

CLUSTERING EVALUATION

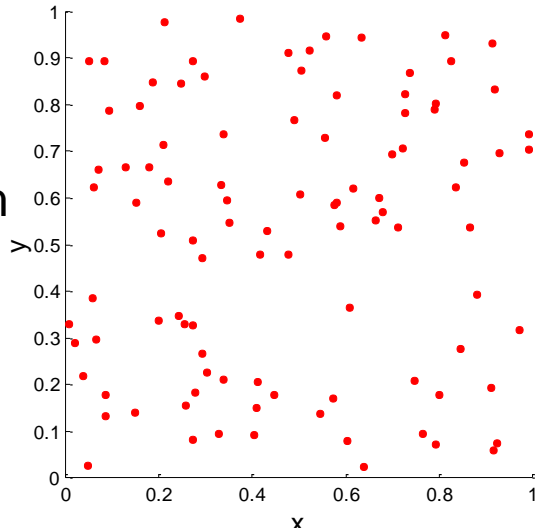
Lecture 05.03

Clustering Evaluation

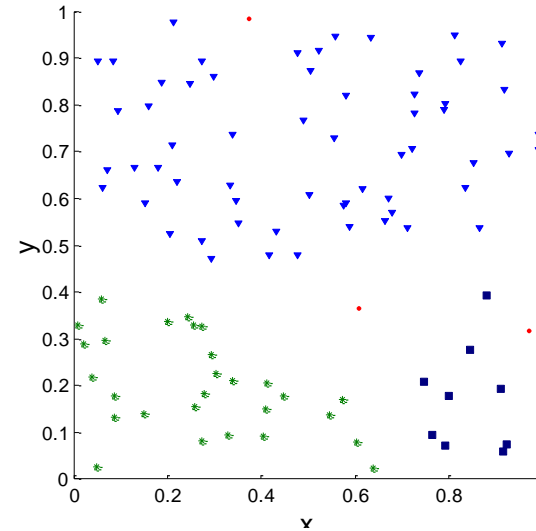
- How do we evaluate the “goodness” of the resulting clusters?
- But “clustering lies in the eye of the beholder”!
- Then why do we want to evaluate them?
 - To avoid finding patterns in noise
 - To compare clusterings, or clustering algorithms
 - To compare against a “ground truth”

Clusters found in Random Data

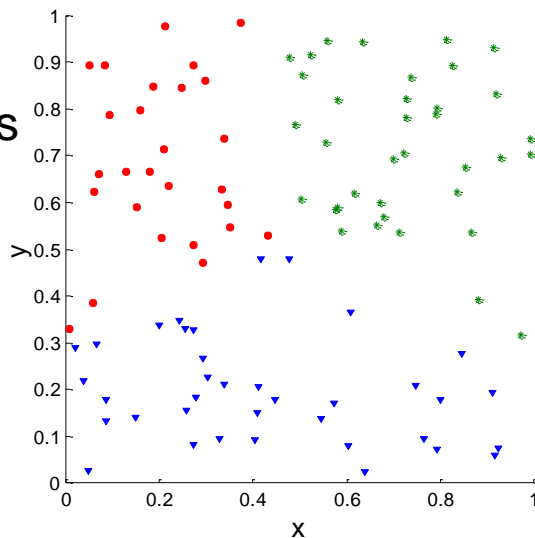
Random
Points



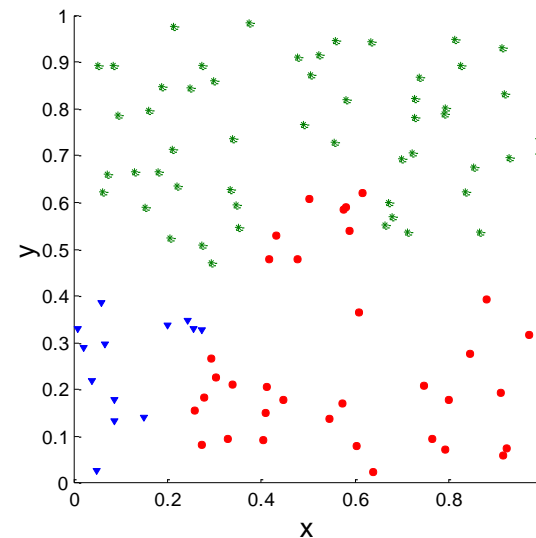
DBSCAN



K-means



Complete
Link



Different Approaches to Cluster Validation

1. Determining the **clustering tendency** 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

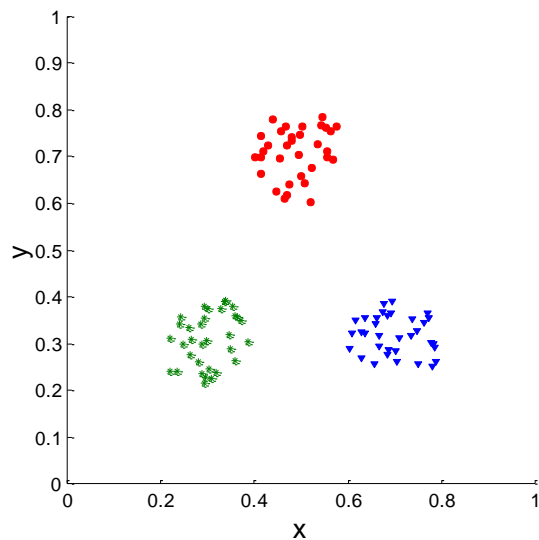
- Numerical measures are classified into the following three types:
 - **External Index:** Used to measure the extent to which cluster labels match **externally supplied class labels**.
 - E.g., entropy, precision, recall
 - **Internal Index:** Used to measure the goodness of a clustering structure **without** reference to external information.
 - E.g., 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
- Criterion is the **general strategy** and index is the **numerical measure** that implements the criterion.

Measuring Cluster Validity Via Correlation

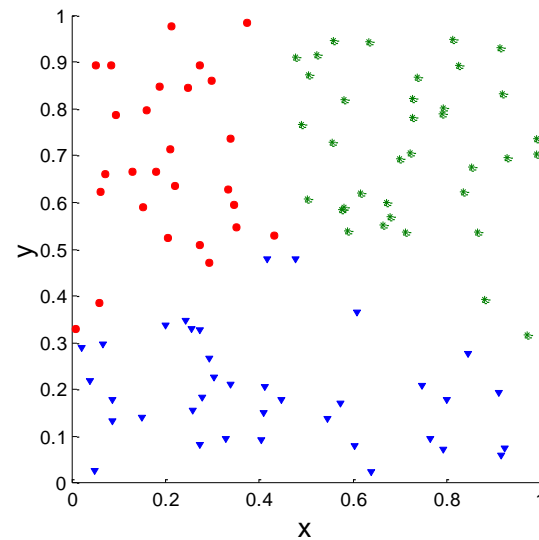
- Two matrices
 - **Similarity** or **Distance** Matrix
 - One row and one column for each data point
 - An entry is the similarity or distance of the associated pair of points
 - **“Incidence” Matrix**
 - One row and one column for each data point
 - An entry is 1 if the associated pair of points belong 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 (**positive** for similarity, **negative** for distance) 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.



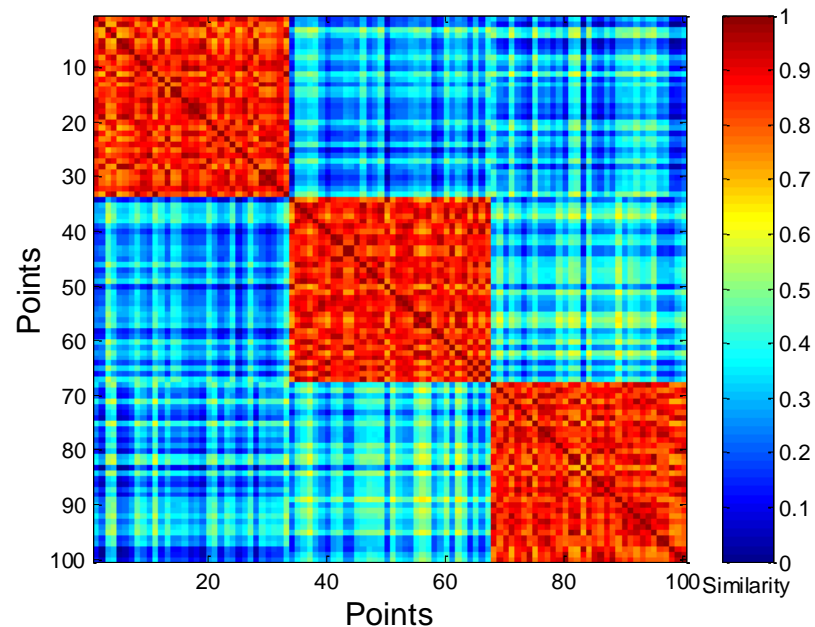
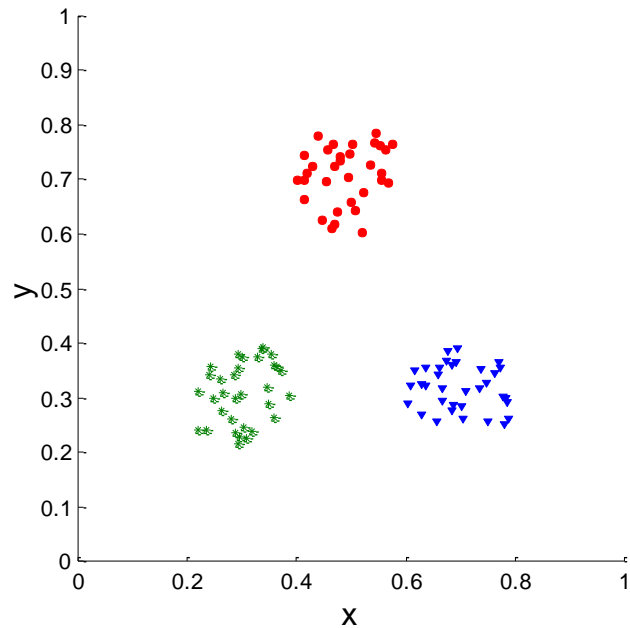
Corr = -0.9235



Corr = -0.5810

Using Similarity Matrix for Cluster Validation

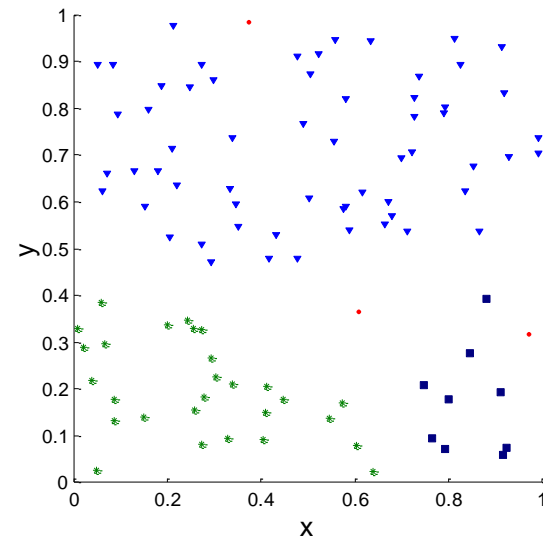
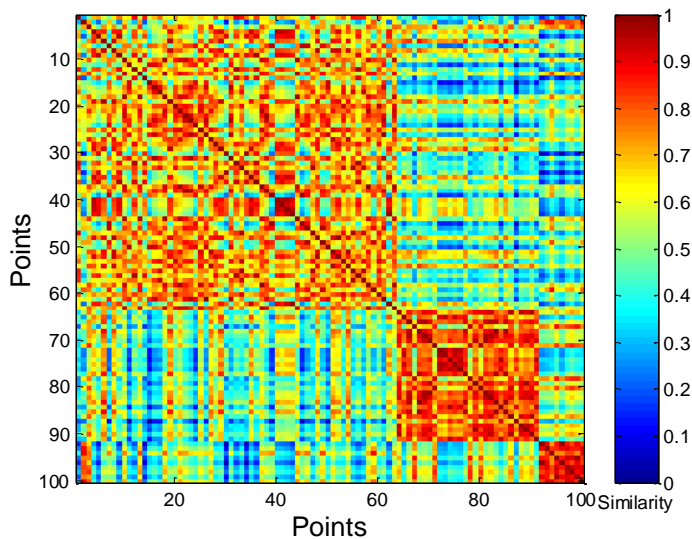
- Order the **similarity** matrix with respect to cluster labels and inspect visually.



$$sim(i,j) = 1 - \frac{d_{ij} - d_{min}}{d_{max} - d_{min}}$$

Using Similarity Matrix for Cluster Validation

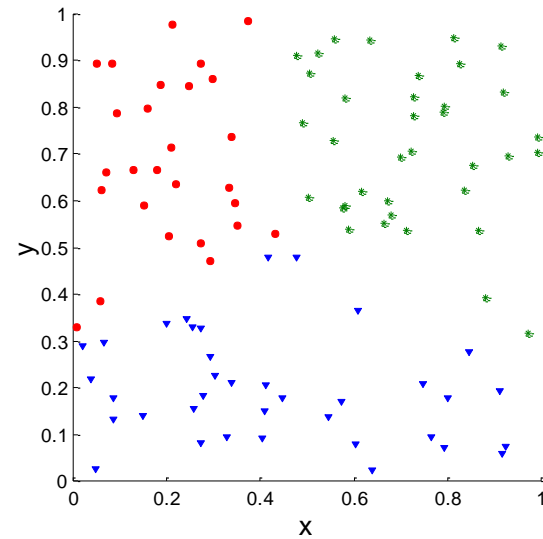
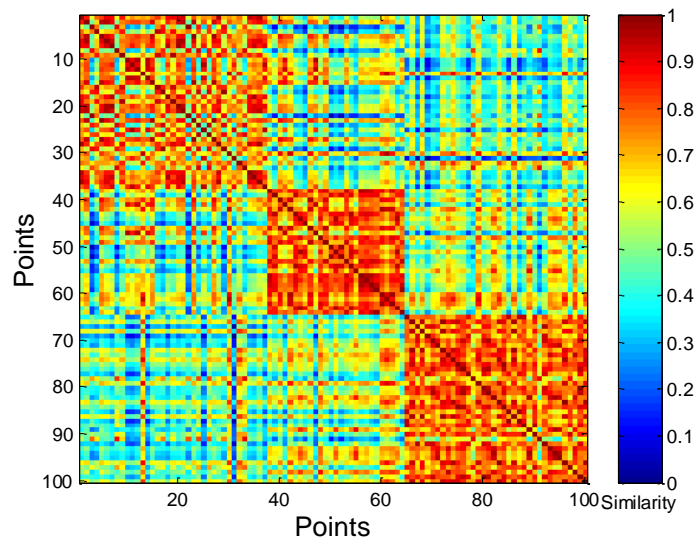
- Clusters in random data are not so crisp



DBSCAN

Using Similarity Matrix for Cluster Validation

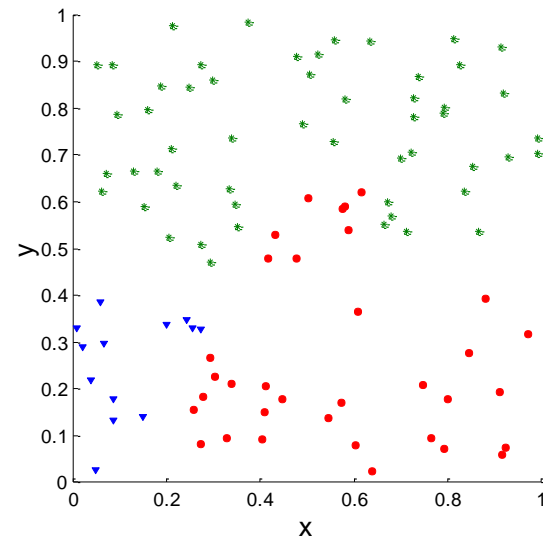
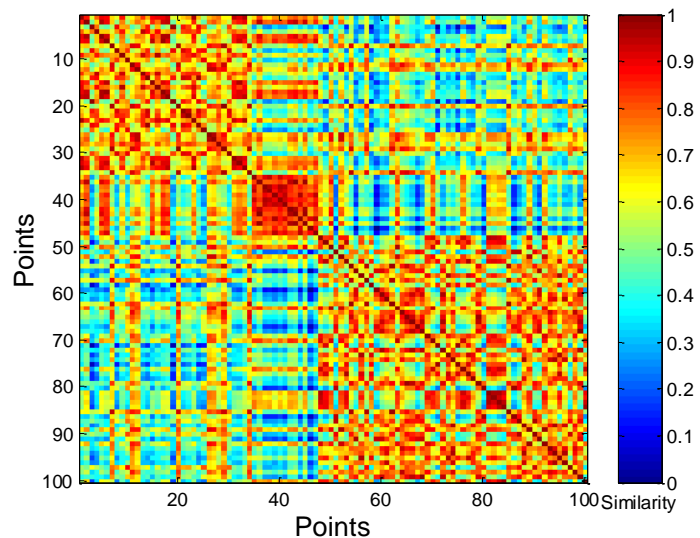
- Clusters in random data are not so crisp



K-means

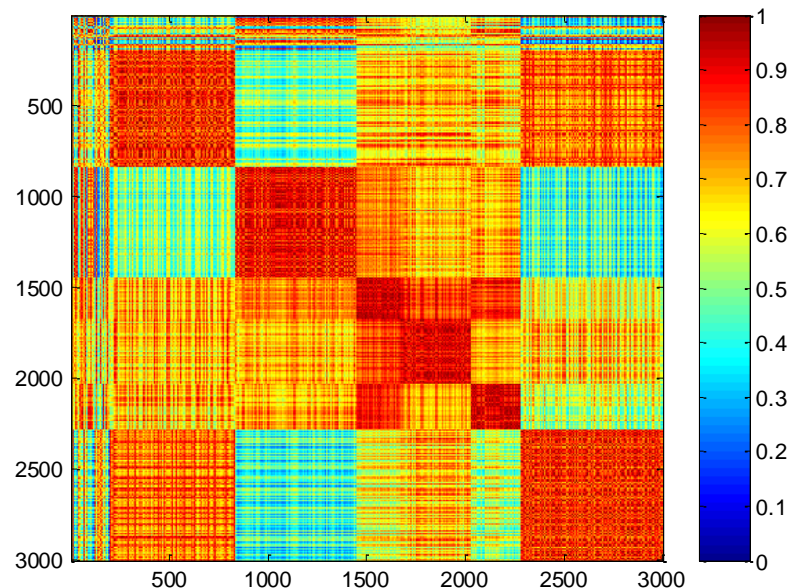
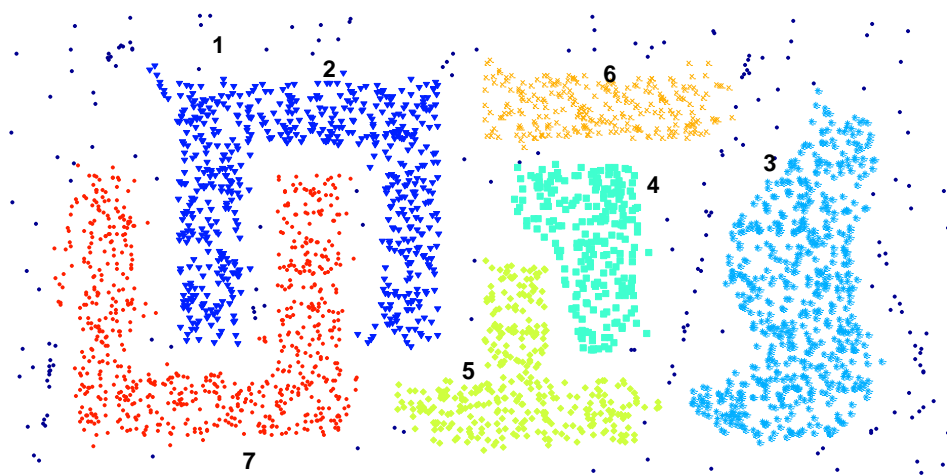
Using Similarity Matrix for Cluster Validation

- Clusters in random data are not so crisp



Complete Link

Using Similarity Matrix for Cluster Validation

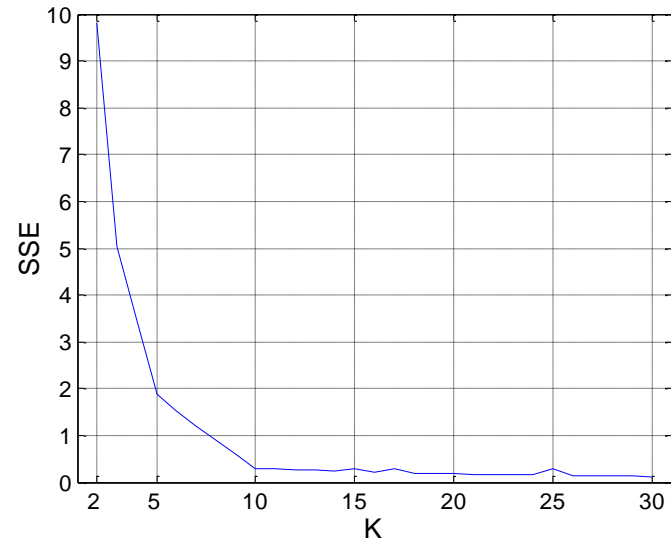
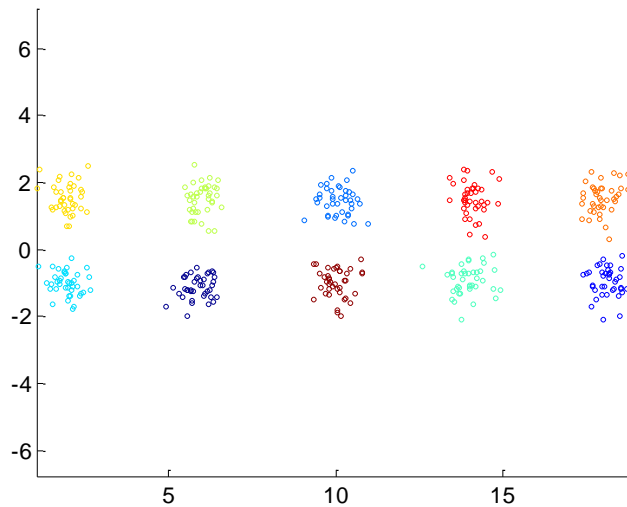


DBSCAN

- Clusters in more complicated figures are not well separated
- This technique can only be used for small datasets since it requires a quadratic computation

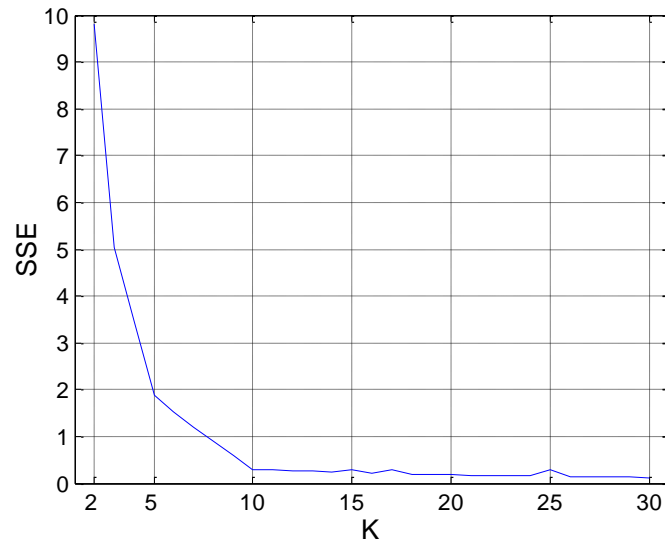
Internal Measures: SSE

- SSE is good for comparing two clusterings or two clusters (average SSE).
- Can also be used to estimate the number of clusters



Estimating the “right” number of clusters

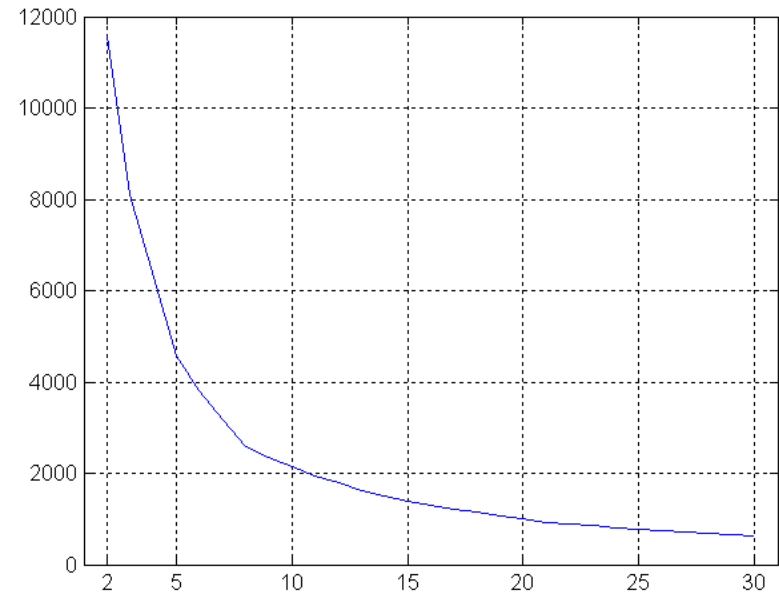
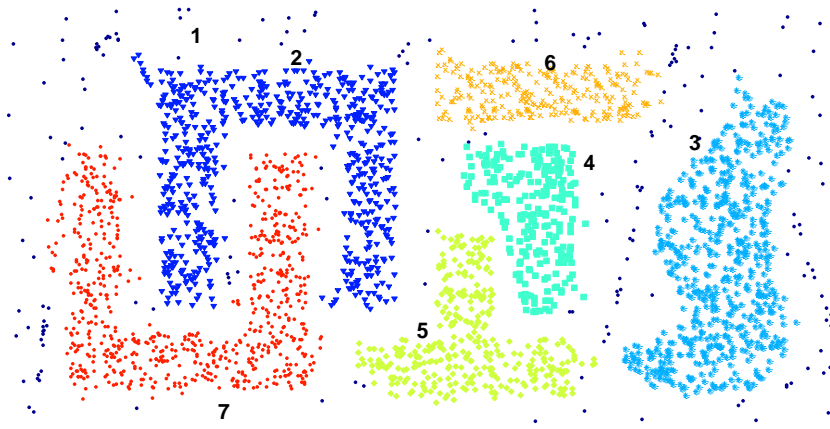
- Typical approach: find a “knee” in an internal measure curve.



- Question: why not the k that **minimizes** the SSE?
- **Desirable property**: the clustering algorithm that does not require the number of clusters to be specified (e.g., DBSCAN)

Internal Measures: SSE

- SSE curve for a more complicated data set



SSE of clusters found using K-means

Internal Measures: Cohesion and Separation

- **Cluster Cohesion:** Measures how closely related are objects in a cluster
- **Cluster Separation:** Measure how distinct or well-separated a cluster is from other clusters
- Example: Squared Error
 - Cohesion is measured by the **within cluster sum of squares** (SSE)

$$WSS = \sum_i \sum_{x \in C_i} (x - c_i)^2$$

We want this to be small

- Separation is measured by the **between cluster sum of squares**

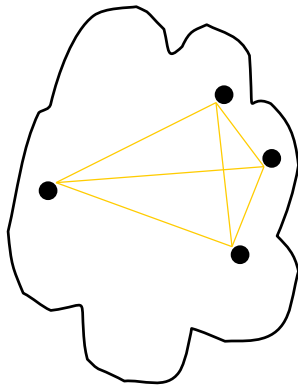
$$BSS = \sum_i m_i (c - c_i)^2$$

We want this to be large

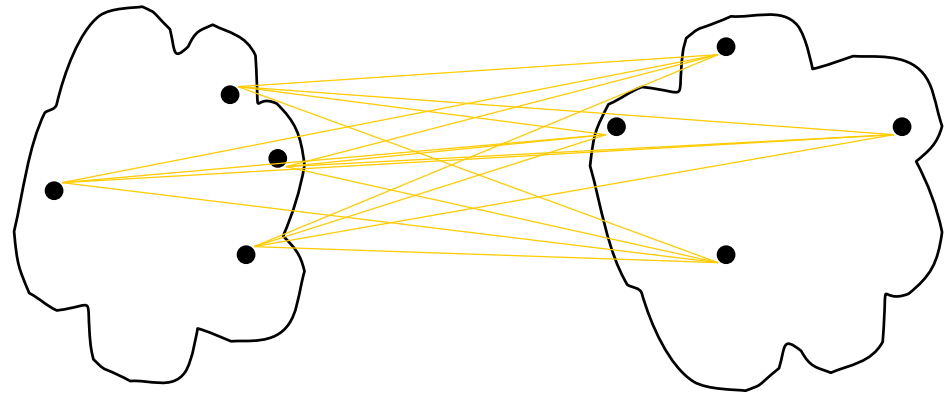
- Where m_i is the size of cluster i

Internal Measures: Cohesion and Separation

- A proximity graph based approach can also be used for cohesion and separation.
 - Cluster cohesion is the sum of the length of all links within a cluster.
 - Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.



cohesion



separation

Internal measures – caveats

- Internal measures have the problem that the clustering algorithm **did not set out to optimize this measure**, so it will not necessarily do well with respect to the measure.
- An internal measure can also be used as an objective function for clustering

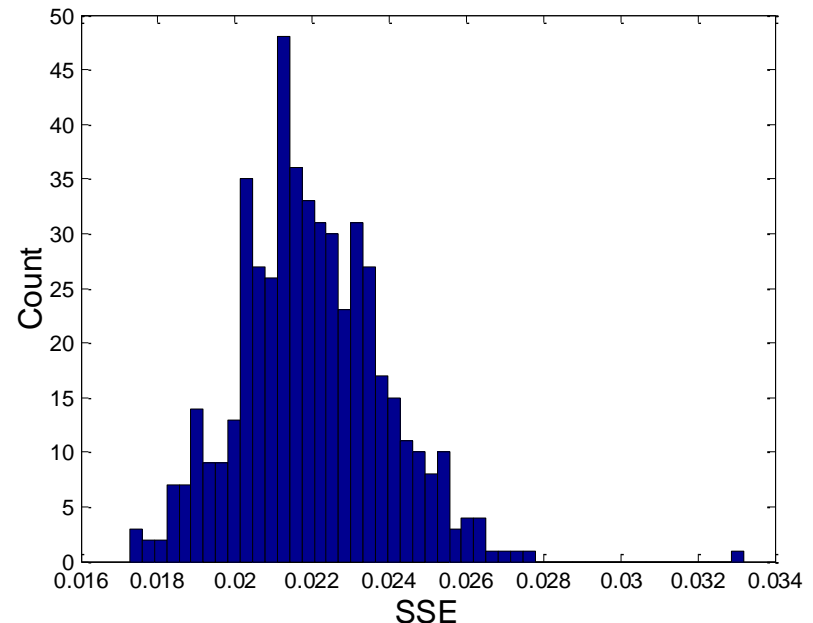
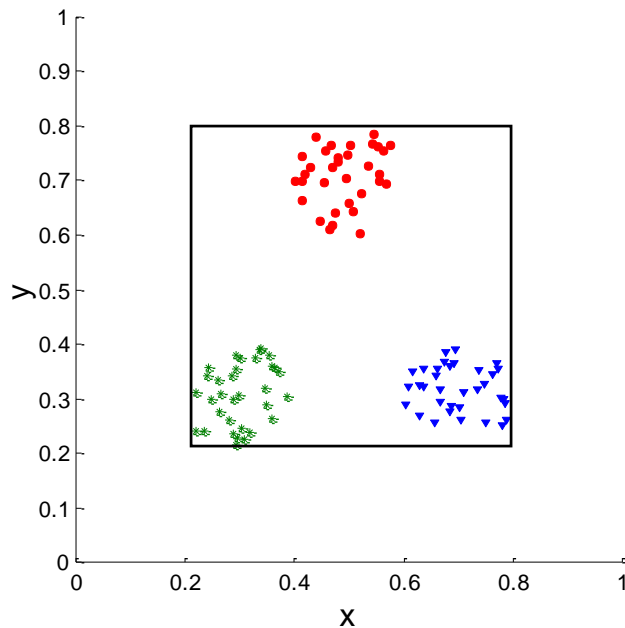
Framework for Cluster Validity

- Need a **framework** to interpret any measure.
 - For example, if our measure of evaluation has the value 10, is that good, fair, or poor?
- **Statistics** provide a framework for cluster validity
 - The more “**non-random**” a clustering result is, the more likely it represents valid structure in the data
 - Can compare the values of an index that result from **random** data or clusterings to those of a clustering result.
 - If the value of the index is **unlikely**, then the cluster results are valid
- For comparing the results of two different sets of cluster analyses, a framework is less necessary.
 - However, there is the question of whether the difference between two index values is **significant**

Statistical Framework for SSE

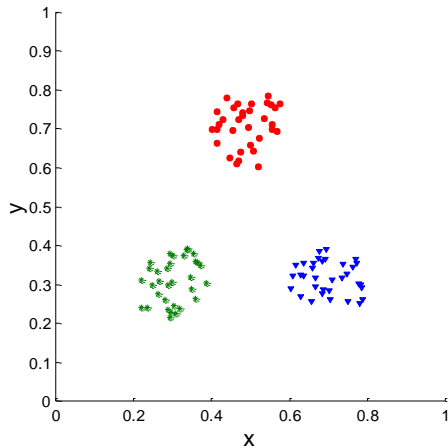
- Example

- Compare SSE of **0.005** against three clusters in random data
- Histogram of SSE for three clusters in 500 random data sets of **100 random points distributed in the range 0.2 – 0.8** for x and y
 - Value 0.005 is very **unlikely**

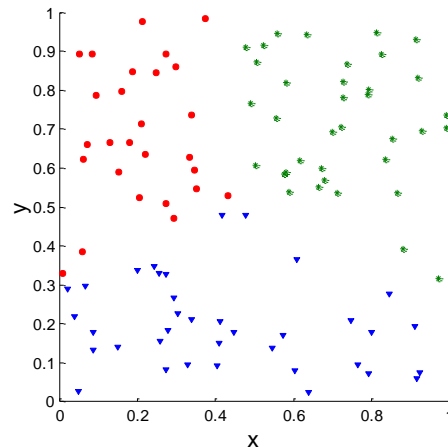


Statistical Framework for Correlation

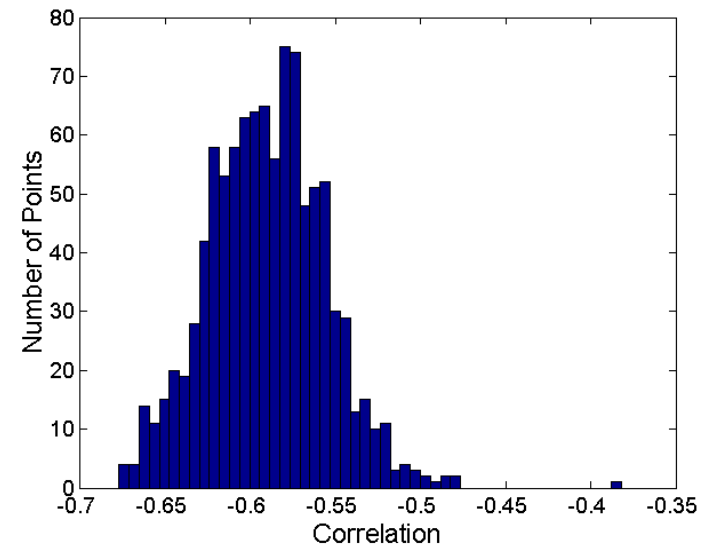
- Correlation of incidence and proximity matrices for the K-means clusterings of the following two data sets.



Corr = -0.9235



Corr = -0.5810



Empirical p-value

- If we have a **measurement v** (e.g., the SSE value)
- ..and we have **N** measurements on **random datasets**
- ...the **empirical p-value** is the **fraction** of measurements in the random data that have value **less or equal** than value **v** (or greater or equal if we want to maximize)
 - i.e., the value in the random dataset is **at least as good** as that in the real data
- We usually require that **$p\text{-value} \leq 0.05$**
- **Hard question**: what is the right notion of a random dataset?

External Measures for Clustering Validity

- Assume that the data is **labeled** with some class labels
 - E.g., **documents** are classified into **topics**, **people** classified according to their **income**, **politicians** classified according to the **political party**.
 - This is called the “**ground truth**”
- In this case we want the clusters to be **homogeneous** with respect to classes
 - **Each cluster** should contain elements of **mostly one class**
 - **Each class** should ideally be assigned to a **single cluster**
- This does not always make sense
 - **Clustering** is not the same as **classification**
 - ...but this is what people use most of the time

Confusion matrix

- n = number of points
- m_i = points in cluster i
- c_j = points in class j
- n_{ij} = points in cluster i coming from class j
- $p_{ij} = n_{ij}/m_i$ = probability of element from cluster i to be assigned in class j

	Class 1	Class 2	Class 3	
Cluster 1	n_{11}	n_{12}	n_{13}	m_1
Cluster 2	n_{21}	n_{22}	n_{23}	m_2
Cluster 3	n_{31}	n_{32}	n_{33}	m_3
	c_1	c_2	c_3	n

	Class 1	Class 2	Class 3	
Cluster 1	p_{11}	p_{12}	p_{13}	m_1
Cluster 2	p_{21}	p_{22}	p_{23}	m_2
Cluster 3	p_{31}	p_{32}	p_{33}	m_3
	c_1	c_2	c_3	n

Measures

	Class 1	Class 2	Class 3	
Cluster 1	p_{11}	p_{12}	p_{13}	m_1
Cluster 2	p_{21}	p_{22}	p_{23}	m_2
Cluster 3	p_{31}	p_{32}	p_{33}	m_3
	c_1	c_2	c_3	n

- **Entropy:**

- Of a **cluster i**: $e_i = -\sum_{j=1}^L p_{ij} \log p_{ij}$
 - Highest when uniform, zero when single class
- Of a clustering: $e = \sum_{i=1}^K \frac{m_i}{n} e_i$

- **Purity:**

- Of a **cluster i**: $p_i = \max_j p_{ij}$
- Of a clustering: $p(C) = \sum_{i=1}^K \frac{m_i}{n} p_i$

Measures

	Class 1	Class 2	Class 3	
Cluster 1	p_{11}	p_{12}	p_{13}	m_1
Cluster 2	p_{21}	p_{22}	p_{23}	m_2
Cluster 3	p_{31}	p_{32}	p_{33}	m_3
	c_1	c_2	c_3	n

- **Precision:**

- Of cluster i with respect to class j : $Prec(i, j) = p_{ij}$

- **Recall:**

- Of cluster i with respect to class j : $Rec(i, j) = \frac{n_{ij}}{c_j}$

- **F-measure:**

- **Harmonic Mean** of Precision and Recall:

$$F(i, j) = \frac{2 * Prec(i, j) * Rec(i, j)}{Prec(i, j) + Rec(i, j)}$$

Measures

Precision/Recall for clusters and clusterings

	Class 1	Class 2	Class 3	
Cluster 1	n_{11}	n_{12}	n_{13}	m_1
Cluster 2	n_{21}	n_{22}	n_{23}	m_2
Cluster 3	n_{31}	n_{32}	n_{33}	m_3
	c_1	c_2	c_3	n

- Assign to cluster i the class k_i such that $k_i = \arg \max_j n_{ij}$

- **Precision:**

- Of cluster i : $Prec(i) = \frac{n_{ik_i}}{m_i}$
- Of the clustering: $Prec(C) = \sum_i \frac{m_i}{n} Prec(i)$

- **Recall:**

- Of cluster i : $Rec(i) = \frac{n_{ik_i}}{c_{k_i}}$
- Of the clustering: $Rec(C) = \sum_i \frac{m_i}{n} Rec(i)$

- **F-measure:**

- **Harmonic Mean** of Precision and Recall

Precision (also called positive predictive value) is the fraction of relevant instances among all positive instances: n of majority class instances / total instances in a cluster

Recall (also known as sensitivity) is the fraction of relevant instances that were positively classified / the total amount of relevant instances – in this case the total number of instances of this class

Good and bad clustering

	Class 1	Class 2	Class 3	
Cluster 1	2	3	85	90
Cluster 2	90	12	8	110
Cluster 3	8	85	7	100
	100	100	100	300

Purity: (0.94, 0.81, 0.85)

– overall 0.86

Precision: (0.94, 0.81, 0.85)

– overall 0.86

Recall: (0.85, 0.9, 0.85)

– overall 0.87

	Class 1	Class 2	Class 3	
Cluster 1	20	35	35	90
Cluster 2	30	42	38	110
Cluster 3	38	35	27	100
	100	100	100	300

Purity: (0.38, 0.38, 0.38)

– overall 0.38

Precision: (0.38, 0.38, 0.38)

– overall 0.38

Recall: (0.35, 0.42, 0.38)

– overall 0.39

Another clustering

	Class 1	Class 2	Class 3	
Cluster 1	0	0	35	35
Cluster 2	50	77	38	165
Cluster 3	38	35	27	100
	100	100	100	300

Cluster 1:
Purity: 1
Precision: 1
Recall: 0.35

Final Comment on Cluster Validity

“The validation of clustering structures is the most difficult and frustrating part of cluster analysis.

Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage.”

Algorithms for Clustering Data, Jain and Dubes