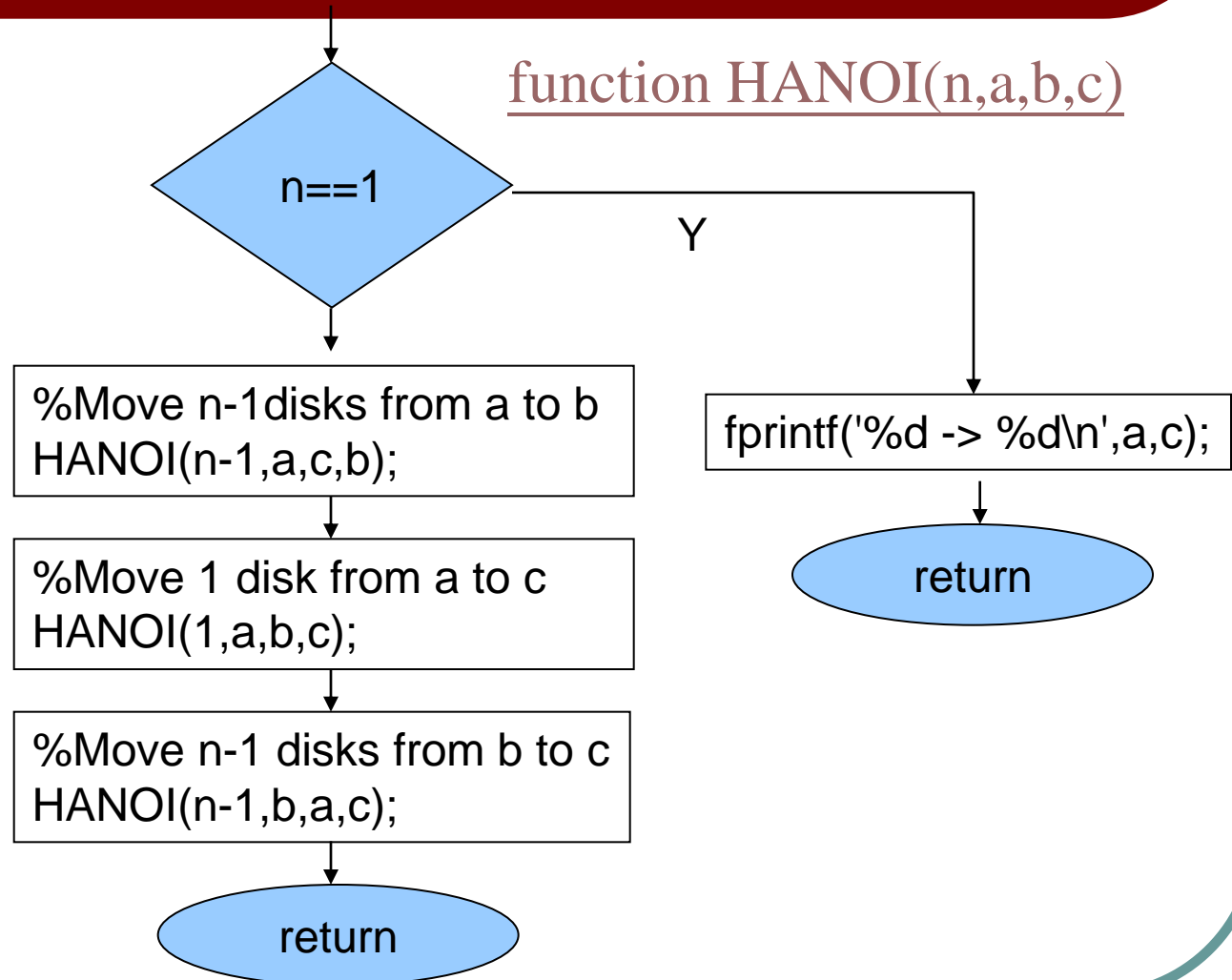


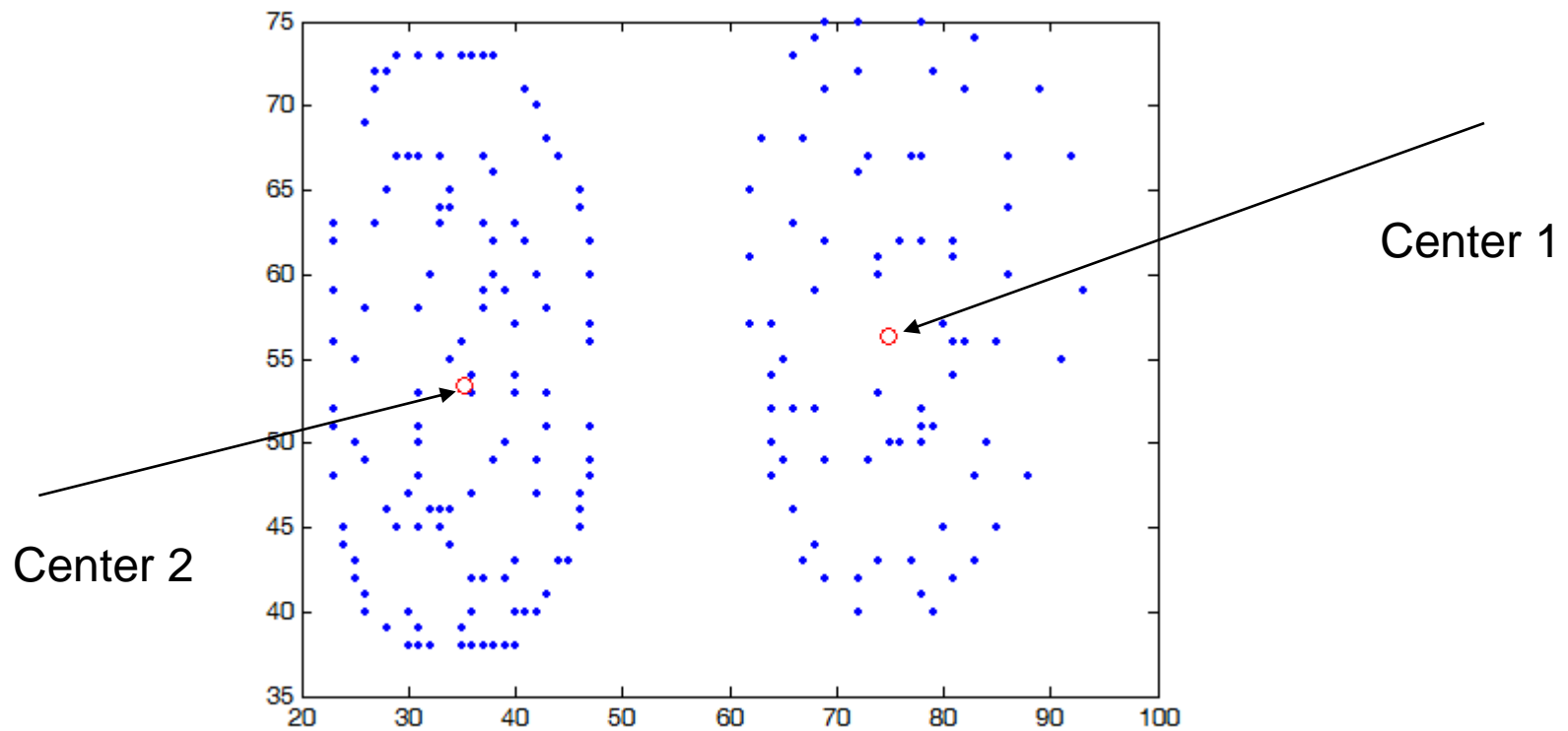
# Lecture 8 K-means for clustering

- Cross distances
- Exclusive memberships
- Clustering
  - An iterative approach

# Flow chart: move n disks from tower a to c



# Two clusters



# Cross distances

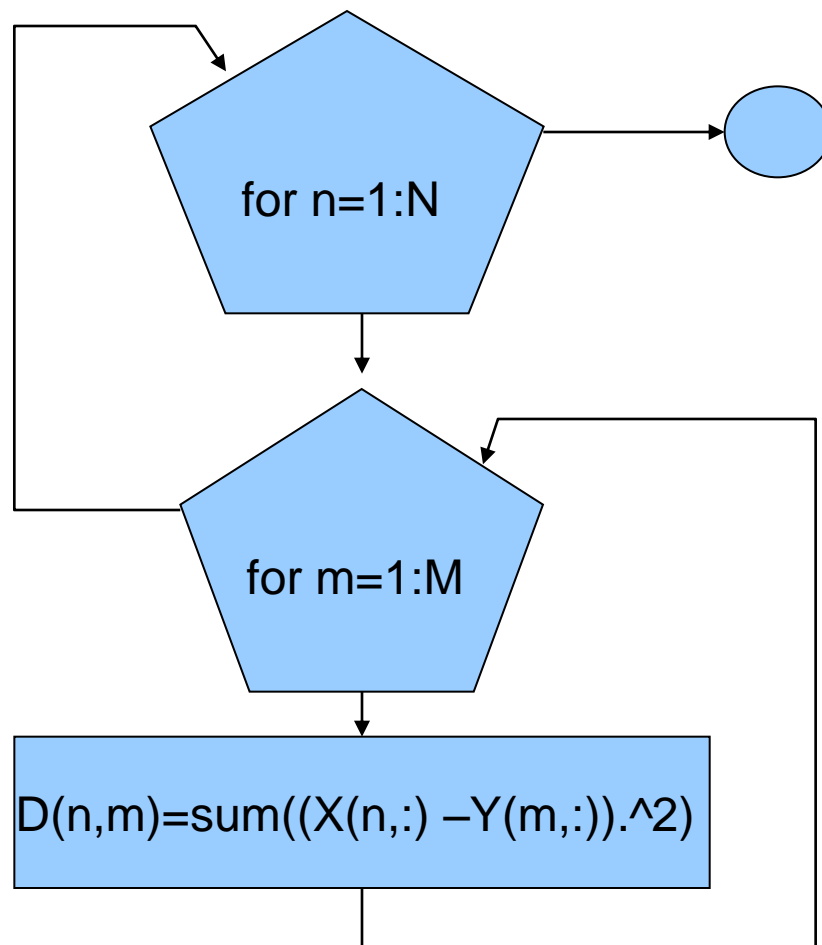
- How to find distances between centers and given points?

# Calculation of Cross distances

- Given  $N$  points  $X: N \times 2$
- $M$  centers  $Y: M \times 2$
- $D: N \times M$
- $D(i,j)$  denotes the distance between  $X(i,:)$  and  $Y(j,:)$
- Given  $X$  and  $Y$ , find  $D$

- It needs to calculate cross distances between  $N$  points and  $M$  centers to determine memberships of  $N$  points

# Nested loops for cross distances



# Matlab codes for nested codes

- ```
for i=1:N
    for j=1:M
        dd=X(i,:)-Y(j,:);
        D(i,j)=sqrt(sum(dd.^2));
    end
end
```



- Straightforward implementation
- Nested looping
  - A loop within a loop
  - MN calculations of the distance between a point and a center
- Time consuming for large M,N and d

# Vector codes

- How to calculate cross distances without using for-looping or while-looping ?
- Vector codes are loop-free
- Vector codes for cross distances can significantly improve efficiency against nested looping in computation

$$\begin{aligned} D_{ij} &= (\mathbf{x}_i - \mathbf{y}_j)(\mathbf{x}_i^T - \mathbf{y}_j^T) \\ &= \mathbf{x}_i \mathbf{x}_i^T - 2\mathbf{x}_i \mathbf{y}_j^T + \mathbf{y}_j \mathbf{y}_j^T \\ &= A_{ij} - 2B_{ij} + C_{ij} \end{aligned}$$

D : cross distances between N points and M centers

Matrix D is decomposed to matrices A, B and C

A : elements in a row are identical

B : multiplication of matrix X and transpose of matrix Y

C : elements in a column are identical

# Performance comparison

demo\_distance2.m

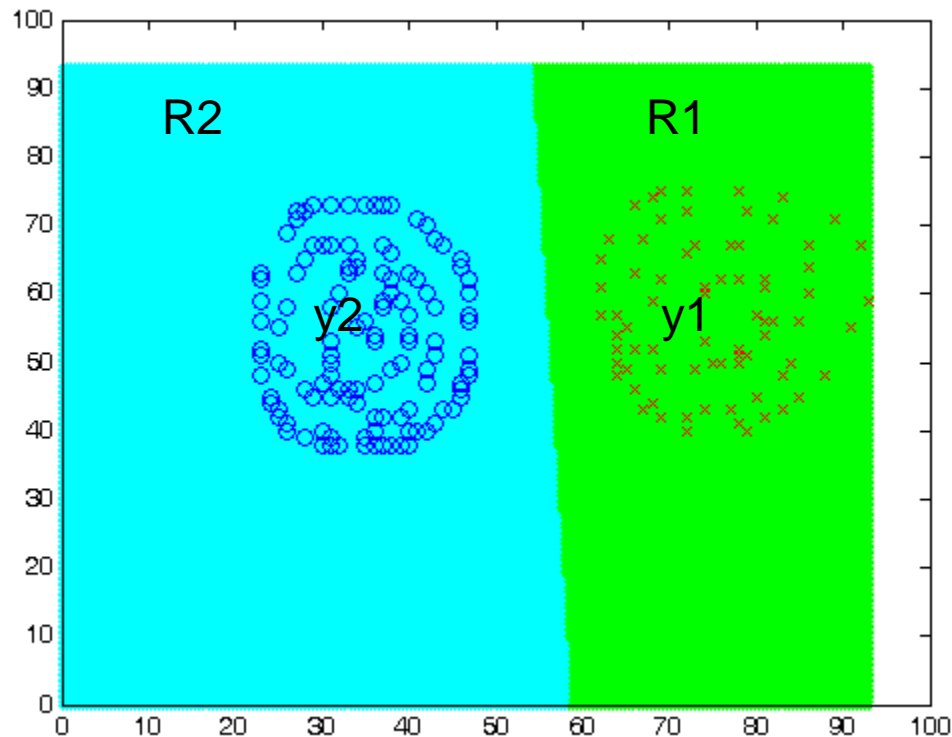
```
7 - for i=1:N
8 -     for j=1:M
9 -         dd=X(i,:)-Y(j,:);
10 -         D(i,j)=sqrt(sum(dd.^2));
11 -     end
12 - end
13 - ss2=cputime;
14 - A=sum(X.^2,2)*ones(1,M);
15 - C=ones(N,1)*sum(Y.^2,2)';
16 - B=X*Y';
17 - DD=sqrt(A-2*B+C);
18 - sum(sum(abs(DD-D)))
```

$$\begin{aligned} D_{ij} &= (\mathbf{x}_i - \mathbf{y}_j)(\mathbf{x}_i^T - \mathbf{y}_j^T) \\ &= \mathbf{x}_i \mathbf{x}_i^T - 2\mathbf{x}_i \mathbf{y}_j^T + \mathbf{y}_j \mathbf{y}_j^T \\ &= A_{ij} - 2B_{ij} + C_{ij} \end{aligned}$$

```
M=size(Y,1);N=size(X,1);  
A=sum(X.^2,2)*ones(1,M);  
C=ones(N,1)*sum(Y.^2,2)';  
B=X*Y';  
D=sqrt(A-2*B+C);
```

# Partition to two regions

Each point has its exclusive membership to non-overlapping regions partitioned by two centers



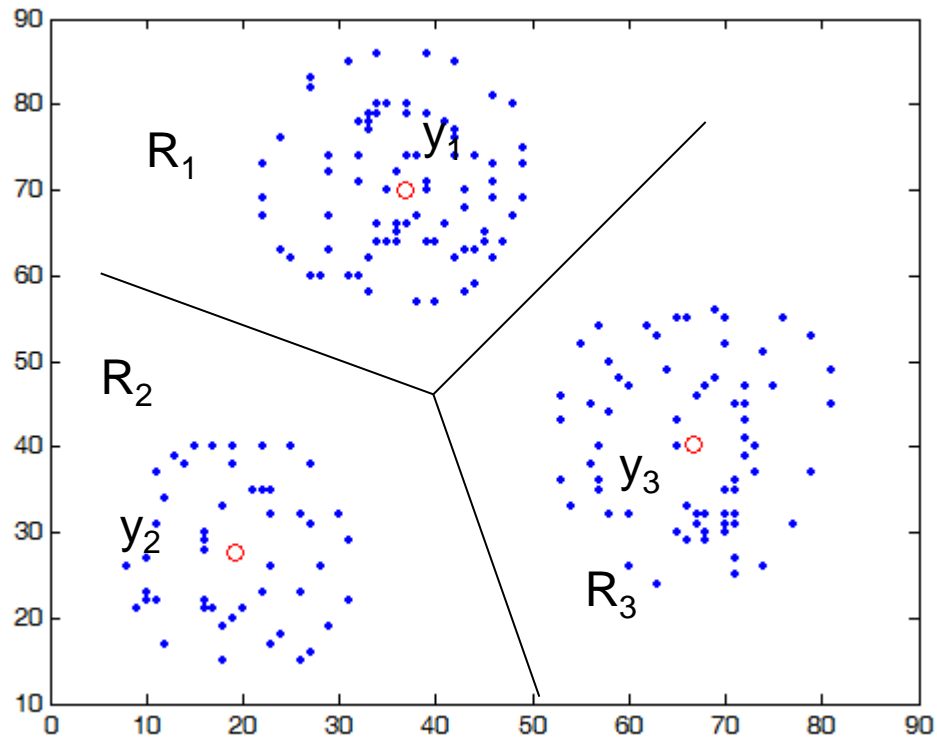
- Step A: calculate cross distances
- Step B: determine memberships
- Step C: determine K means

# Exclusive membership

- $y_1$  and  $y_2$  denote two centers
- $R_1$  and  $R_2$  denote two regions partitioned by  $y_1$  and  $y_2$
- A point belongs to  $R_1$  if it is closer to  $y_1$
- A point belongs to  $R_2$  if it is closer to  $y_2$

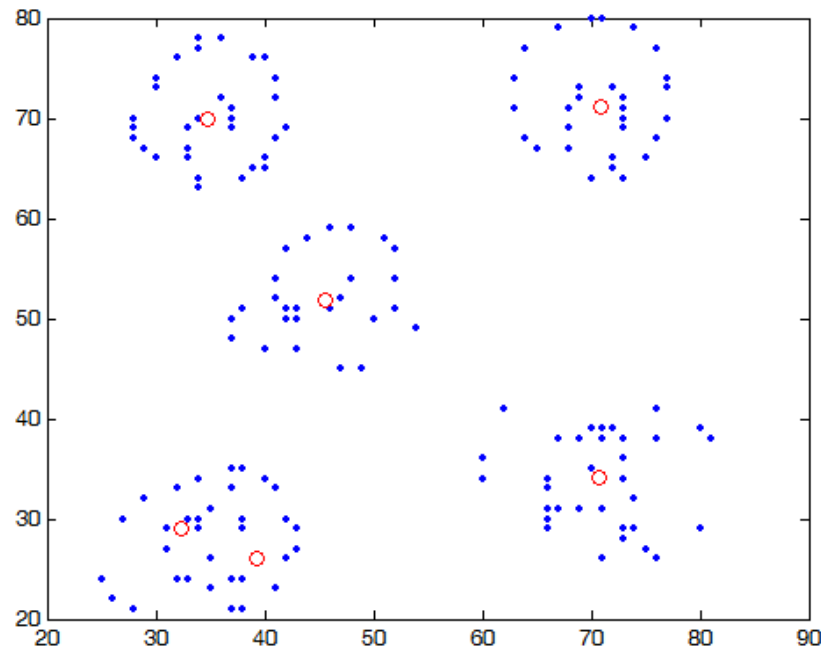


# Three clusters



# K clusters

- Locating K centers
- Significant geometric features of points in  $\mathbb{R}^d$



# Exclusive membership

- $y_1, y_2, \dots, y_K$  denote  $K$  distinct centers
- $R_1, R_2, \dots, R_K$  denote  $K$  regions partitioned by  $y_1, y_2, \dots, y_K$
- A point belongs to  $R_i$  if it is closest to  $y_i$  among  $K$  centers

- Let  $D$  denote cross distances between  $N$  given points and  $M$  centers
- Find points nearest to the  $j$ th center

Point  $x(i, :)$  belongs the  $j$ th cluster  
if

$$D(i, j) = \min_k D(i, k)$$

# Exclusive memberships (step B)

- Given  $D$ , exclusive memberships  $v$  of  $N$  points can be determined by

$$[xx \ v] = \min(D');$$

# Updating K-means (step C)

- Determine who belong the  $j$ th cluster

```
ind=find(v == j);
```

- Determine their mean

```
Y(j,:) = mean(X(ind,:))
```

# Clustering

- Where are  $K$  centers ?



# K-means

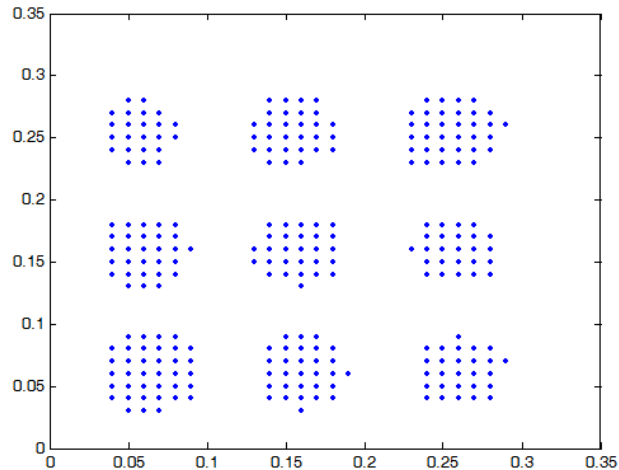
- A popular heuristic approach for clustering analysis
- An iterative approach
  - Step A : cross distance  $D$
  - Step B : exclusive memberships  $v$
  - Step C : updating centers  $Y$

# Data and MATLAB codes

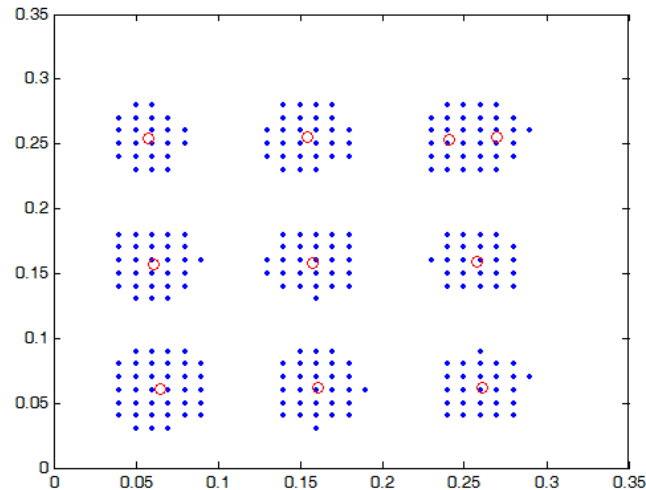
[data\\_9.zip](#)  
[demo\\_kmeans.m](#)

```
load data_9.mat
plot(X(:,1),X(:,2),'.');
[clidx, Y] = kmeans(X,10);
hold on;
plot(Y(:,1),Y(:,2),'ro');
```

# Data Clustering



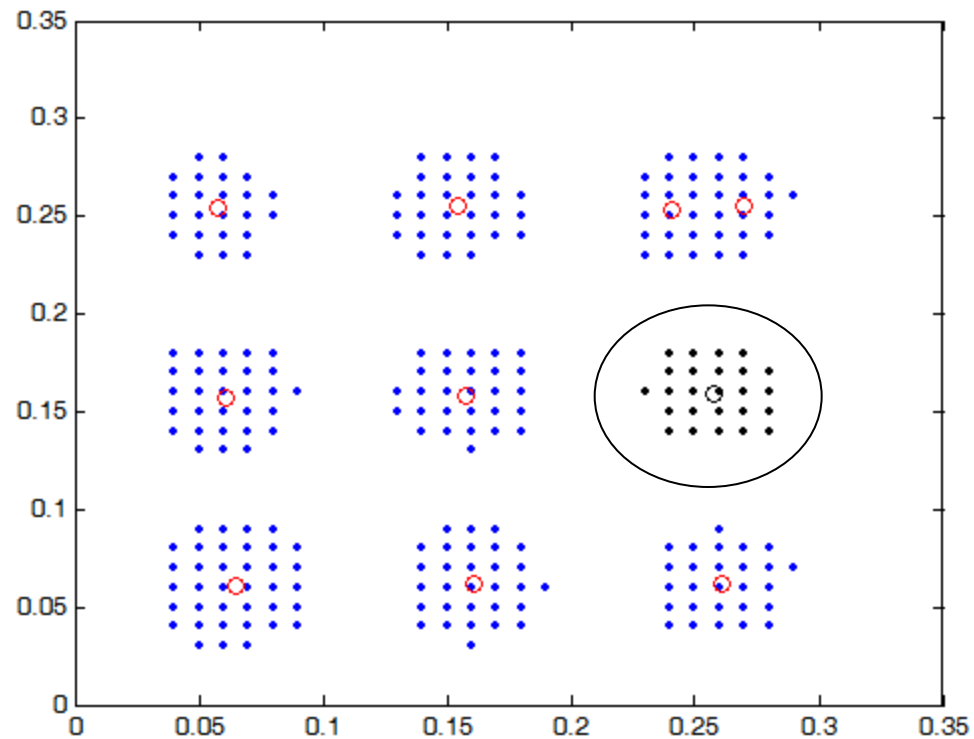
```
[cidx, Y] = kmeans(X,10);
```



- Partition given data to K clusters
- The K-means algorithm aims to find means (centers) of K clusters

# Memberships

Black points belong to the cluster centered at black circle



# A math tool for data clustering

ClusteringTest.rar

# Data Clustering

MATH PROGRAMMING  
AM NDHU

New

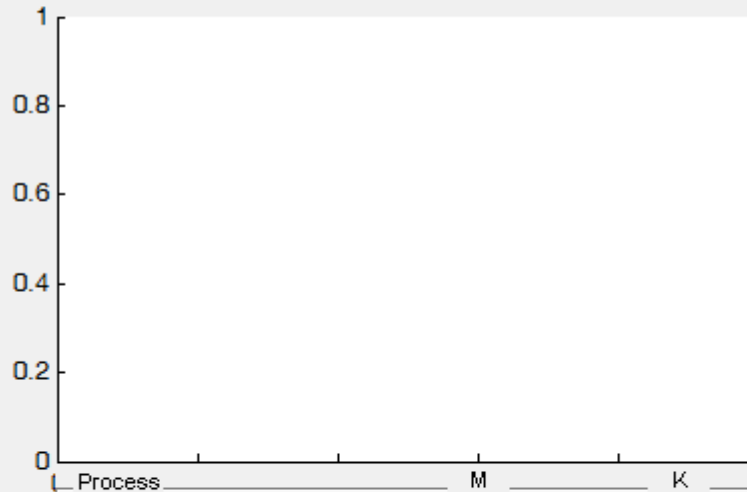
PenData

OK

Filing

LOAD

SAVE



KMEANS

Linear Separation

err rate

Steps for data clustering:

New  
PenData  
Enter M  
KMEANS

# Data Clustering

MATH PROGRAMMING  
AM NDHU

New

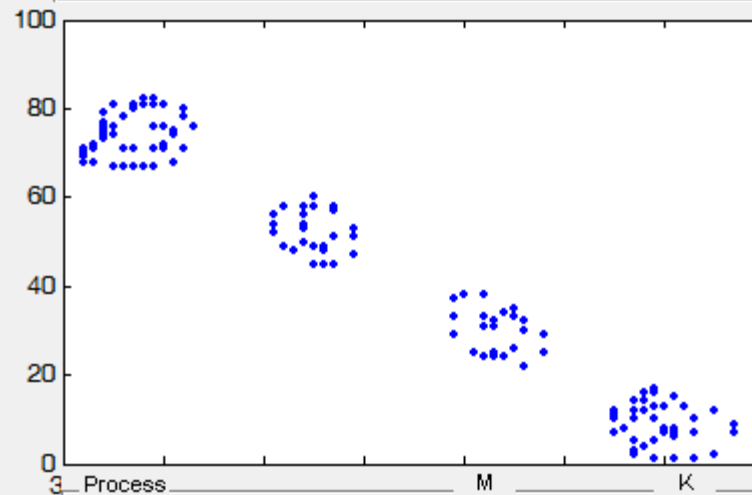
PenData

OK

Filing

LOAD

SAVE

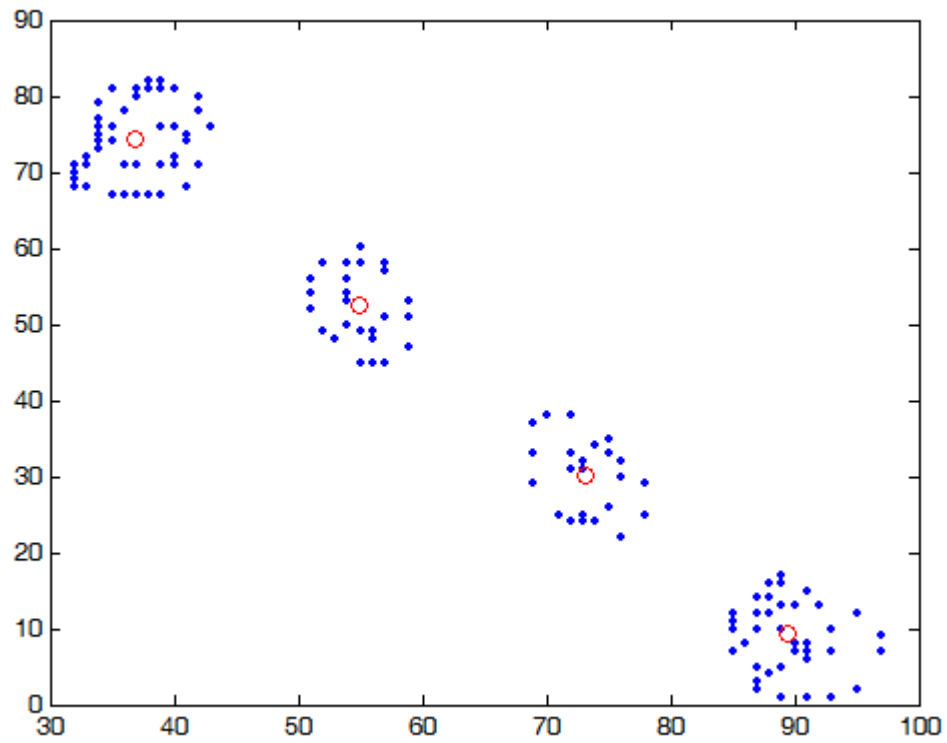


KMEANS

4

Linear Separation

err rate



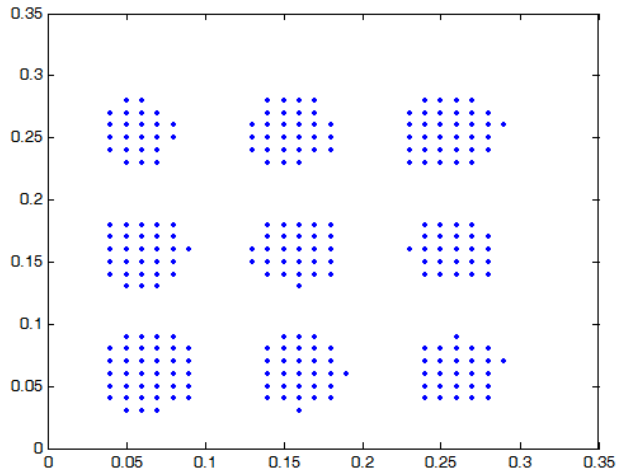


# Main steps of K-means clustering

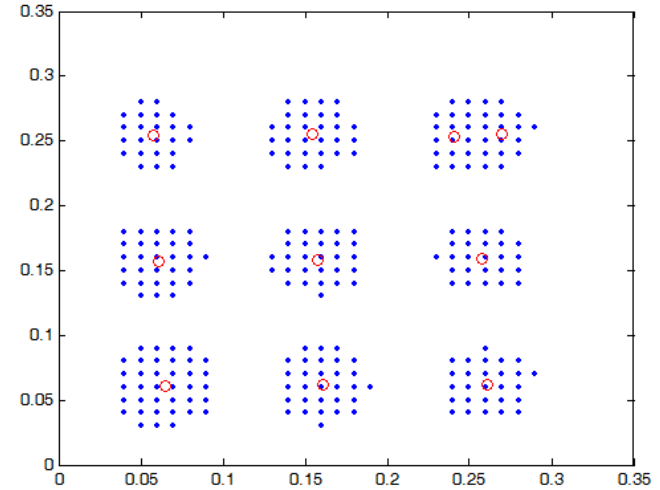
- Iterative execution of steps A, B and C until K means converge
  - A : cross distances  $D$
  - B : exclusive memberships  $v$
  - C : update K means  $Y$

# Step A: Calculation of Cross distances

- $X: N \times 2$
- $Y: M \times 2$
- $D: N \times M$
- $D(i,j)$  denotes the distance between  $X(i,:)$  and  $Y(j,:)$
- Given  $X$  and  $Y$ , find  $D$



$[cidx, Y] = kmeans(X, 10);$

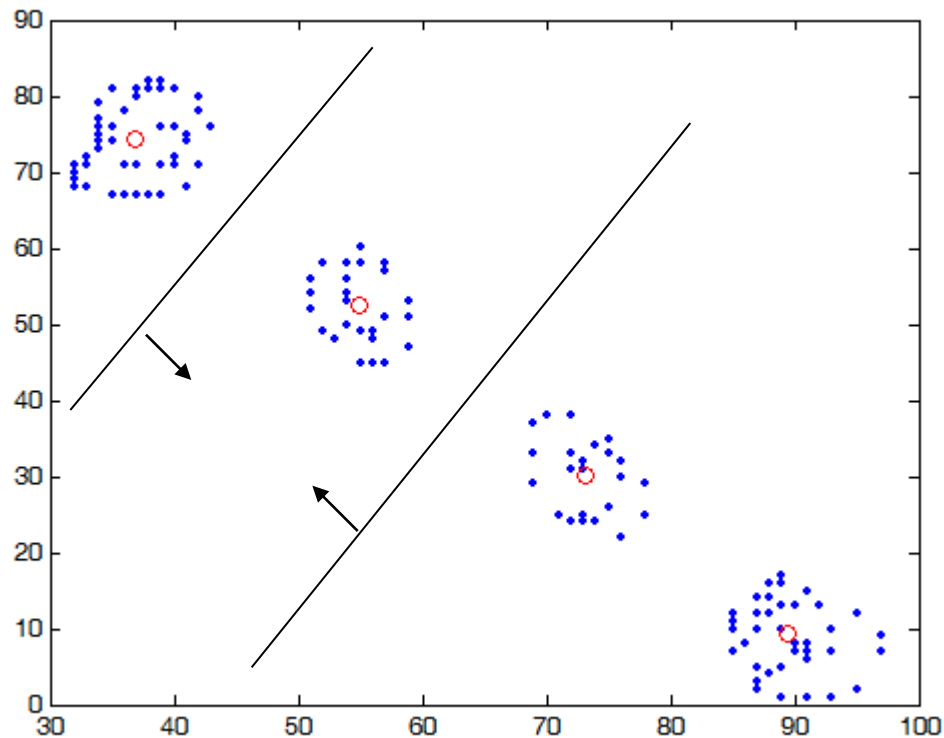


Cross  
Distances ←

```
M=size(Y,1);N=size(X,1);
A=sum(X.^2,2)*ones(1,M);
C=ones(N,1)*sum(Y.^2,2)';
B=X*Y';
D=sqrt(A-2*B+C);
```

# Membership

- $X(:,i)$  belongs a cluster



# Exclusive memberships (step B)

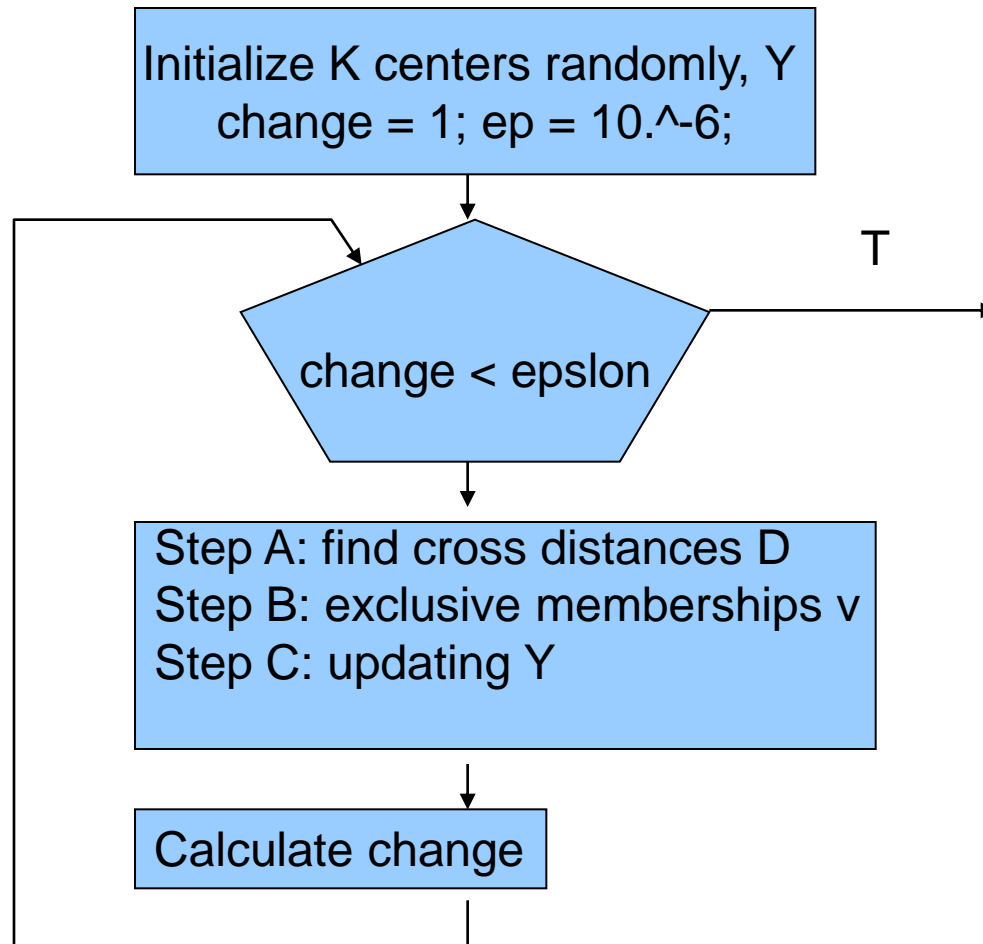
- Given  $D$ , exclusive memberships  $v$  of  $N$  points can be determined by

$$[xx \ v] = \min(D');$$

# Updating K-means (step C)

```
ind=find(v == j);
```

```
Y(j,:) = mean(X(ind,:))
```



# Initialization

- Calculate the mean of N points
- $\text{mean\_x} = \text{mean}(X)$
- $Y = \text{rand}(M, 2) * 0.1 - 0.05 + \text{mean\_x}$



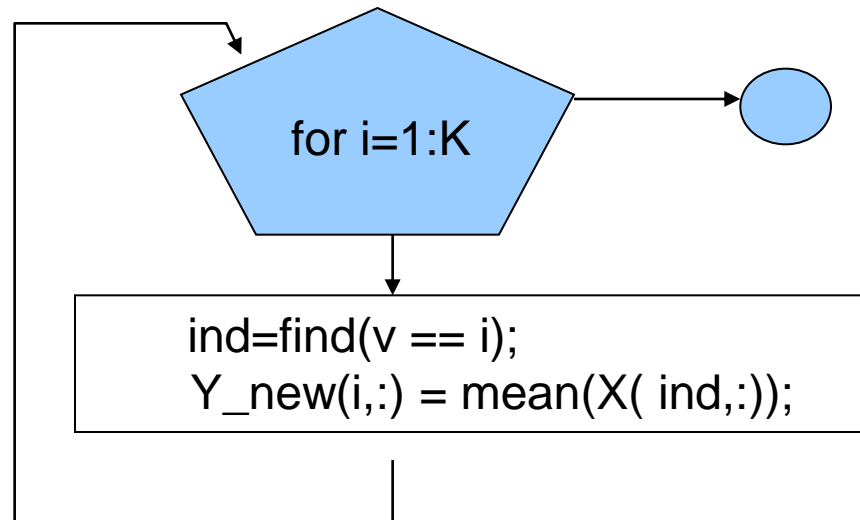
# Step A: Cross distances

$D = \text{cross\_distances}(X, Y)$

## Step B: Exclusive memberships $v$

- $[xx \ v]=\min(D')$ ;

# Step C: Update centers



# Partition N points to K clusters

```
D = cross_distances(X,Y)
[xx v]=min(D');

%
for i=1:K
    ind=find(v == i);
    Y_new(i,:) = mean(X( ind,:));
end
```

# Halting

- $\text{change} = \text{mean}(\text{mean}(\text{abs}(Y - Y_{\text{new}})))$
- Halting condition  
 $\text{chang} < \text{ep}$

# my\_kmeans

## my\_kmeans.m

```
load data_9.mat  
>> plot(X(:,1),X(:,2),'!');  
>> Y=my_kmeans(X,10);  
>> hold on  
>> plot(Y(:,1),Y(:,2),'or');
```

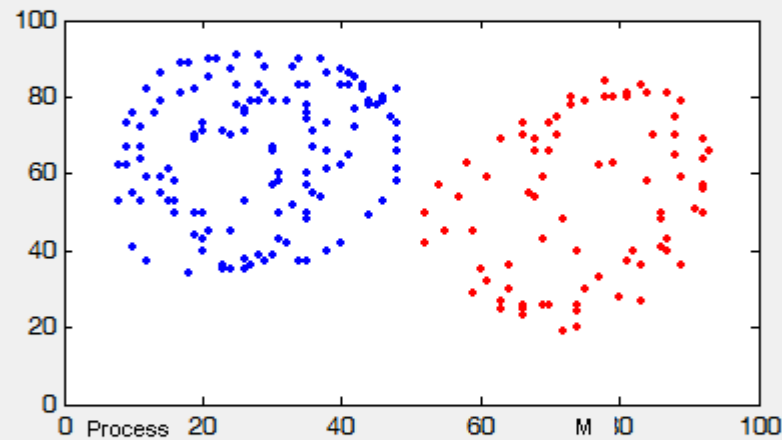
ClusteringTest2.fig  
ClusteringTest2.m

A version that is able to show stepwise execution of partitioning and updating of the kmeans algorithm

# Data Clustering

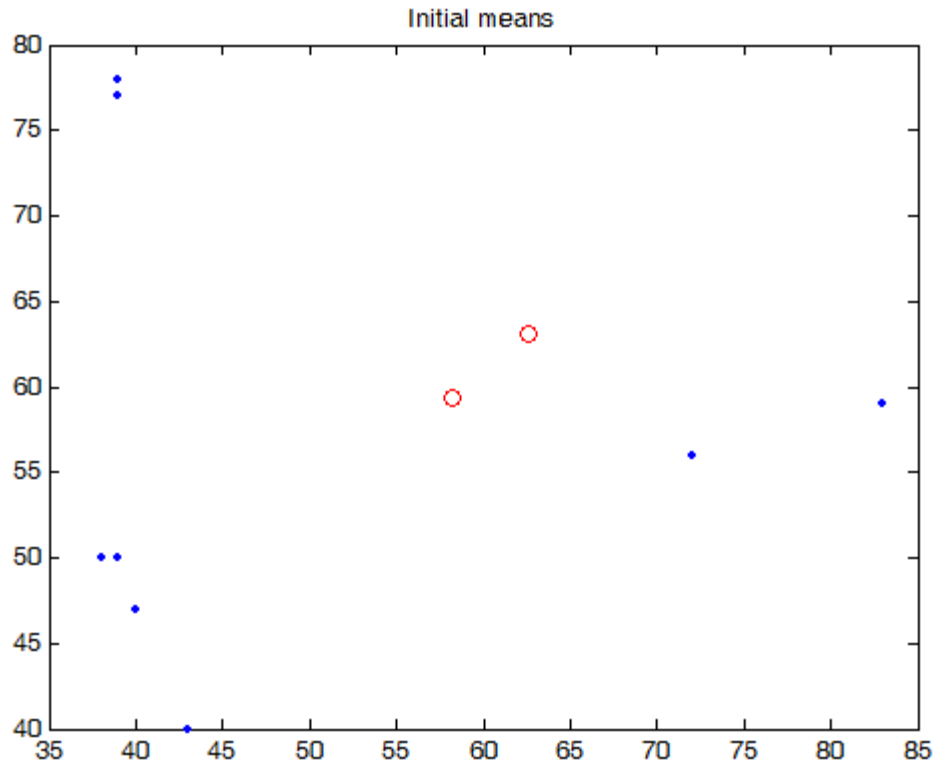
MATH PROGRAMMING  
AM NDHU

Filing





# Initialization



X =

|    |    |
|----|----|
| 39 | 78 |
| 39 | 77 |
| 43 | 40 |
| 40 | 47 |
| 38 | 50 |
| 72 | 56 |
| 83 | 59 |
| 39 | 50 |

Y =

|         |         |
|---------|---------|
| 49.0603 | 57.2121 |
| 49.1061 | 57.2084 |

# Cross distance

X =

|    |    |
|----|----|
| 39 | 78 |
| 39 | 77 |
| 43 | 40 |
| 40 | 47 |
| 38 | 50 |
| 72 | 56 |
| 83 | 59 |
| 39 | 50 |

Y =

|         |         |
|---------|---------|
| 49.0603 | 57.2121 |
| 49.1061 | 57.2084 |

D = cross\_distances(X,Y);

D =

|         |         |
|---------|---------|
| 23.0943 | 23.1176 |
| 22.1984 | 22.2226 |
| 18.2478 | 18.2596 |
| 13.6519 | 13.6797 |
| 13.2039 | 13.2404 |
| 22.9717 | 22.9257 |
| 33.9868 | 33.9412 |
| 12.3783 | 12.4135 |

[xx v]=min(D');

v =

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 2 | 2 | 1 |
|---|---|---|---|---|---|---|---|

# Partition & Update

ind =

1 1 1 1 1 2 2 1



```
for i=1:M
    ind=find(v == i);
    Y_new(i,:) =mean(X(ind,:));
end
```

Initialize K centers randomly, Y  
change = 1; ep = 10.<sup>-6</sup>;

function Y=my\_kmeans(X,M)

change < epsilon

T

Step A: find cross distances D  
Step B: exclusive memberships v  
Step C: updating Y

```
D = cross_distances(X,Y);  
[xx v]=min(D');  
for i=1:M  
    ind=find(v == i);  
    Y_new(i,:) =mean(X(ind,:));  
end  
change = mean(mean(abs(Y-Y_new)));  
Y=Y_new;
```

Calculate change

# k-means clustering - Wikipedia

# Description

Given a set of observations  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ , where each observation is a  $d$ -dimensional real vector, the  $k$ -means clustering aims to partition the  $n$  observations into  $k$  sets ( $k < n$ )  $\mathbf{S} = \{S_1, S_2, \dots, S_k\}$  so as to minimize the within-cluster sum of squares (WCSS):

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

where  $\boldsymbol{\mu}_i$  is the mean of points in  $S_i$ .

# History

The term "*k*-means" was first used by James MacQueen in 1967,<sup>[1]</sup> though the idea goes back to Hugo Steinhaus in 1956.<sup>[2]</sup> The *standard algorithm* was first proposed by Stuart Lloyd in 1957 as a technique for *pulse-code modulation*, though it wasn't published until 1982.<sup>[3]</sup>

# Standard Algorithm

**Assignment step:** Assign each observation to the cluster with the closest mean (i.e. partition the observations according to the [Voronoi diagram](#) generated by the means).

$$S_i^{(t)} = \left\{ \mathbf{x}_j : \|\mathbf{x}_j - \mathbf{m}_i^{(t)}\| \leq \|\mathbf{x}_j - \mathbf{m}_{i^*}^{(t)}\| \text{ for all } i^* = 1, \dots, k \right\}$$

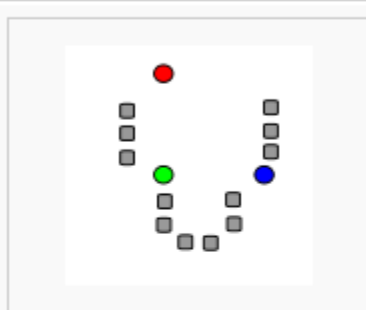
**Update step:** Calculate the new means to be the centroid of the observations in the cluster.

$$\mathbf{m}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x}_j \in S_i^{(t)}} \mathbf{x}_j$$

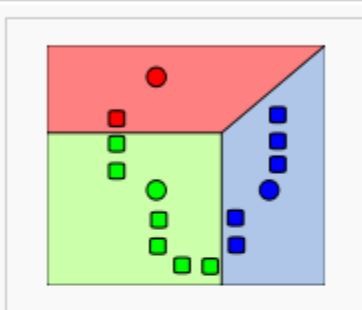


The algorithm is deemed to have converged when the assignments no longer change.

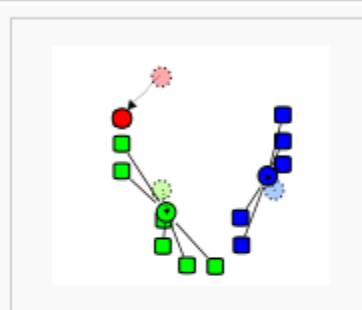
### Demonstration of the standard algorithm



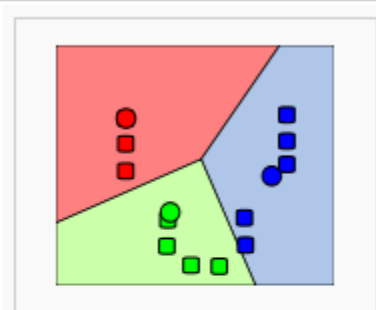
1)  $k$  initial "means" (in this case  $k=3$ ) are randomly selected from the data set (shown in color).



2)  $k$  clusters are created by associating every observation with the nearest mean. The partitions here represent the **Voronoi diagram** generated by the means.



3) The **centroid** of each of the  $k$  clusters becomes the new means.



4) Steps 2 and 3 are repeated until convergence has been reached.

# Advanced topics based on K-means

- Classification
- Function approximation
- Density estimation