# CS 484
# Data Mining
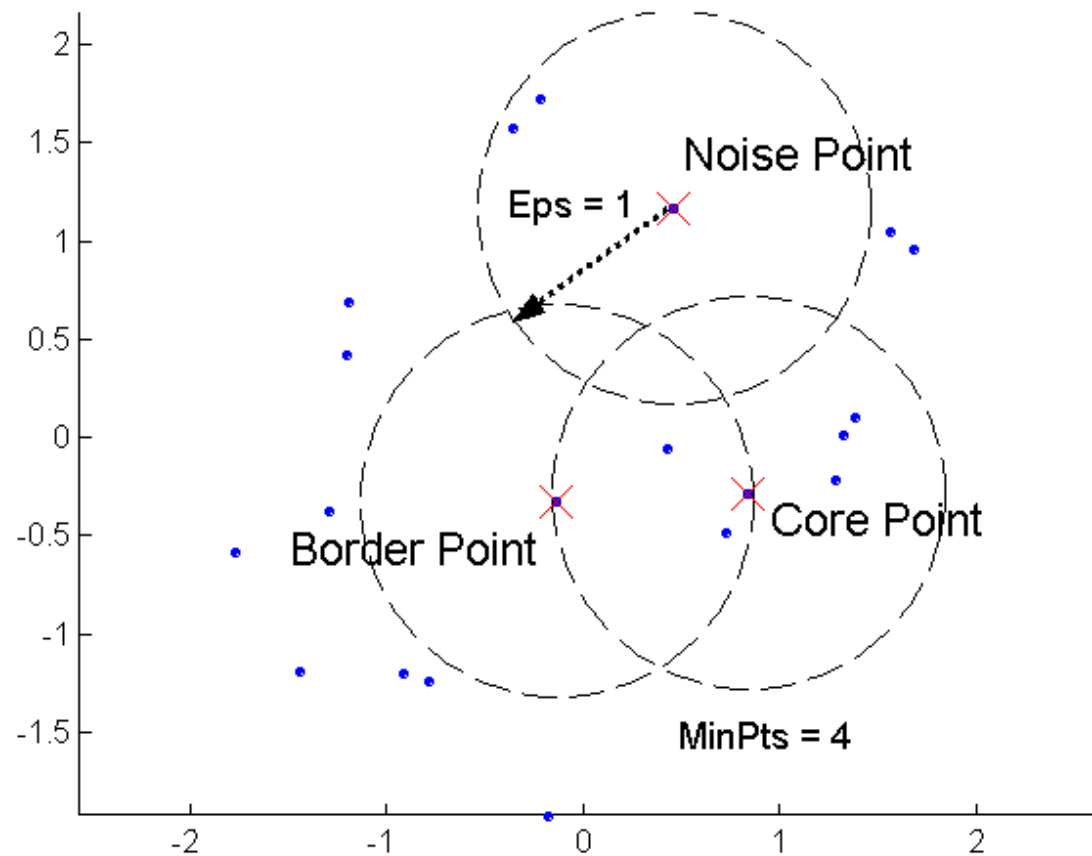
Clustering 5

# DBSCAN

- DBSCAN is a density-based algorithm.
  - Density = number of points within a specified radius (Eps)
  - A point is a core point if it has more than a specified number of points (MinPts) within Eps
    - These are points that are at the interior of a cluster
  - A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point
  - A noise point is any point that is not a core point or a border point.
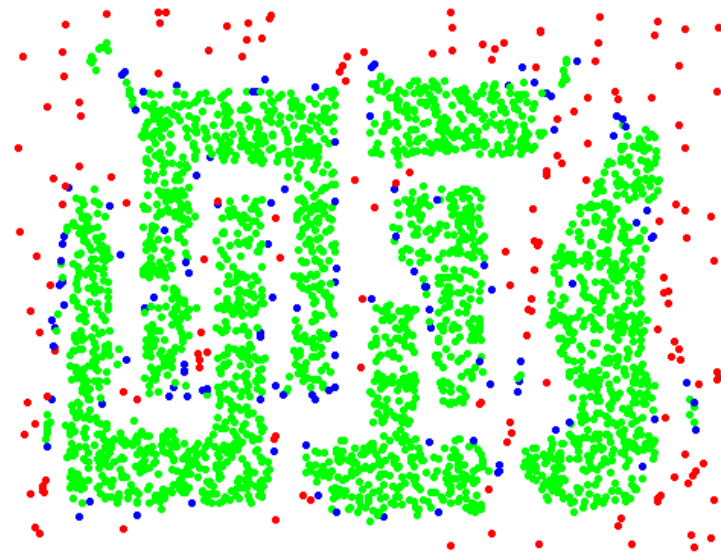
# DBSCAN: Core, Border, and Noise Points

# DBSCAN Algorithm

- Label all points as core, border or noise
- Eliminate noise points
- Put an edge between all core points that are within Eps of each other.
- Make each group of connected points into a separate cluster.
- Assign each border point to one of the clusters of its associated core points.

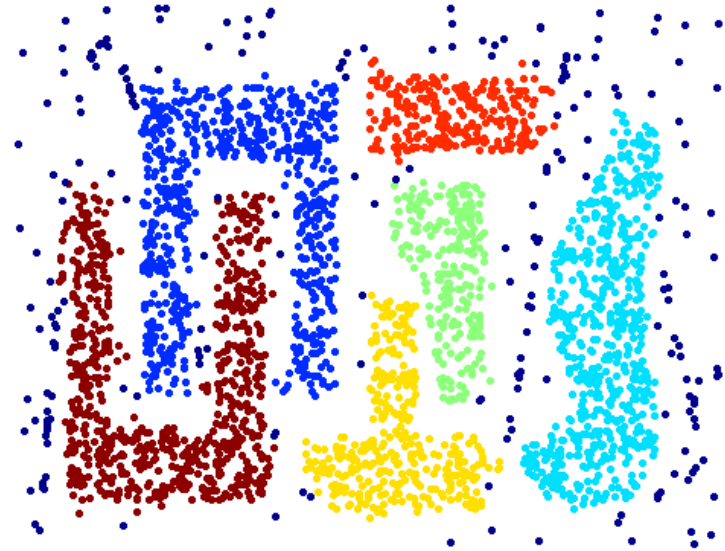# DBSCAN: Core, Border and Noise Points



**Original Points**

**Eps = 10, MinPts = 4**

**Point types: core, border and noise**

# When DBSCAN Works Well
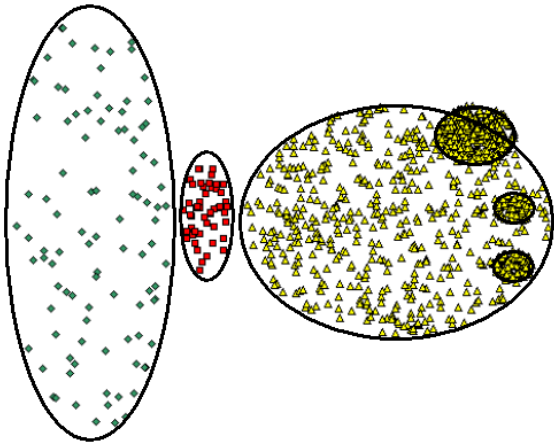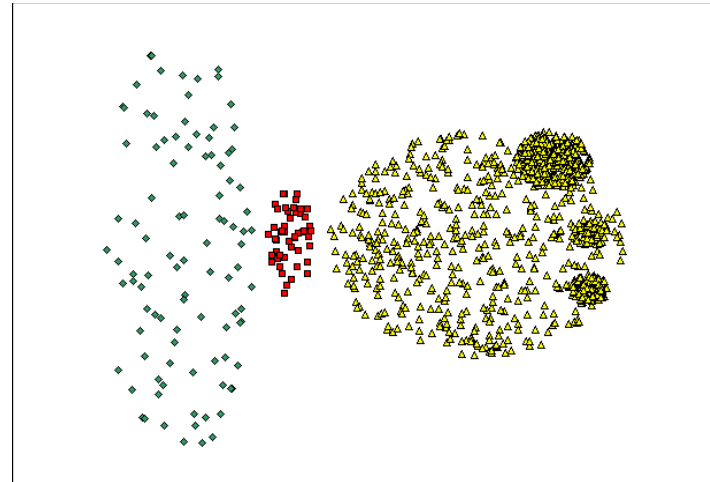


Original Points

Clusters

- **Resistant to Noise**

- **Can handle clusters of different shapes and sizes**
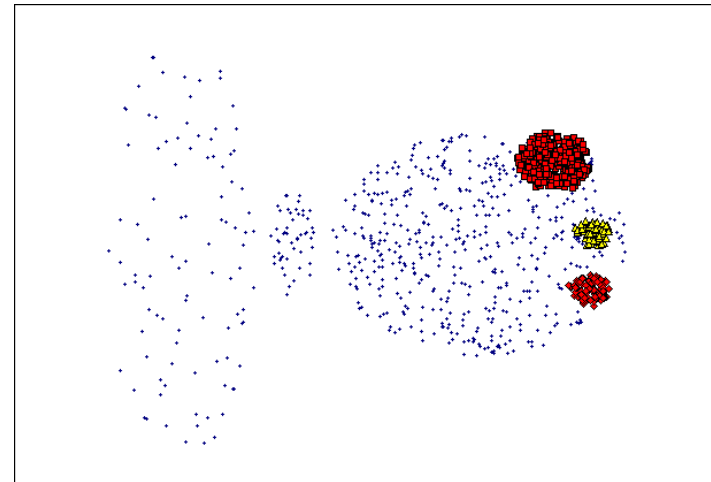
# When DBSCAN Does NOT Work Well



**Original Points**

(MinPts=4, Eps=9.75).

(MinPts=4, Eps=9.92)

- Varying densities
- High-dimensional data

# Shared Near Neighbor Approach

SNN graph: the weight of an edge is the number of shared neighbors between vertices given that the vertices are connected

# DBSCAN using SNN

1. **Find the SNN density of each Point.**
   Using a user specified parameters, *Eps*, find the number points that have an SNN similarity of *Eps* or greater to each point. This is the SNN density of the point

5. **Find the core points**
   Using a user specified parameter, *MinPts*, find the core points, i.e., all points that have an SNN density greater than *MinPts*

6. **Form clusters from the core points**
   If two core points are within a radius, *Eps*, of each other they are place in the same cluster
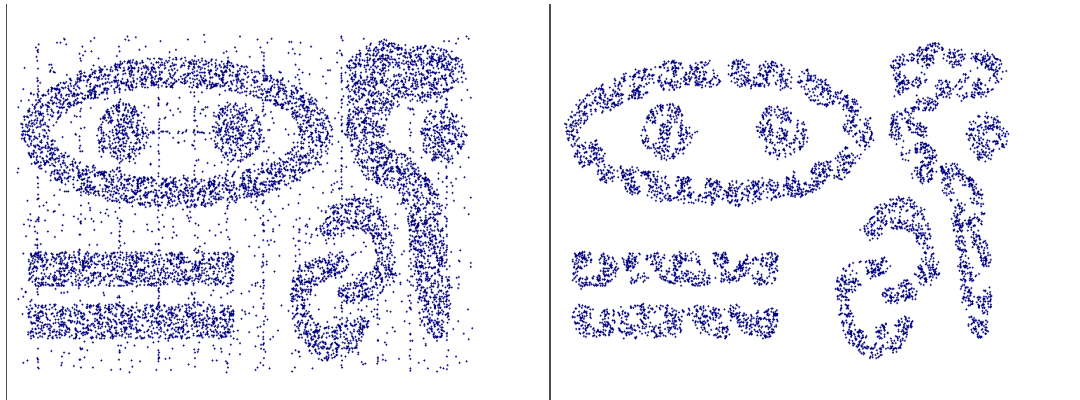
7. **Discard all noise points**
   All non-core points that are not within a radius of *Eps* of a core point are discarded

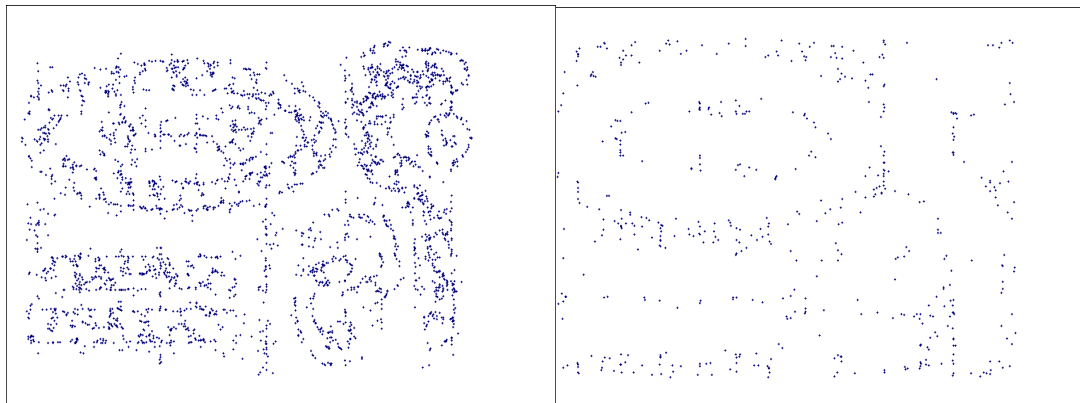8. **Assign all non-noise, non-core points to clusters**
   This can be done by assigning such points to the nearest core point
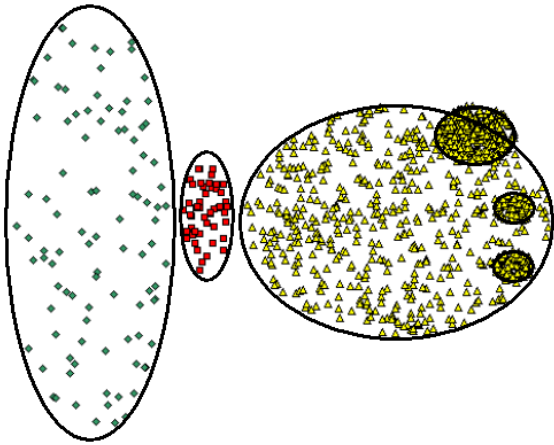
# SNN Density
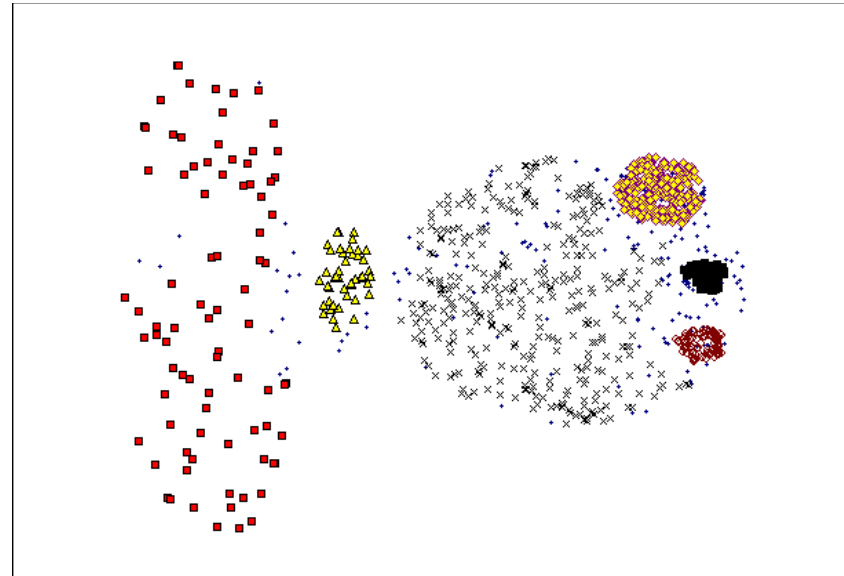


a) All Points

b) High SNN Density

c) Medium SNN Density

d) Low SNN Density

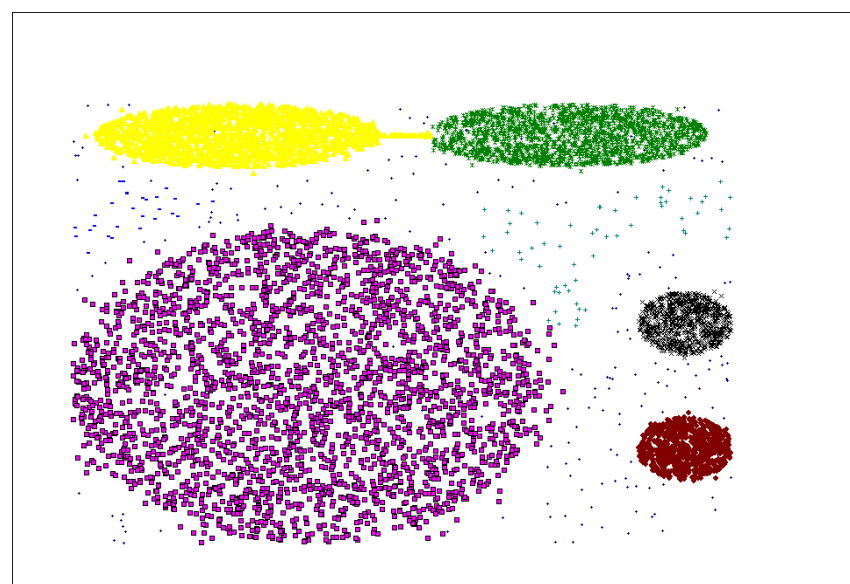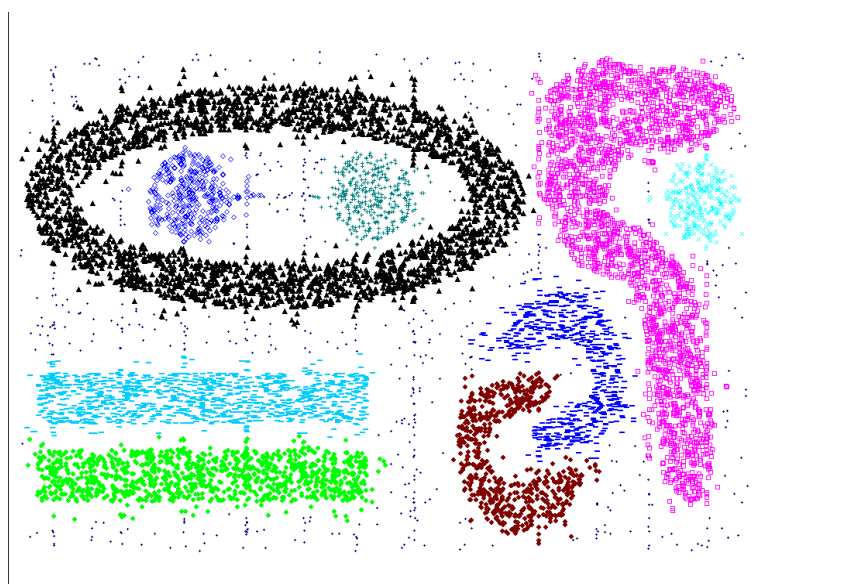# SNN Clustering Can Handle Differing Densities



Original Points



SNN Clustering

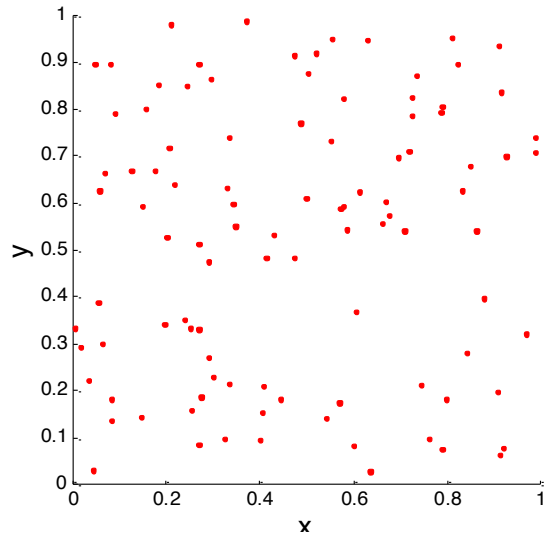# SNN Clustering Can Handle Other Difficult Situations
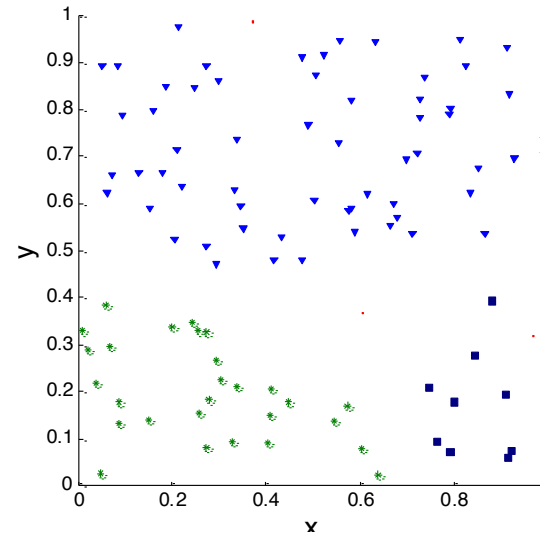
# Cluster Validity

- For supervised classification we have a variety of measures to evaluate how good our model is
  - Accuracy, precision, recall

- For cluster analysis, the analogous question is how to evaluate the "goodness" of the resulting clusters?

- But "clusters are in the eye of the beholder"!

- Then why do we want to evaluate them?
  - To avoid finding patterns in noise
  - To compare clustering algorithms
  - To compare two sets of clusters
  - To compare two clusters
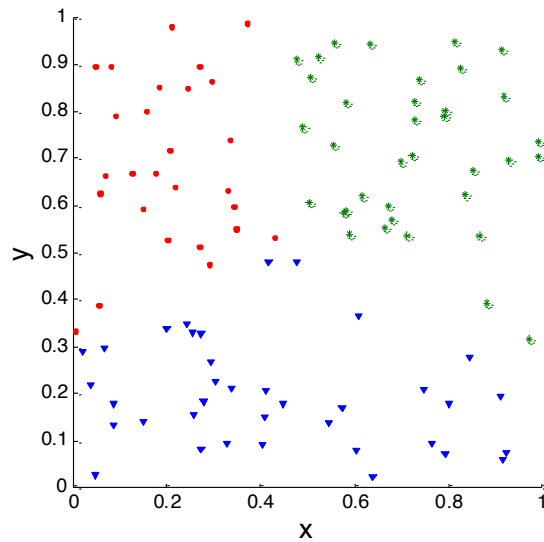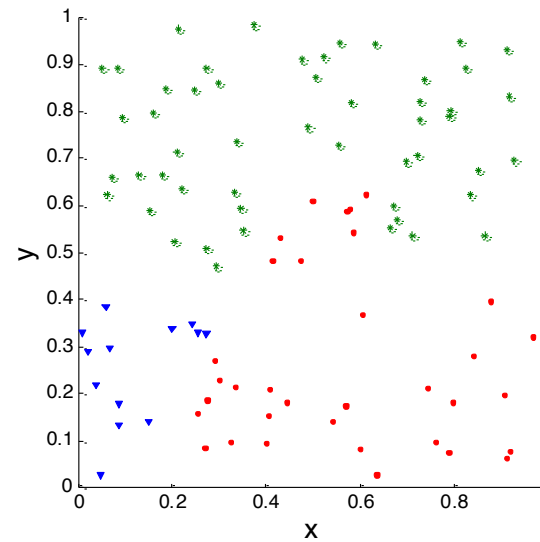
# Clusters found in Random Data



**Random Points**

**DBSCAN**

**K-means**

**Complete Link**

# Different Aspects of Cluster Validation

1.  Determining the <span style="color:red">clustering tendency</span> of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.

2.  Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.

3.  Evaluating how well the results of a cluster analysis fit the data *without* reference to external information.

    - Use only the data

4.  Comparing the results of two different sets of cluster analyses to determine which is better.

5.  Determining the 'correct' number of clusters.

    For 2, 3, and 4, we can further distinguish whether we want to evaluate the entire clustering or just individual clusters.

# Measures of Cluster Validity

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.
    - External Index: Used to measure the extent to which cluster labels match externally supplied class labels.
        - Entropy
    - Internal Index:  Used to measure the goodness of a clustering structure *without* respect to external information.
        - Sum of Squared Error (SSE)
    - Relative Index: Used to compare two different clusterings or clusters.
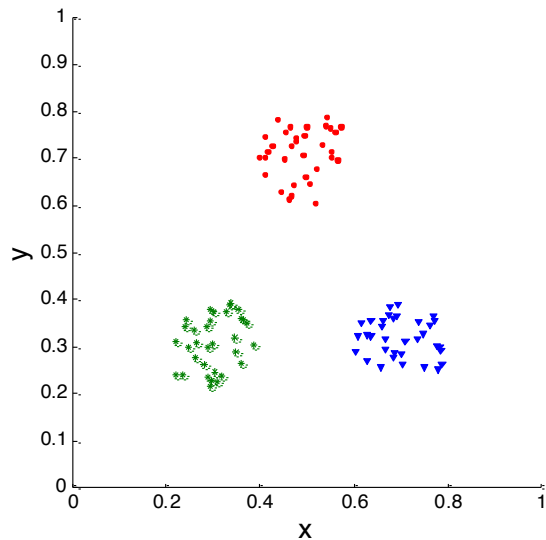        - Often an external or internal index is used for this function, e.g., SSE or entropy
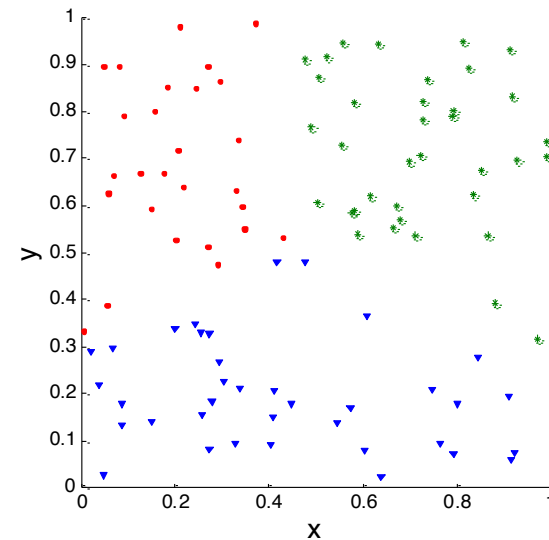
# Measuring Cluster Validity Via Correlation

- Two matrices
  - Proximity Matrix
  - "Incidence" Matrix
    - One row and one column for each data point
    - An entry is 1 if the associated pair of points belong to the same cluster
    - An entry is 0 if the associated pair of points belongs to different clusters

- Compute the correlation between the two matrices
  - Since the matrices are symmetric, only the correlation between $n(n-1) / 2$ entries needs to be calculated.

- High correlation indicates that points that belong to the same cluster are close to each other.

- Not a good measure for some density or contiguity based clusters.

# Measuring Cluster Validity Via Correlation

- Correlation of incidence and proximity matrices for the K-means clusterings of the following two data sets.
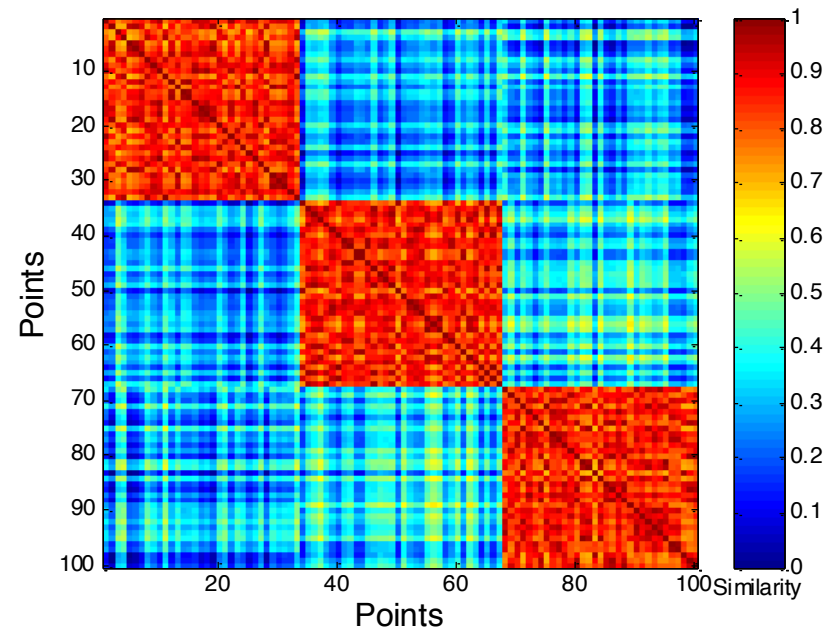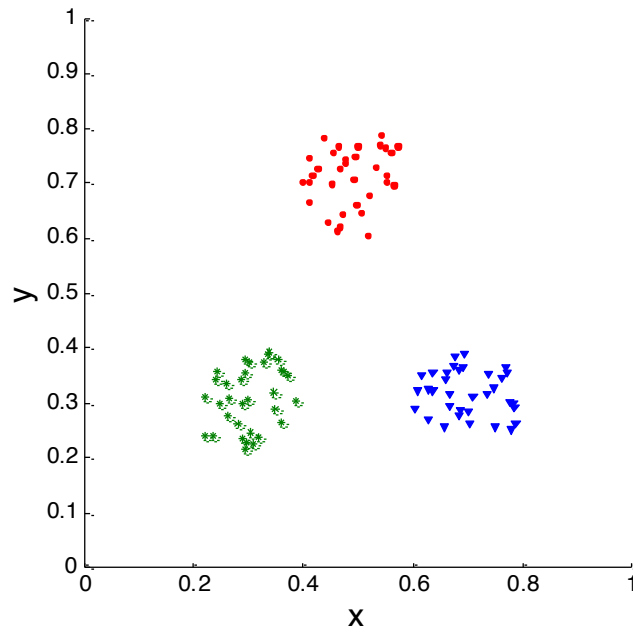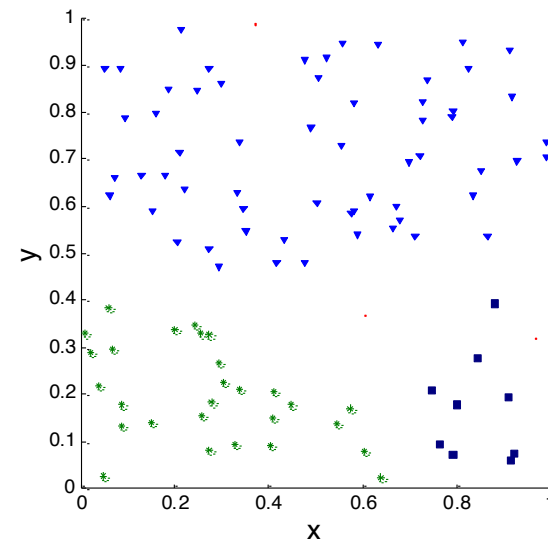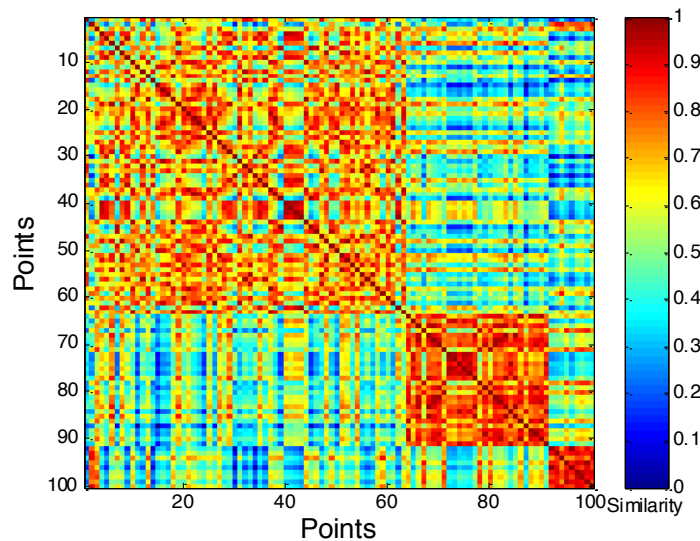


Corr = -0.9235



Corr = -0.5810

# Using Similarity Matrix for Cluster Validation

- Order the similarity matrix with respect to cluster labels and inspect visually.

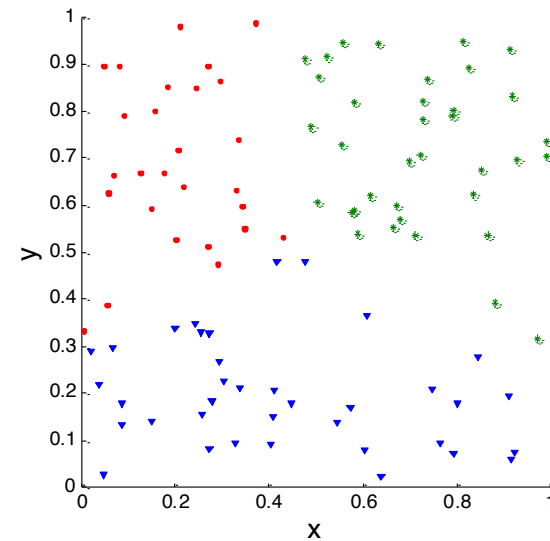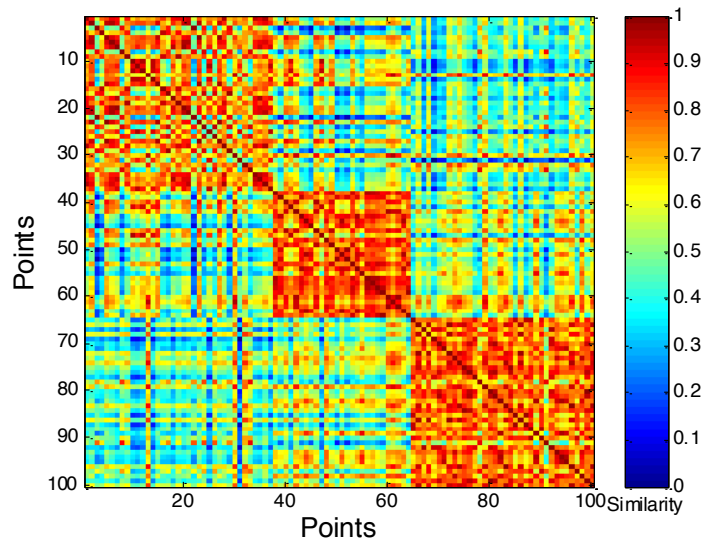# Using Similarity Matrix for Cluster Validation

- Clusters in random data are not so crisp



**DBSCAN**

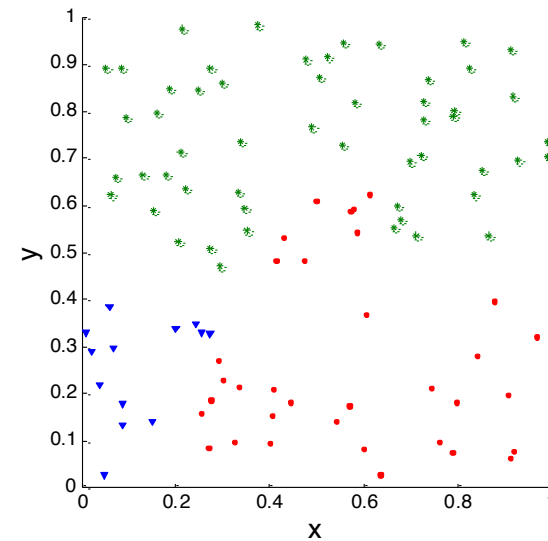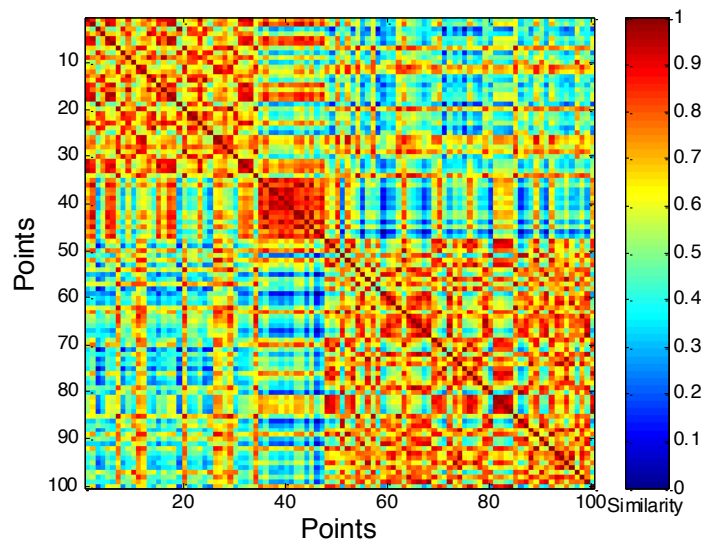# Using Similarity Matrix for Cluster Validation

- Clusters in random data are not so crisp
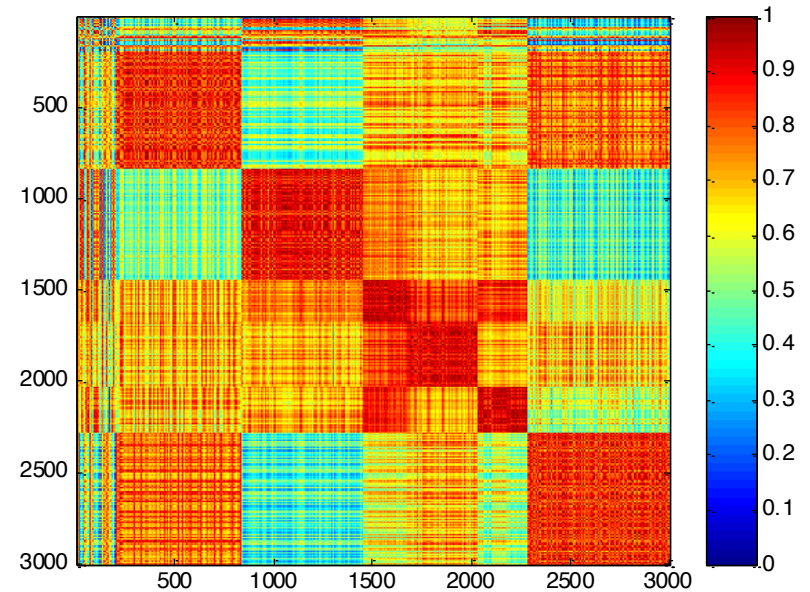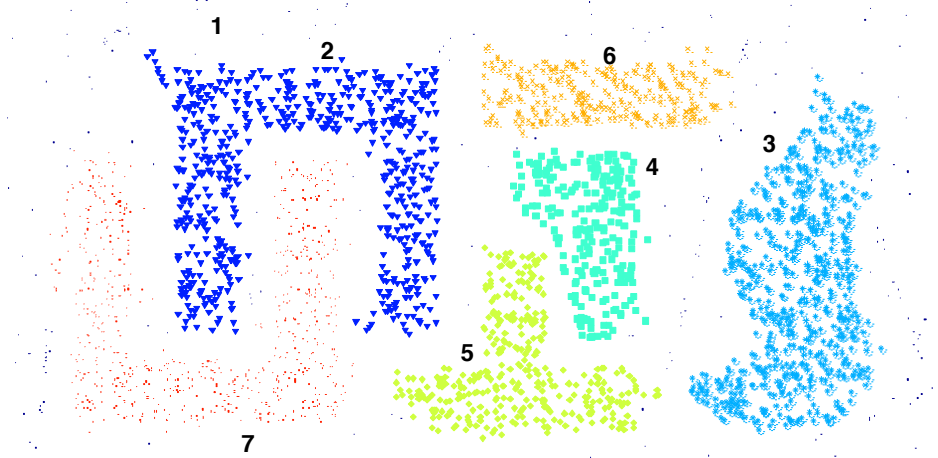


**K-means**

# Using Similarity Matrix for Cluster Validation

- Clusters in random data are not so crisp
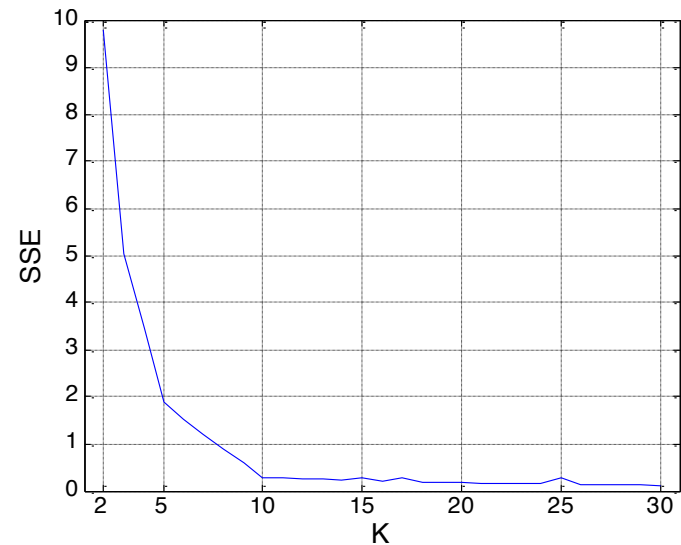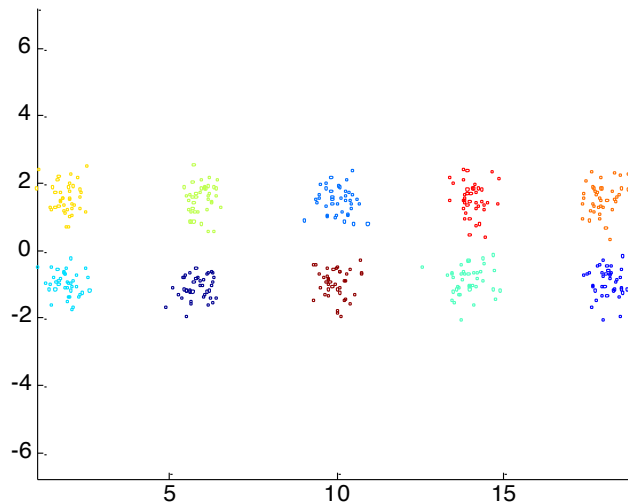


**Complete Link**

# Using Similarity Matrix for Cluster Validation



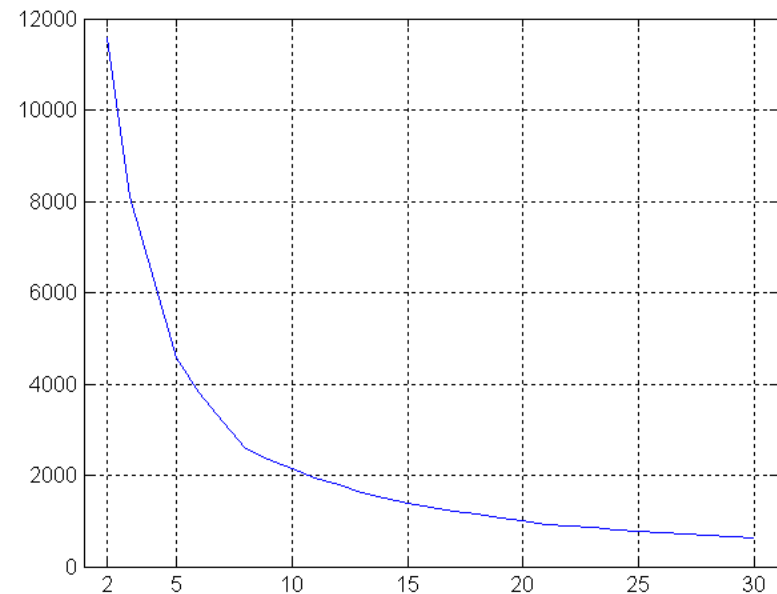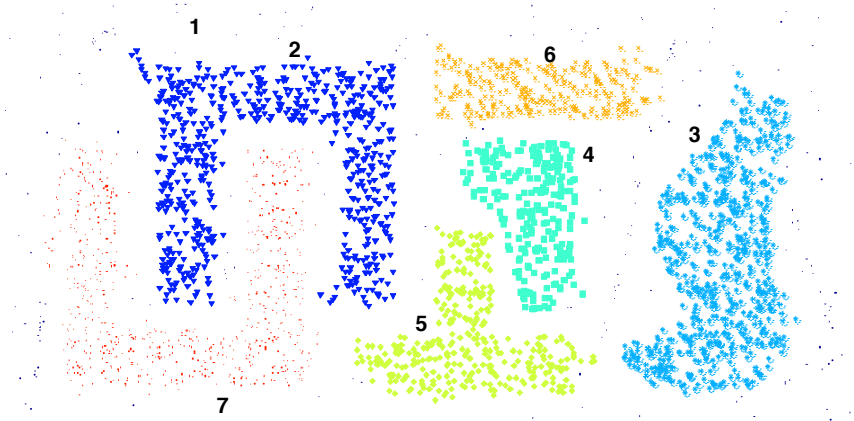**DBSCAN**

# Internal Measures: SSE

- Clusters in more complicated figures aren't well separated

- Internal Index: Used to measure the goodness of a clustering structure without respect to external information
  - SSE

- SSE is good for comparing two clusterings or two clusters (average SSE).

- Can also be used to estimate the number of clusters

# Internal Measures: SSE

- SSE curve for a more complicated data set



**SSE of clusters found using K-means**