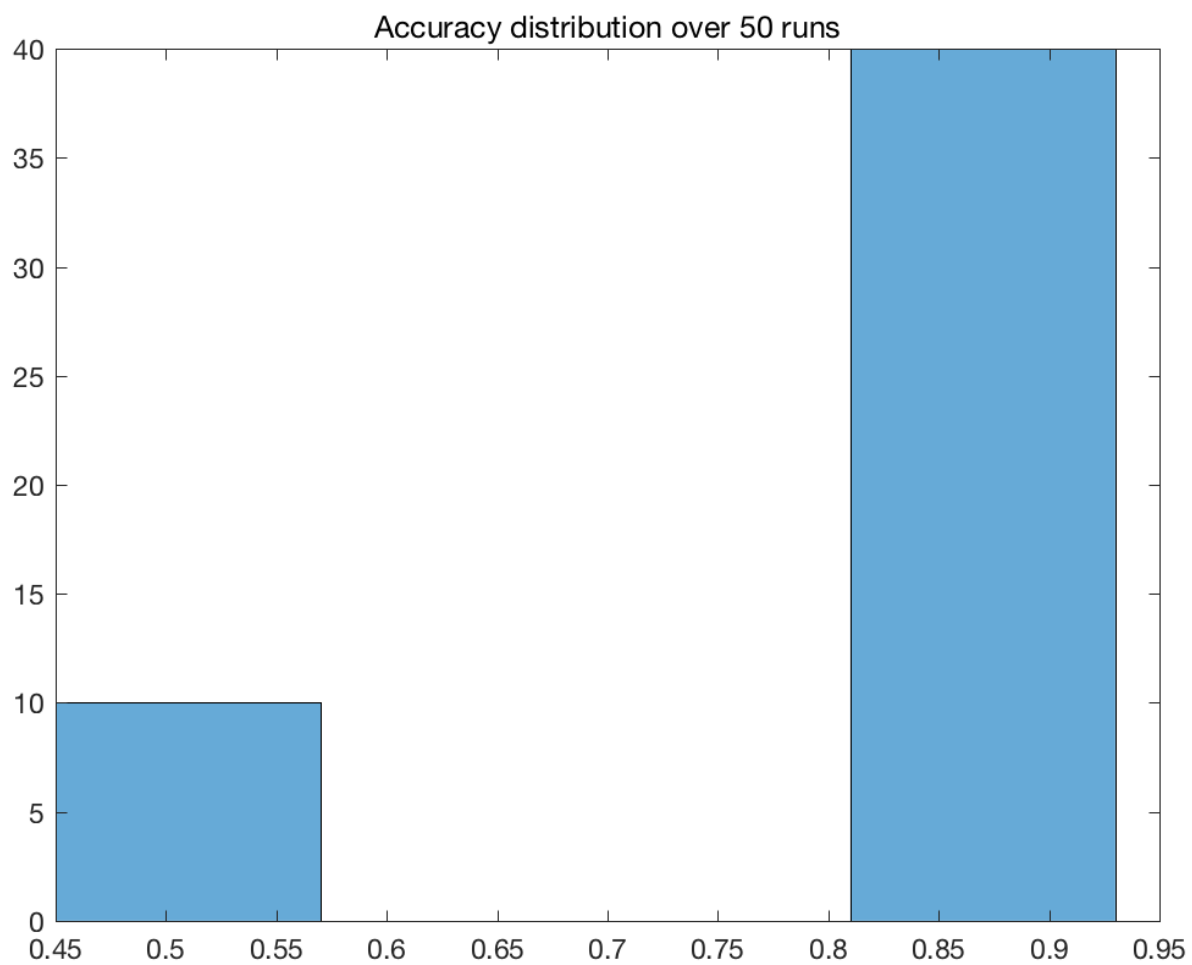


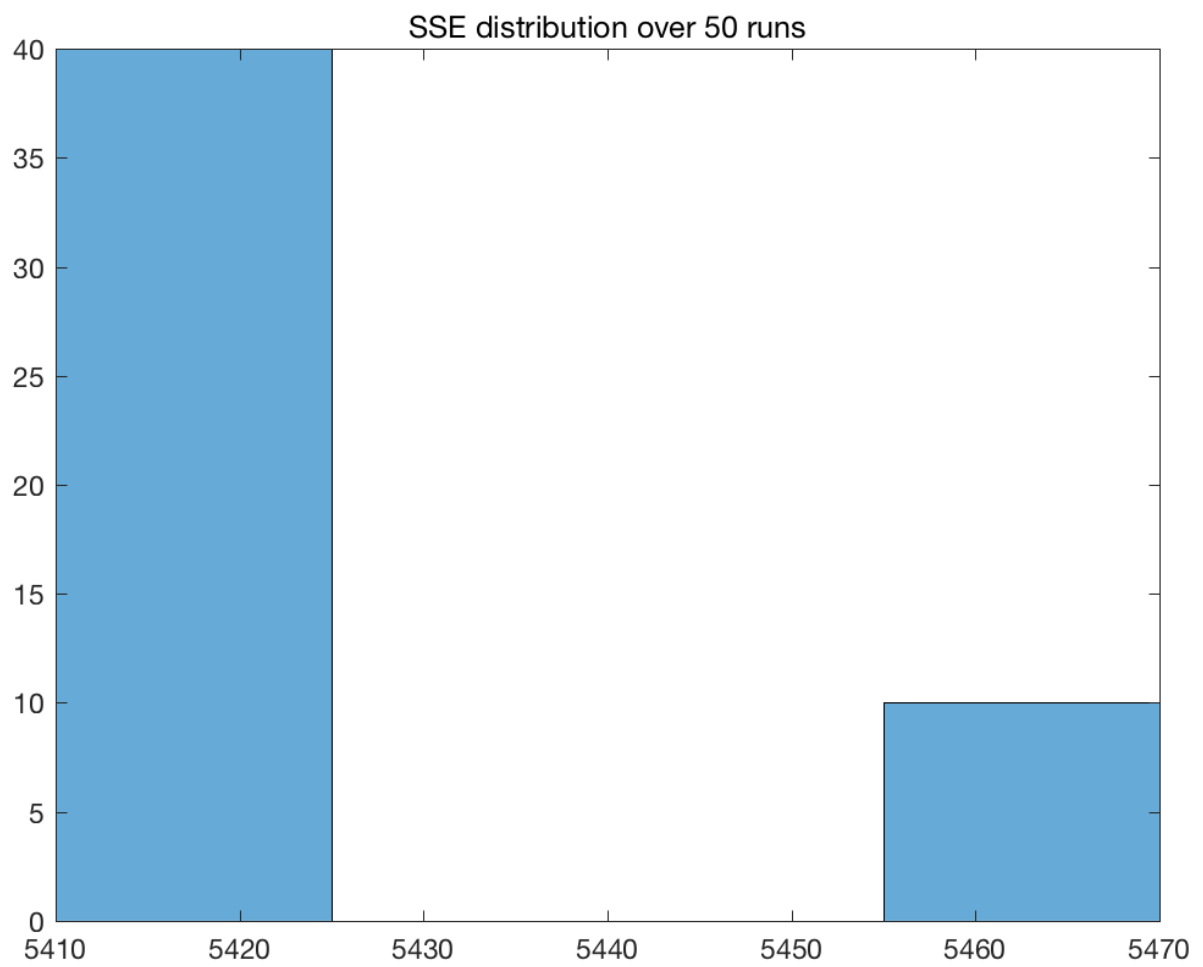


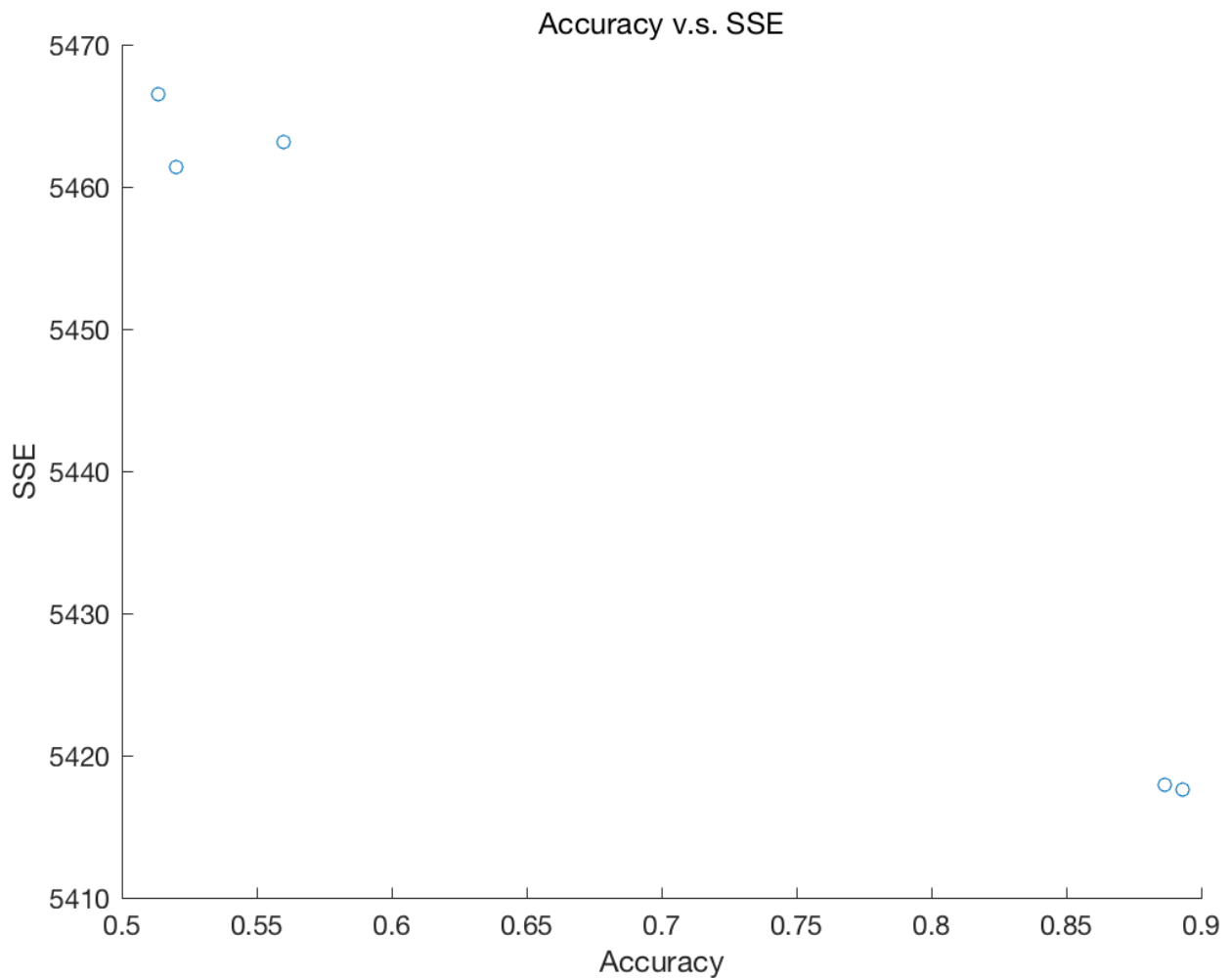
Test clustering quality on the full iris data

- Three accuracy numbers: 0.6667, 0.8867, 0.8867;
- Three RandIndex numbers: 0.8737, 0.8797, 0.8797.

Compare accuracy, RandIndex, and SSE as cluster quality measures





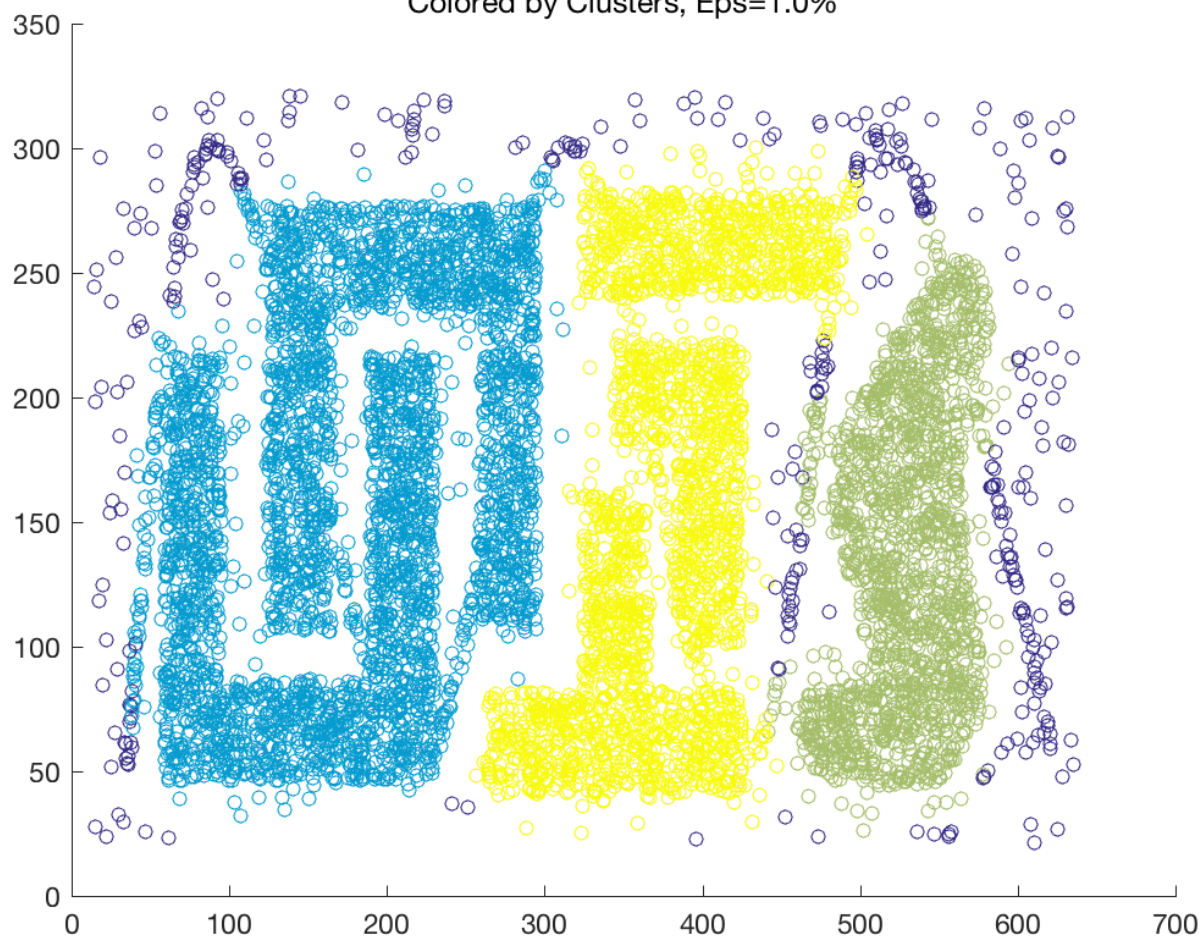


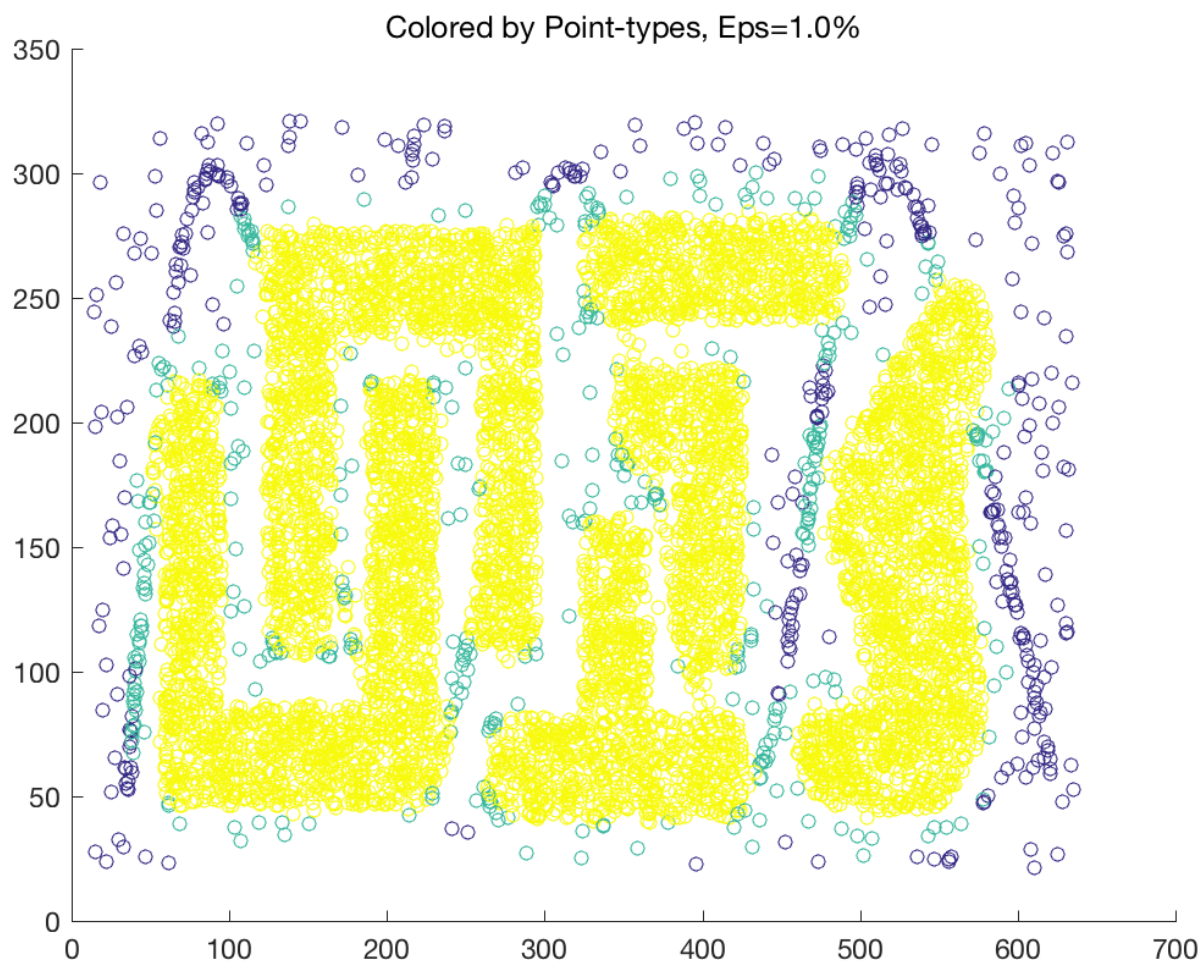
Conclusions

SSE is as good a metric as accuracy in measuring performance of unsupervised learning problems such as clustering. When SSE is low, the points tend to gather close to their centroids, forming a well-separated shape of clusters, leading to high accuracy, and often vice versa. SSE is even better to use when there are no actual labels provided, under what circumstances accuracy cannot be calculated, but we can still use SSE to obtain the concept of how points congregate (loose or tight, etc.).

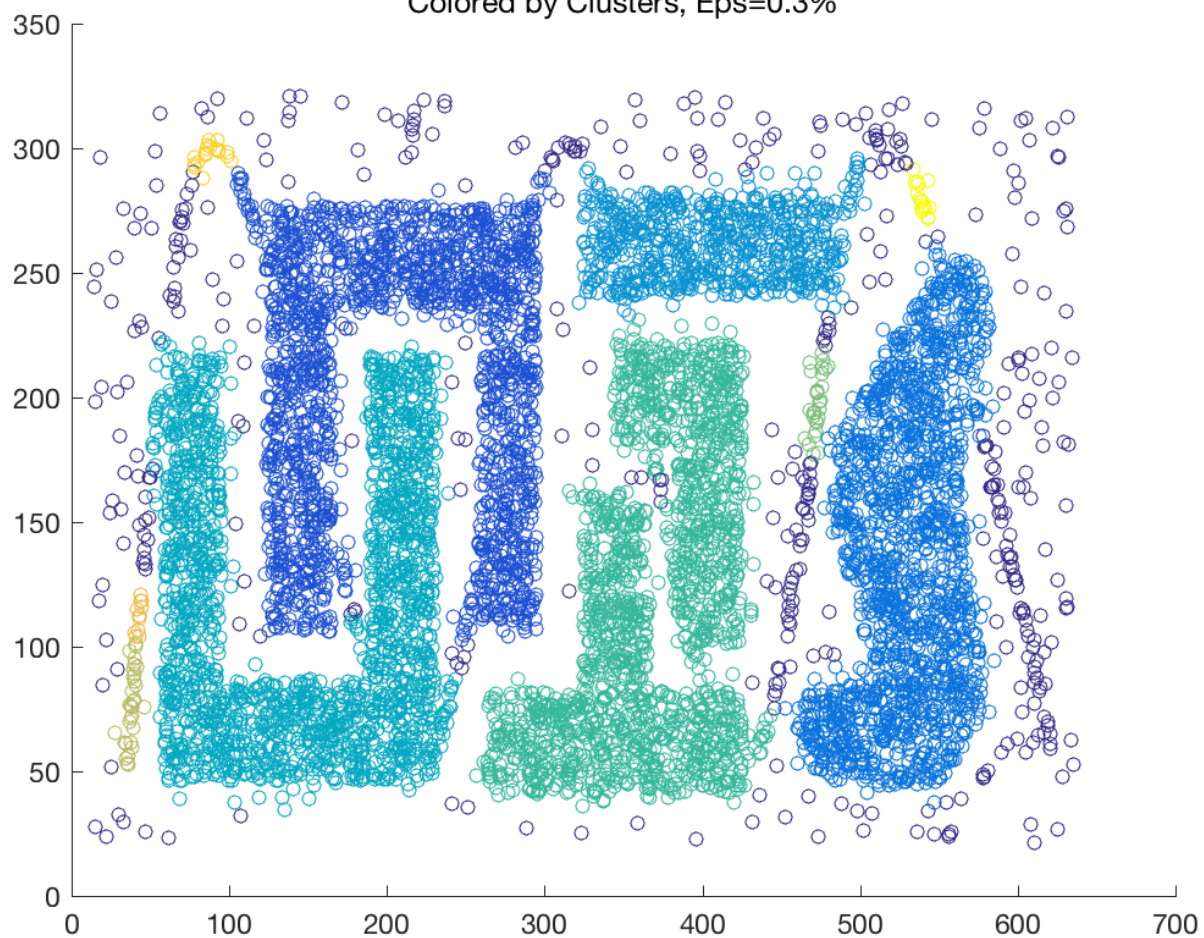
Problem 6

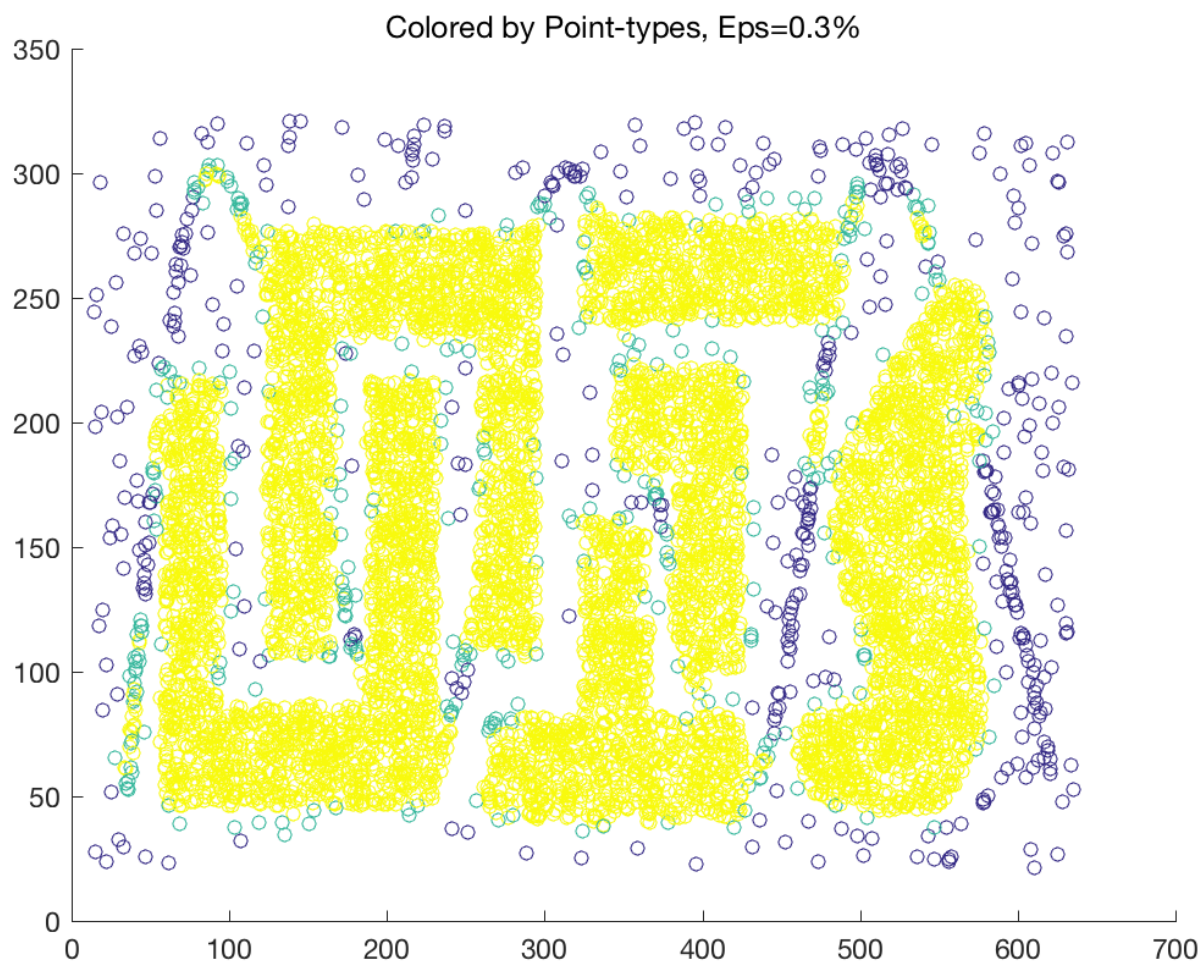
Colored by Clusters, Eps=1.0%



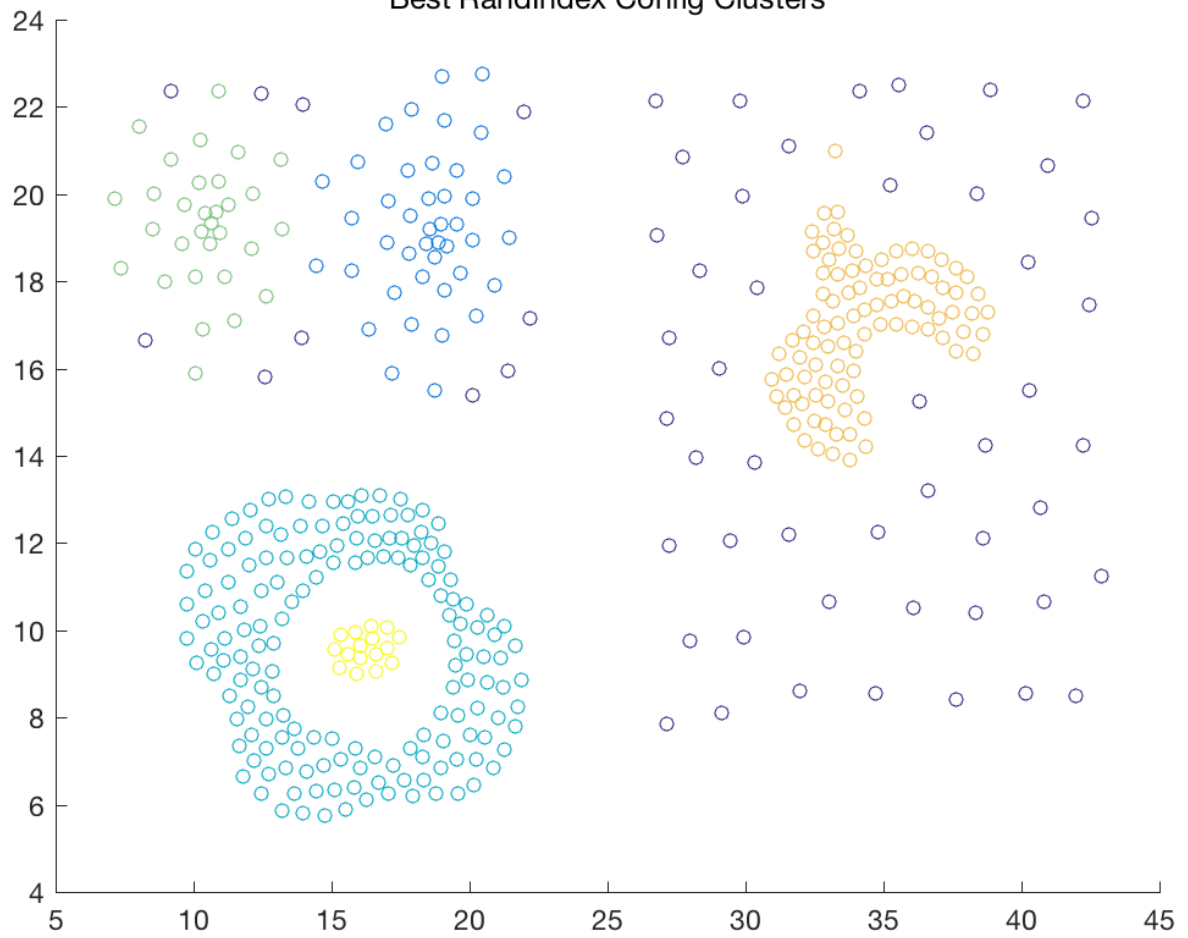


Colored by Clusters, Eps=0.3%





Best RandIndex Config Clusters



Worst RandIndex Config Clusters

