# Clustering
# Lecture 1: Basics

## Jing Gao
SUNY Buffalo

# Outline

- **Basics**
  - Motivation, definition, evaluation
- **Methods**
  - Partitional
  - Hierarchical
  - Density-based
  - Mixture model
  - Spectral methods
- **Advanced topics**
  - Clustering in MapReduce
  - Clustering ensemble
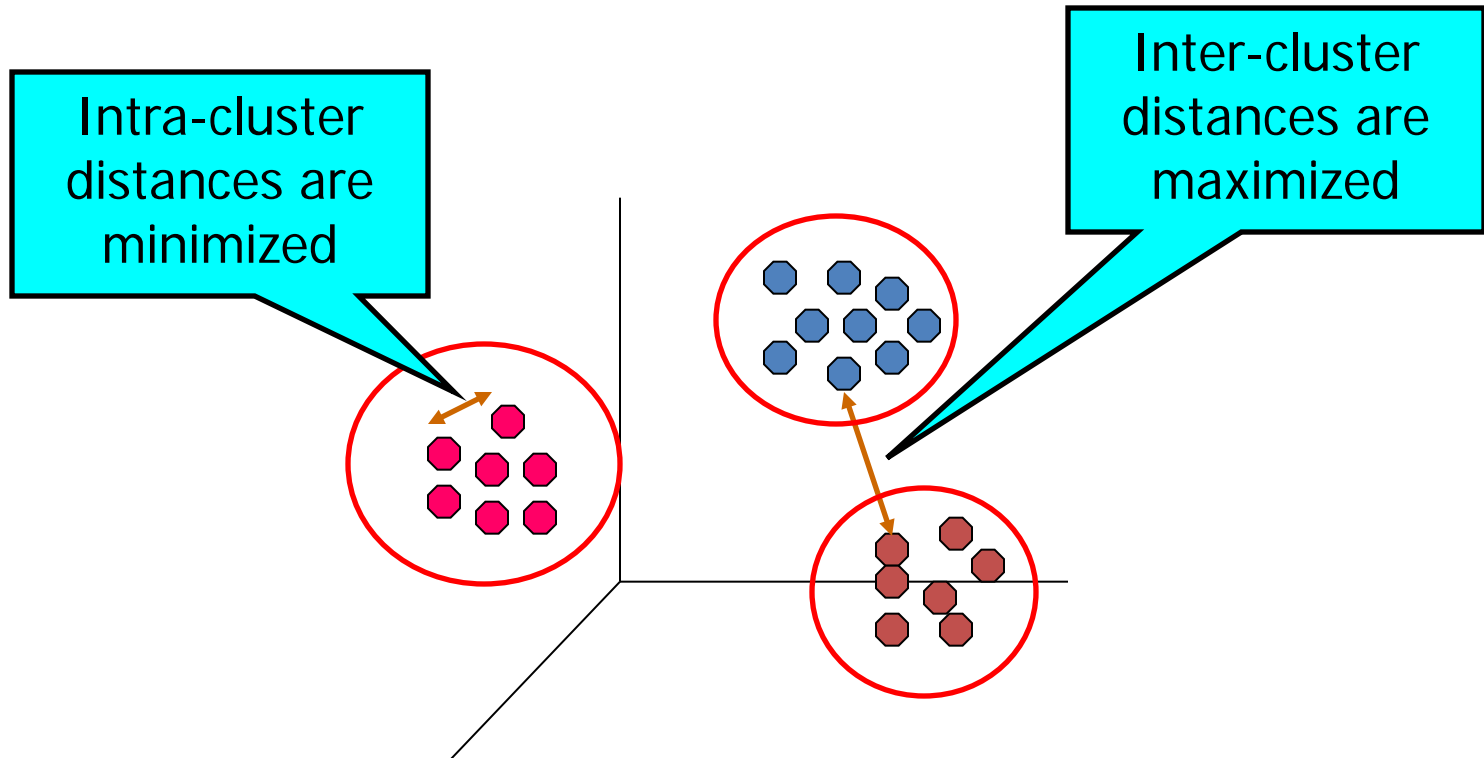  - Semi-supervised clustering, subspace clustering, co-clustering, etc.

# Readings

- Tan, Steinbach, Kumar, Chapters 8 and 9.
- Han, Kamber, Pei. Data Mining: Concepts and Techniques. Chapters 10 and 11.
- Additional readings posted on website

# Clustering Basics

- Definition and Motivation
- Data Preprocessing and Similarity Computation
- Objective of Clustering
- Clustering Evaluation

# Clustering

- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups

Intra-cluster distances are minimized

Inter-cluster distances are maximized
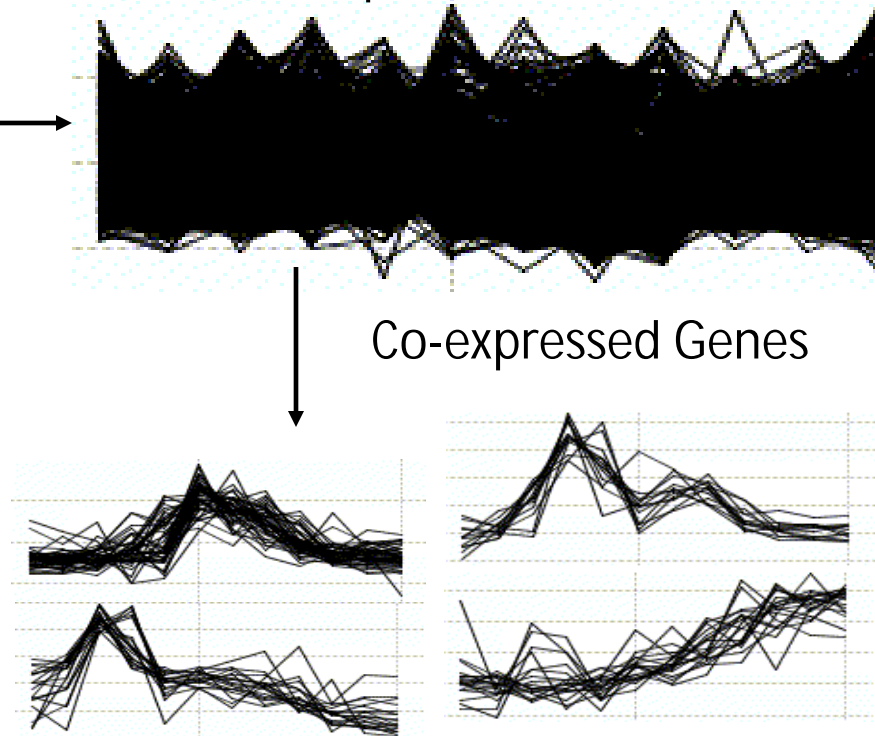
# **Application Examples**

- A stand-alone tool: explore data distribution
- A preprocessing step for other algorithms
- Pattern recognition, spatial data analysis, image processing, market research, WWW, …
  - Cluster documents
  - Cluster web log data to discover groups of similar access patterns

# Clustering Co-expressed Genes

Gene Expression Data Matrix



Gene Expression Patterns



Co-expressed Genes



Why looking for co-expressed genes?
- ¾ *Co-expression indicates co-function;*
- ¾ *Co-expression also indicates co-regulation.*

# Gene-based Clustering



Iyer's data [2]

Examples of co-expressed genes and coherent patterns in gene expression data

¢ [2] Iyer, V.R. et al. The transcriptional program in the response of human fibroblasts to serum. *Science*, 283:83–87, 1999.

# Other Applications

- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs

- City-planning: Identifying groups of houses according to their house type, value, and geographical location

- Climate: understanding earth climate, find patterns of atmosphere and ocean

# Two Important Aspects

- **Properties of input data**
  - Define the similarity or dissimilarity between points
- **Requirement of clustering**
  - Define the objective and methodology

# **Clustering Basics**

- Definition and Motivation
- Data Preprocessing and Distance computation
- Objective of Clustering
- Clustering Evaluation

# Data Representation

- Data: Collection of data objects and their attributes

- An attribute is a property or characteristic of an object
  - Examples: eye color of a person, temperature, etc.
  - Attribute is also known as dimension, variable, field, characteristic, or feature

- A collection of attributes describe an object
  - Object is also known as record, point, case, sample, entity, or instance

Attributes

Objects

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |

# Data Matrix

- ## Represents *n* objects with *p* attributes
  - An *n* by *p* matrix

Attributes

Objects

$$
\begin{bmatrix}
x_{11} & \cdots & x_{1f} & \cdots & x_{1p} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
x_{i1} & \cdots & x_{if} & \cdots & x_{ip} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
x_{n1} & \cdots & x_{nf} & \cdots & x_{np}
\end{bmatrix}
$$

The value of the *i*-th object on the *f*-th attribute

# Gene Expression Data

| | condition 1 | condition 2 | condition 3 | condition 4 | condition… |
|---|---|---|---|---|---|
| gene 1 | 0.13 | 0.72 | 0.1 | 0.57 | |
| gene 2 | 0.34 | 1.58 | 1.05 | 1.15 | |
| gene 3 | 0.43 | 1.1 | 0.97 | 1 | |
| gene 4 | 1.22 | 0.97 | 1 | 0.85 | |
| gene 5 | -0.89 | 1.21 | 1.29 | 1.08 | |
| gene 6 | 1.1 | 1.45 | 1.44 | 1.12 | |
| gene 7 | 0.83 | 1.15 | 1.1 | 1 | |
| gene 8 | 0.87 | 1.32 | 1.35 | 1.13 | |
| gene 9 | -0.33 | 1.01 | 1.38 | 1.21 | |
| gene 10 | 0.10 | 0.85 | 1.03 | 1 | |
| gene … | | | | | |

- ## Clustering genes

- Genes are objects

- Experiment conditions are attributes

- Find genes with similar behavior

# Similarity and Dissimilarity

- ## Similarity
  - Numerical measure of how alike two data objects are
  - Is higher when objects are more alike
  - Often falls in the range [0,1]
- ## Dissimilarity
  - Numerical measure of how different are two data objects
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies

# Types of Attributes

- ## Discrete
  - Has only a finite or countably infinite set of values
  - Examples: zip codes, counts, or the set of words in a collection of documents
  - Note: binary attributes are a special case of discrete attributes

- ## Ordinal
  - Has only a finite or countably infinite set of values
  - Order of values is important
  - Examples: rankings (e.g., pain level 1-10), grades (A, B, C, D)

- ## Continuous
  - Has real numbers as attribute values
  - Examples: temperature, height, or weight
  - Continuous attributes are typically represented as floating-point variables

# Similarity/Dissimilarity for Simple Attributes

$p$ and $q$ are the attribute values for two data objects.

| Attribute Type | Dissimilarity | Similarity |
|---|---|---|
| Discrete | $d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$ | $s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$ |
| Ordinal | $d = \frac{\|p-q\|}{n-1}$ <br> (values mapped to integers 0 to $n-1$, where $n$ is the number of values) | $s = 1 - \frac{\|p-q\|}{n-1}$ |
| Continuous | $d = \|p - q\|$ | $s = -d,\ s = \frac{1}{1+d}$ or <br> $s = 1 - \frac{d - min\_d}{max\_d - min\_d}$ |

Dissimilarity and similarity between $p$ and $q$

# Distance Matrix

- ## Represents pairwise distance in *n* objects
  - An *n* by *n* matrix
  - $d(i,j)$: distance or dissimilarity between objects *i* and *j*
  - Nonnegative
  - Close to 0: similar

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \cdots & \cdots & 0 \end{bmatrix}$$

# Data Matrix -> Distance Matrix

| | s 1 | s 2 | s 3 | s 4 | ... |
|---|---|---|---|---|---|
| g 1 | 0.13 | 0.72 | 0.1 | 0.57 | |
| g 2 | 0.34 | 1.58 | 1.05 | 1.15 | |
| g 3 | 0.43 | 1.1 | 0.97 | 1 | |
| g 4 | 1.22 | 0.97 | 1 | 0.85 | |
| g 5 | -0.89 | 1.21 | 1.29 | 1.08 | |
| g 6 | 1.1 | 1.45 | 1.44 | 1.12 | |
| g 7 | 0.83 | 1.15 | 1.1 | 1 | |
| g 8 | 0.87 | 1.32 | 1.35 | 1.13 | |
| g 9 | -0.33 | 1.01 | 1.38 | 1.21 | |
| g 10 | 0.10 | 0.85 | 1.03 | 1 | |
| ... | | | | | |

Original Data Matrix

| | g 1 | g 2 | g 3 | g 4 | ... |
|---|---|---|---|---|---|
| g 1 | 0 | $d(1,2)$ | $d(1,3)$ | $d(1,4)$ | |
| g 2 | | 0 | $d(2,3)$ | $d(2,4)$ | |
| g 3 | | | 0 | $d(3,4)$ | |
| g 4 | | | | 0 | |
| ... | | | | | |

Distance Matrix

# Minkowski Distance—Continuous Attribute

- Minkowski distance: a generalization

$$d(i,j) = \sqrt[q]{|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + ... + |x_{ip} - x_{jp}|^q} \quad (q > 0)$$

- If $q = 2$, d is Euclidean distance
- If $q = 1$, d is Manhattan distance

$X_i$ (1,7)

8.48

$q=2$

$X_j$ (7,1)

$X_i$

12

6

6

$q=1$

$X_j$

# Standardization

- Calculate the mean absolute deviation

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + ... + x_{nf})$$

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + ... + |x_{nf} - m_f|)$$

- Calculate the standardized measurement (z-score)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

# Mahalanobis Distance

$$d(p,q) = (p - q) \, \Sigma^{-1} \, (p - q)^T$$



$\Sigma$ is the covariance matrix of the input data $X$

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^{n} (X_{ij} - \overline{X}_j)(X_{ik} - \overline{X}_k)$$

For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.

# Mahalanobis Distance



Covariance Matrix:

$$S = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

A: (0.5, 0.5)

B: (0, 1)

C: (1.5, 1.5)

Mahal(A,B) = 5

Mahal(A,C) = 4

# Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties

  1. $d(p, q) \geq 0$ for all $p$ and $q$ and $d(p, q) = 0$ only if $p = q$. (Positive definiteness)
  2. $d(p, q) = d(q, p)$ for all $p$ and $q$. (Symmetry)
  3. $d(p, r) \leq d(p, q) + d(q, r)$ for all points $p$, $q$, and $r$. (Triangle Inequality)

  where $d(p, q)$ is the distance (dissimilarity) between points (data objects), $p$ and $q$.

- A distance that satisfies these properties is a metric

# Similarity for Binary Attributes

- Common situation is that objects, *p* and *q*, have only binary attributes

- Compute similarities using the following quantities
  $M_{01}$ = the number of attributes where p was 0 and q was 1
  $M_{10}$ = the number of attributes where p was 1 and q was 0
  $M_{00}$ = the number of attributes where p was 0 and q was 0
  $M_{11}$ = the number of attributes where p was 1 and q was 1

- Simple Matching and Jaccard Coefficients

  SMC = number of matches / total number of attributes
  $$= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$$

  J = number of matches / number of not-both-zero attributes values
  $$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$

# SMC versus Jaccard: Example

$p =$ 1 0 0 0 0 0 0 0 0 0

$q =$ 0 0 0 0 0 0 1 0 0 1

$M_{01} = 2$  (the number of attributes where p was 0 and q was 1)
$M_{10} = 1$  (the number of attributes where p was 1 and q was 0)
$M_{00} = 7$  (the number of attributes where p was 0 and q was 0)
$M_{11} = 0$  (the number of attributes where p was 1 and q was 1)

SMC = $(M_{11} + M_{00})/(M_{01} + M_{10} + M_{11} + M_{00})$ = (0+7) / (2+1+0+7) = 0.7

J = $(M_{11}) / (M_{01} + M_{10} + M_{11})$ = 0 / (2 + 1 + 0) = 0

# Document Data

- Each document becomes a `term' vector,
  - each term is a component (attribute) of the vector,
  - the value of each component is the number of times the corresponding term occurs in the document.

| | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

# Cosine Similarity

- If $d_1$ and $d_2$ are two document vectors, then

$$\cos(d_1, d_2) = (d_1 \cdot d_2) / \|d_1\| \, \|d_2\| \,,$$

  where $\cdot$ indicates vector dot product and $\| d \|$ is the length of vector $d$.

- Example:

$$d_1 = 3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0$$
$$d_2 = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2$$

$d_1 \cdot d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$

$\|d_1\| = (3*3+2*2+0*0+5*5+0*0+0*0+0*0+2*2+0*0+0*0)^{0.5} = (42)^{0.5} = 6.481$

$\|d_2\| = (1*1+0*0+0*0+0*0+0*0+0*0+0*0+1*1+0*0+2*2)^{0.5} = (6)^{0.5} = 2.245$

$$\cos(d_1, d_2) = .3150$$

# Correlation

- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects, p and q, and then take their dot product (continuous attributes)

$$p'_k = (p_k - mean(p)) / std(p)$$

$$q'_k = (q_k - mean(q)) / std(q)$$

$$s(p,q) = p' \cdot q'$$

# Common Properties of a Similarity

- ## Similarities, also have some well known properties.

    1.  *$s(p, q) = 1$ (or maximum similarity) only if $p = q$.*

    2.  *$s(p, q) = s(q, p)$* for all *$p$* and *$q$*. (Symmetry)

    where *$s(p, q)$* is the similarity between points (data objects), *$p$* and *$q$*.

# Characteristics of the Input Data Are Important

- Sparseness
- Attribute type
- Type of Data
- Dimensionality
- Noise and Outliers
- Type of Distribution
- => Conduct preprocessing and select the appropriate dissimilarity or similarity measure
- => Determine the objective of clustering and choose the appropriate method

# Clustering Basics

- Definition and Motivation
- Data Preprocessing and Distance computation
- Objective of Clustering
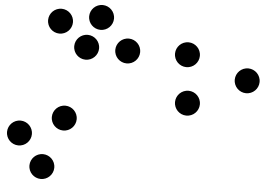- Clustering Evaluation

# Considerations for Cluster Analysis

- **Partitioning criteria**
  - Single level vs. hierarchical partitioning (often, multi-level hierarchical partitioning is desirable)

- **Separation of clusters**
  - Exclusive (e.g., one customer belongs to only one region) vs. overlapping (e.g., one document may belong to more than one topic)

- **Hard versus fuzzy**

  - In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
  - Weights must sum to 1
  - Probabilistic clustering has similar characteristics

- **Similarity measure and data types**

- **Heterogeneous versus homogeneous**

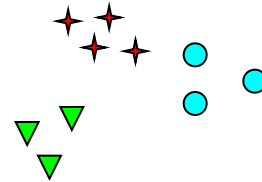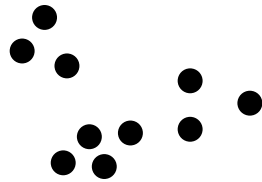  - Cluster of widely different sizes, shapes, and densities

# **Requirements of Clustering**

- Scalability
- Ability to deal with different types of attributes
- Minimal requirements for domain knowledge to determine input parameters
- Able to deal with noise and outliers
- Discovery of clusters with arbitrary shape
- Insensitive to order of input records
- High dimensionality
- Incorporation of user-specified constraints
- Interpretability and usability
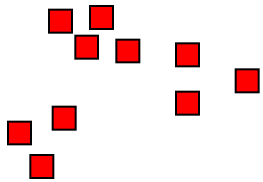- <span style="color:red">**What clustering results we want to get?**</span>
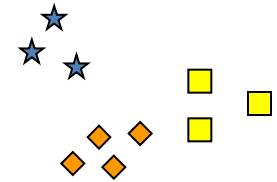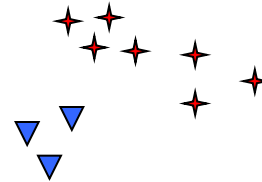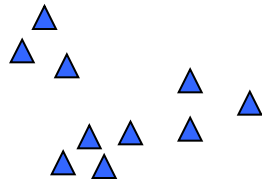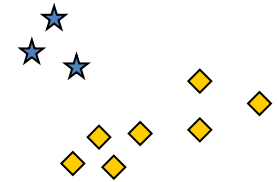
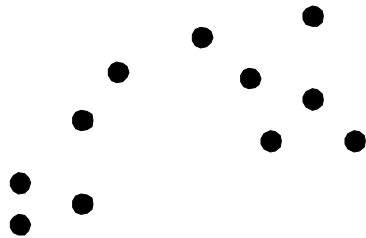# Notion of a Cluster can be Ambiguous
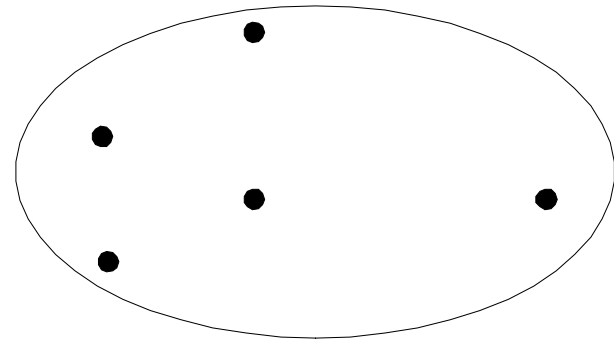
How many clusters?

Six Clusters

Two Clusters

Four Clusters

# Partitional Clustering

Input Data

A Partitional Clustering

# Hierarchical Clustering



Clustering Solution 1



Clustering Solution 2

# Types of Clusters: Center-Based

- ## Center-based

  - A cluster is a set of objects such that an object in a cluster is closer (more similar) to the "center" of a cluster, than to the center of any other cluster

  - The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most "representative" point of a cluster

4 center-based clusters

# Types of Clusters: Density-Based

- ## Density-based
  - A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
  - Used when the clusters are irregular or intertwined, and when noise and outliers are present.

6 density-based clusters

# Clustering Basics

- Definition and Motivation
- Data Preprocessing and Distance computation
- Objective of Clustering
- Clustering Evaluation

# Cluster Validation

- **Cluster validation**
  - Quality: "goodness" of clusters
  - Assess the quality and reliability of clustering results

- **Why validation?**
  - To avoid finding clusters formed by chance
  - To compare clustering algorithms
  - To choose clustering parameters
    - e.g., the number of clusters

# Aspects of Cluster Validation

- Comparing the clustering results to *ground truth* (externally known results)

    – External Index

- Evaluating the quality of clusters *without* reference to external information

    – Use only the data

    – Internal Index

- Determining the *reliability* of clusters

    – To what confidence level, the clusters are not formed by chance

    – Statistical framework

# Comparing to Ground Truth

- ## Notation
  - N: number of objects in the data set
  - $P=\{P_1,...,P_s\}$: the set of "ground truth" clusters
  - $C=\{C_1,...,C_t\}$: the set of clusters reported by a clustering algorithm
- ## The "incidence matrix"
  - $N \times N$ (both rows and columns correspond to objects)
  - $P_{ij} = 1$ if $O_i$ and $O_j$ belong to the same "ground truth" cluster in $P$; $P_{ij}=0$ otherwise
  - $C_{ij} = 1$ if $O_i$ and $O_j$ belong to the same cluster in $C$; $C_{ij}=0$ otherwise

# Rand Index and Jaccard Coefficient

- A pair of data object ($O_i$, $O_j$) falls into one of the following categories
    - SS: $C_{ij}$=1 and $P_{ij}$=1;        (agree)
    - DD: $C_{ij}$=0 and $P_{ij}$=0;        (agree)
    - SD: $C_{ij}$=1 and $P_{ij}$=0;        (disagree)
    - DS: $C_{ij}$=0 and $P_{ij}$=1;        (disagree)

- **Rand index**
$$Rand = \frac{|Agree|}{|Agree|+|Disagree|} = \frac{|SS|+|DD|}{|SS|+|SD|+|DS|+|DD|}$$

    - may be dominated by DD

- **Jaccard Coefficient**
$$Jaccard\ coefficient = \frac{|SS|}{|SS|+|SD|+|DS|}$$

Clustering

| | g 1 | g 2 | g 3 | g 4 | g 5 |
|---|---|---|---|---|---|
| g 1 | 1 | 1 | 1 | 0 | 0 |
| g 2 | 1 | 1 | 1 | 0 | 0 |
| g 3 | 1 | 1 | 1 | 0 | 0 |
| g 4 | 0 | 0 | 0 | 1 | 1 |
| g 5 | 0 | 0 | 0 | 1 | 1 |

Groundtruth

| | g 1 | g 2 | g 3 | g 4 | g 5 |
|---|---|---|---|---|---|
| g 1 | 1 | 1 | 0 | 0 | 0 |
| g 2 | 1 | 1 | 0 | 0 | 0 |
| g 3 | 0 | 0 | 1 | 1 | 1 |
| g 4 | 0 | 0 | 1 | 1 | 1 |
| g 5 | 0 | 0 | 1 | 1 | 1 |

Ground truth →

Clustering

| | | Same Cluster | Different Cluster |
|---|---|---|---|
| | Same Cluster | 9 | 4 |
| | Different Cluster | 4 | 8 |

$$Rand = \frac{|SS|+|DD|}{|SS|+|SD|+|DS|+|DD|} = \frac{17}{25}$$

$$Jaccard = \frac{|SS|}{|SS|+|SD|+|DS|} = \frac{9}{17}$$

# Entropy and Purity

- ## Notation

  - $|C_k \cap P_j|$  the number of objects in both the *k*-th cluster of the clustering solution and *j*-th cluster of the groundtruth
  - $|C_k|$ the number of objects in the *k*-th cluster of the clustering solution
  - $|P_j|$ the number of objects in the *j*-th cluster of the groundtruth

- ## Purity

$$Purity = \frac{1}{N} \sum_k \max_j |C_k \cap P_j|$$

- ## Normalized Mutual Information

$$NMI = \frac{I(C,P)}{\sqrt{H(C)H(P)}} \qquad I(C,P) = \sum_k \sum_j \frac{|C_k \cap P_j|}{N} \log \frac{N \times |C_k \cap P_j|}{|C_k||P_j|}$$

$$H(C) = \sum_k \frac{|C_k|}{N} \log \frac{|C_k|}{N} \qquad H(P) = \sum_j \frac{|P_j|}{N} \log \frac{|P_j|}{N}$$

# Example

|       | P 1 | P 2 | P 3 | P 4 | P5 | P6  | Total |
|-------|-----|-----|-----|-----|-----|-----|-------|
| C1    | 3   | 5   | 40  | 506 | 96  | 27  | 677   |
| C 2   | 4   | 7   | 280 | 29  | 39  | 2   | 361   |
| C 3   | 1   | 1   | 1   | 7   | 4   | 671 | 685   |
| C 4   | 10  | 162 | 3   | 119 | 73  | 2   | 369   |
| C 5   | 331 | 22  | 5   | 70  | 13  | 23  | 464   |
| C 6   | 5   | 358 | 12  | 212 | 48  | 13  | 648   |
| total | 354 | 555 | 341 | 943 | 273 | 738 | 3204  |

$$Purity = \frac{1}{N} \sum_k \max_j |C_k \cap P_j|$$

$$Purity = \frac{506 + 280 + 671 + 162 + 331 + 358}{3204}$$

$$= 0.7203$$

$$NMI = \frac{I(C,P)}{\sqrt{H(C)H(P)}} \qquad I(C,P) = \sum_k \sum_j \frac{|C_k \cap P_j|}{N} \log \frac{N \times |C_k \cap P_j|}{|C_k||P_j|}$$

$$H(C) = \sum_k \frac{|C_k|}{N} \log \frac{|C_k|}{N} \qquad H(P) = \sum_j \frac{|P_j|}{N} \log \frac{|P_j|}{N}$$

# Internal Index

- "Ground truth" may be unavailable

- Use only the data to measure cluster quality
  - Measure the "*cohesion*" and "*separation*" of clusters
  - Calculate the *correlation* between clustering results and distance matrix

# Cohesion and Separation

- **Cohesion** is measured by the within cluster sum of squares

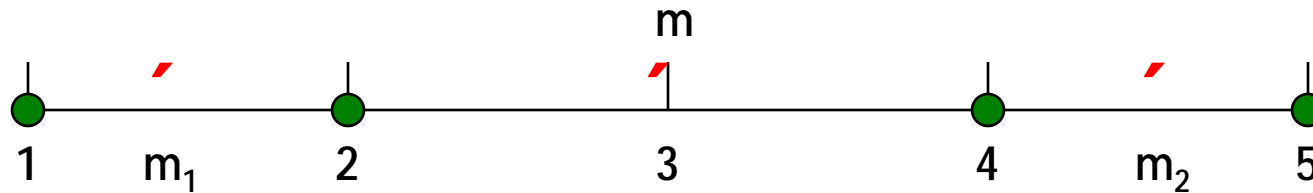$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

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

$$BSS = \sum_i |C_i|(m - m_i)^2$$

  where $|Ci|$ is the size of cluster $i$, m is the centroid of the whole data set

- BSS + WSS = constant
- WSS (Cohesion) measure is called Sum of Squared Error (SSE)—a commonly used measure
- A larger number of clusters tend to result in smaller SSE

# Example



K=1 :

$$WSS = (1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2 = 10$$
$$BSS = 4 \times (3-3)^2 = 0$$
$$Total = 10 + 0 = 10$$

K=2 :

$$WSS = (1-1.5)^2 + (2-1.5)^2 + (4-4.5)^2 + (5-4.5)^2 = 1$$
$$BSS = 2 \times (3-1.5)^2 + 2 \times (4.5-3)^2 = 9$$
$$Total = 1 + 9 = 10$$

K=4:

$$WSS = (1-1)^2 + (2-2)^2 + (4-4)^2 + (5-5)^2 = 0$$
$$BSS = 1 \times (1-3)^2 + 1 \times (2-3)^2 + 1 \times (4-3)^2 + 1 \times (5-3)^2 = 10$$
$$Total = 0 + 10 = 10$$

# Silhouette Coefficient

- Silhouette Coefficient combines ideas of both cohesion and separation

- For an individual point, $i$
    - Calculate $a$ = average distance of $i$ to the points in its cluster
    - Calculate $b$ = min (average distance of $i$ to points in another cluster)
    - The **silhouette coefficient** for a point is then given by

      s = 1 – a/b   if a < b,   (s = b/a - 1   if a ³ b, not the usual case)

    - Typically between 0 and 1
    - The closer to 1 the better



- Can calculate the Average Silhouette width for a cluster or a clustering

# Correlation with Distance Matrix

- Distance Matrix
  - $D_{ij}$ is the similarity between object $O_i$ and $O_j$
- Incidence Matrix
  - $C_{ij}=1$ if $O_i$ and $O_j$ belong to the same cluster, $C_{ij}=0$ otherwise
- Compute the correlation between the two matrices
  - Only $n(n\text{-}1)/2$ entries needs to be calculated
- High correlation indicates good clustering

# Correlation with Distance Matrix

- Given Distance Matrix D = $\{d_{11}, d_{12}, ..., d_{nn}\}$ and Incidence Matrix $C = \{c_{11}, c_{12}, ..., c_{nn}\}$.

- Correlation $r$ between $D$ and $C$ is given by

$$r = \frac{\sum\limits_{i=1, j=1}^{n} (d_{ij} - \bar{d})(c_{ij} - \bar{c})}{\sqrt{\sum\limits_{i=1, j=1}^{n} (d_{ij} - \bar{d})^2} \sqrt{\sum\limits_{i=1, j=1}^{n} (c_{ij} - \bar{c})^2}}$$
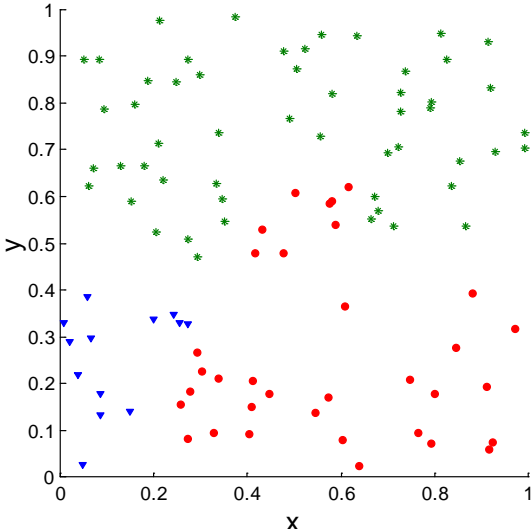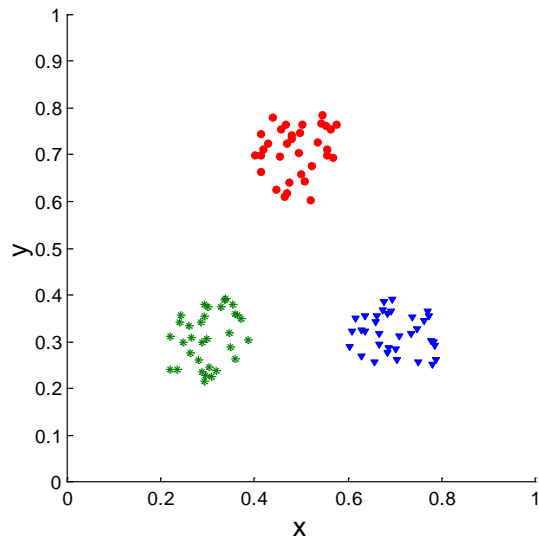
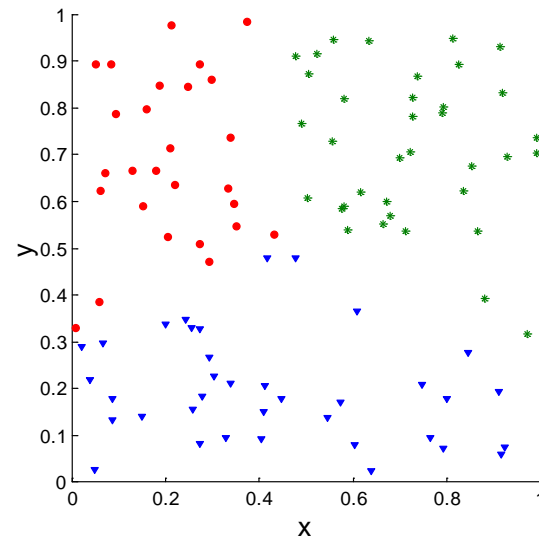# Are There Clusters in the Data?



Random Points

DBSCAN

K-means

Complete Link

# Measuring Cluster Validity Via Correlation

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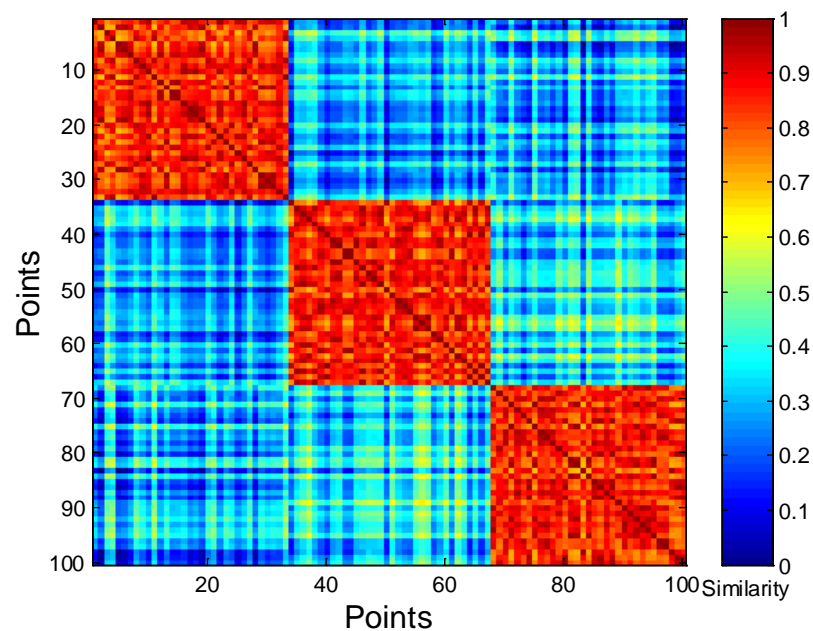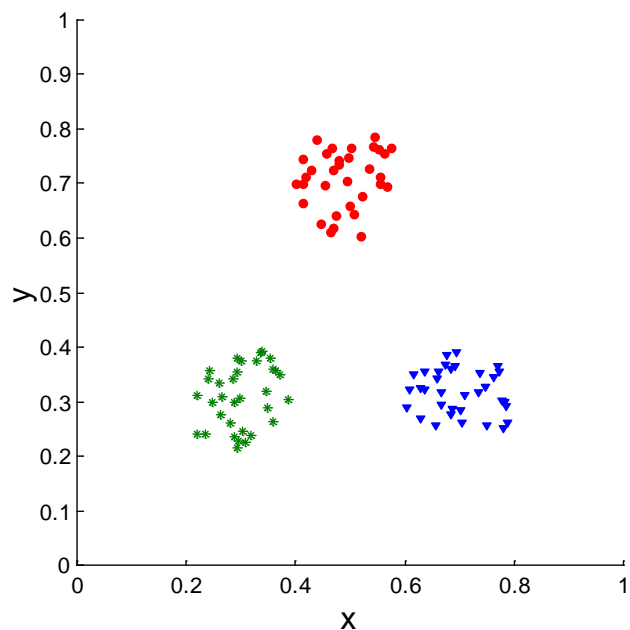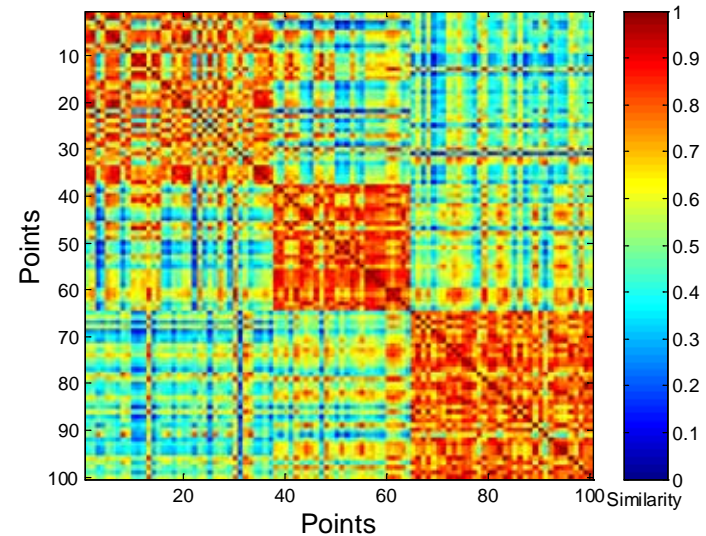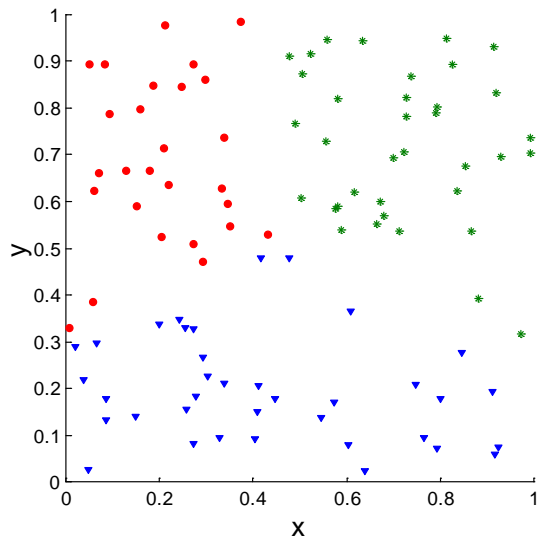


Corr = -0.9235

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

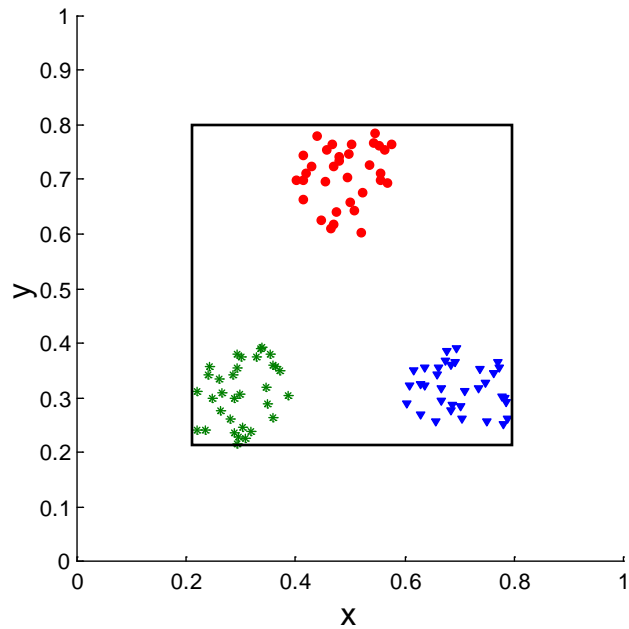- Clusters in random data are not so crisp

# Reliability of Clusters

- 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 "atypical" a clustering result is, the more likely it represents valid structure in the data
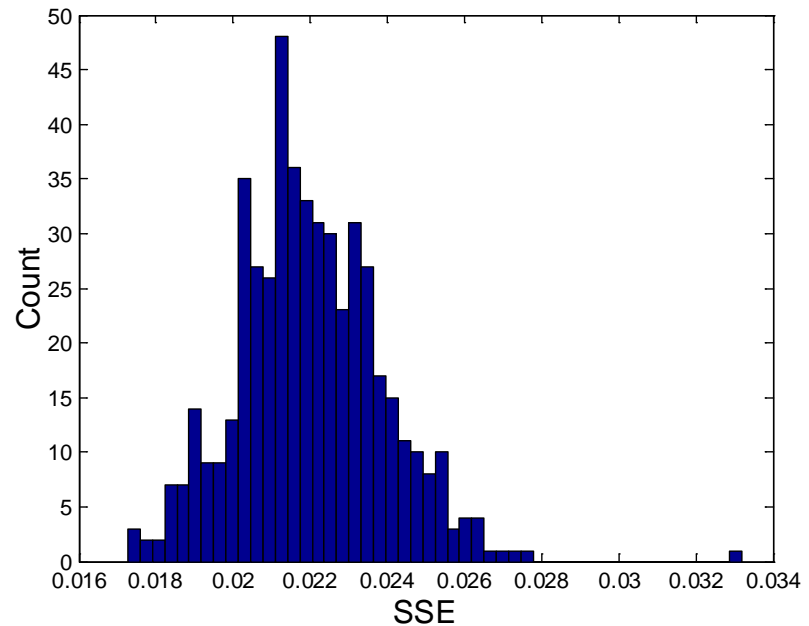
# Statistical Framework for SSE

- ## Example
  - – Compare SSE of 0.005 against three clusters in random data
  - – SSE Histogram of 500 sets of random data points of size 100—lowest SSE is 0.0173
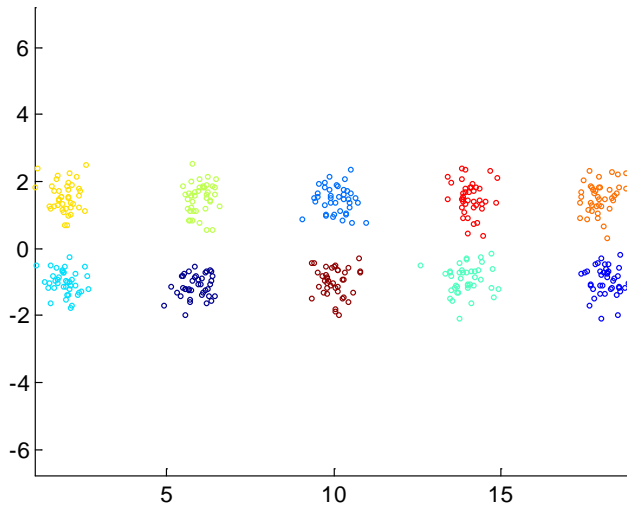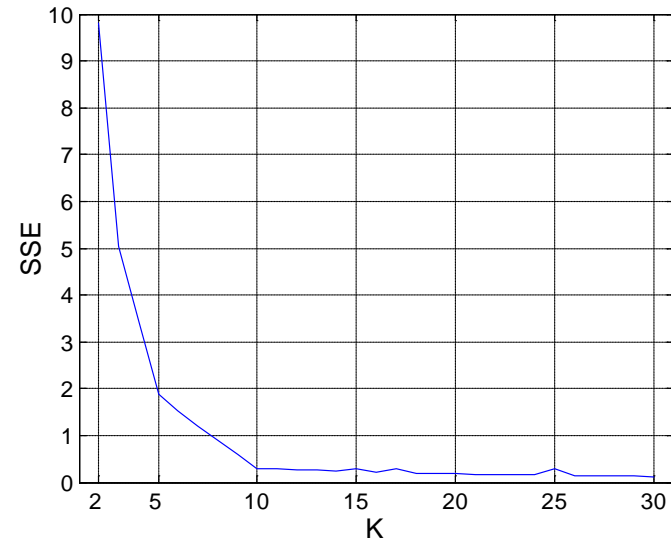


SSE = 0.005

# Determine the Number of Clusters Using SSE

- SSE curve



Clustering of Input Data



SSE wrt K

# Take-away Message

- What's clustering?
- Why clustering is important?
- How to preprocess data and compute dissimilarity/similarity from data?
- What's a good clustering solution?
- How to evaluate the clustering results?