

### 1.1. What is Cluster Analysis

#### **What is a cluster?**

- A cluster is a collection of data objects which are
  - Similar (or related) to one another within the same group (i.e., cluster)
  - Dissimilar (or unrelated) to the objects in other groups (i.e., clusters)

#### **Cluster analysis (or *clustering, data segmentation, ...*)**

- Given a set of data points, partition them into a set of groups (i.e., clusters) which are as similar as possible
- Cluster analysis is **unsupervised learning** (i.e., no predefined classes)
  - This contrasts with *classification* (i.e., *supervised learning*)

### 1.2. Applications of Cluster Analysis

- A key intermediate step for other data mining tasks
  - Generating a compact summary of data for classification, pattern discovery, hypothesis generation and testing, etc.
  - Outlier detection: Outliers—those “far away” from any cluster
- Data summarization, compression, and reduction
  - Ex. Image processing: Vector quantization
- Collaborative filtering, recommendation systems, or customer segmentation
  - Find like-minded users or similar products

### 1.4 A Multi-Dimensional Categorization

#### **Technique-Centered**

- Distance-based methods
- Density-based and grid-based methods
- Probabilistic and generative models
- Leveraging dimensionality reduction methods
- High-dimensional clustering
- Scalable techniques for cluster analysis

#### **Data Type-Centered**

- Clustering numerical data, categorical data, text data, multimedia data, time-series data, sequences, stream data, networked data, uncertain data

### 1.5 An Overview of Typical Clustering Methodologies

**❑ Distance-based methods**

- ❑ Partitioning algorithms: K-Means, K-Medians, K-Medoids
- ❑ Hierarchical algorithms: Agglomerative vs. divisive methods

**❑ Density-based and grid-based methods**

- ❑ Density-based: Data space is explored at a high-level of granularity and then post-processing to put together dense regions into an arbitrary shape
- ❑ Grid-based: Individual regions of the data space are formed into a grid-like structure

**❑ Probabilistic and generative models:** Modeling data from a generative process

- ❑ Assume a specific form of the generative model (e.g., mixture of Gaussians)
- ❑ Model parameters are estimated with the Expectation-Maximization (EM) algorithm (using the available dataset, for a maximum likelihood fit)
- ❑ Then estimate the generative probability of the underlying data points

**❑ High-dimensional clustering**

- ❑ Subspace clustering: Find cluster on various subspaces
  - ❑ Bottom-up, top-down, correlation-based methods vs.  $\delta$ -cluster methods
- ❑ Dimensionality reduction: A vertical form (i.e., columns) of clustering
  - ❑ Columns are clustered; may cluster rows and columns together (co-clustering)
- ❑ Probabilistic latent semantic indexing (PLSI) then LDA: Topic modeling of text data
  - ❑ A cluster (i.e., topic) is associated with a set of words (i.e., dimensions) and a set of documents (i.e., rows) simultaneously
- ❑ Nonnegative matrix factorization (NMF) (as one kind of co-clustering)
  - ❑ A nonnegative matrix  $A$  (e.g., word frequencies in documents) can be approximately factorized two non-negative low rank matrices  $U$  and  $V$

## 1.6 An Overview of Clustering Different Types of Data

**❑ Numerical data**

- ❑ Most earliest clustering algorithms were designed for numerical data

**❑ Categorical data (including binary data)**

- ❑ Discrete data, no natural order (e.g., sex, race, zip-code, and market-basket)

**❑ Text data:** Popular in social media, Web, and social networks

- ❑ Features: High-dimensional, sparse, value corresponding to word frequencies
- ❑ Methods: Combination of k-means and agglomerative; topic modeling; co-clustering

**❑ Multimedia data:** Image, audio, video (e.g., on Flickr, YouTube)

- ❑ Multi-modal (often combined with text data)
- ❑ Contextual: Containing both behavioral and contextual attributes
  - ❑ Images: Position of a pixel represents its context, value represents its behavior
  - ❑ Video and music data: Temporal ordering of records represents its meaning

- ❑ **Time-series data:** Sensor data, stock markets, temporal tracking, forecasting, etc.
  - ❑ Data are temporally dependent
  - ❑ Time: contextual attribute; data value: behavioral attribute
  - ❑ Correlation-based online analysis (e.g., online clustering of stock to find stock tickers)
  - ❑ Shape-based offline analysis (e.g., cluster ECG based on overall shapes)
- ❑ **Sequence data:** Weblogs, biological sequences, system command sequences
  - ❑ Contextual attribute: Placement (rather than time)
  - ❑ Similarity functions: Hamming distance, edit distance, longest common subsequence
  - ❑ Sequence clustering: Suffix tree; generative model (e.g., Hidden Markov Model)

## 2.1 Basic Concepts: Measuring Similarity between Objects

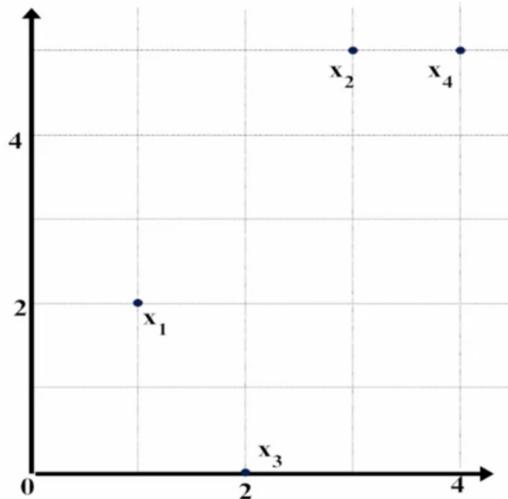
- ❑ A good clustering method will produce high quality clusters which should have
  - ❑ **High intra-class similarity:** **Cohesive** within clusters
  - ❑ **Low inter-class similarity:** **Distinctive** between clusters
- ❑ **Quality function**
  - ❑ There is usually a separate “quality” function that measures the “goodness” of a cluster
  - ❑ It is hard to define “similar enough” or “good enough”
    - ❑ The answer is typically highly subjective
- ❑ There exist many similarity measures and/or functions for different applications

## 2.2 Distance on Numeric Data Minkowski Distance

- ❑ **Data matrix**
  - ❑ A data matrix of  $n$  data points with  $l$  dimensions 
- ❑ **Dissimilarity (distance) matrix**
  - ❑  $n$  data points, but registers only the distance  $d(i, j)$  (typically metric)
  - ❑ Usually symmetric, thus a triangular matrix 
  - ❑ **Distance functions** are usually different for real, boolean, categorical, ordinal, ratio, and vector variables
  - ❑ Weights can be associated with different variables based on applications and data semantics

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nl} \end{pmatrix}$$

$$\begin{pmatrix} 0 & & & \\ d(2,1) & 0 & & \\ \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & 0 \end{pmatrix}$$



**Data Matrix**

point	attribute1	attribute2
x1	1	2
x2	3	5
x3	2	0
x4	4	5

**Dissimilarity Matrix (by Euclidean Distance)**

	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

- **Minkowski distance:** A popular distance measure

$$d(i, j) = \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{il} - x_{jl}|^p}$$

where  $i = (x_{i1}, x_{i2}, \dots, x_{il})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jl})$  are two  $l$ -dimensional data objects, and  $p$  is the order (the distance so defined is also called  $L_p$  norm)

- Properties

- $d(i, j) > 0$  if  $i \neq j$ , and  $d(i, i) = 0$  (Positivity)
- $d(i, j) = d(j, i)$  (Symmetry)
- $d(i, j) \leq d(i, k) + d(k, j)$  (Triangle Inequality)

- $p = 1$ : ( $L_1$  norm) **Manhattan (or city block) distance**

- E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{il} - x_{jl}|$$

- $p = 2$ : ( $L_2$  norm) **Euclidean distance**

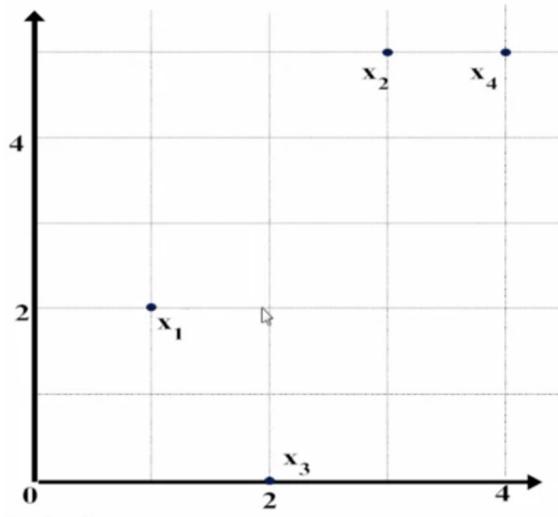
$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{il} - x_{jl}|^2}$$

- $p \rightarrow \infty$ : ( $L_{\max}$  norm,  $L_{\infty}$  norm) **"supremum" distance**

- The maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{p \rightarrow \infty} \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{il} - x_{jl}|^p} = \max_{f=1}^l |x_{if} - x_{jf}|$$

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5



### Manhattan ( $L_1$ )

L	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

### Euclidean ( $L_2$ )

L2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

### Supremum ( $L_\infty$ )

$L_\infty$	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

## 2.3 Proximity Measure for Symmetric vs Asymmetric Binary Variables

- A contingency table for binary data

		Object $j$		sum
		1	0	
Object $i$	1	$q$	$r$	$q + r$
	0	$s$	$t$	$s + t$
sum		$q + s$	$r + t$	$p$

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

- Distance measure for symmetric binary variables:

$$d(i, j) = \frac{r + s}{q + r + s}$$

- Jaccard coefficient (*similarity* measure for *asymmetric*

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

Example (usually, the asymmetric attributes are more important)

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

		Mary		
		1	0	$\Sigma_{\text{row}}$
Jack		1	2	2
		0	1	4
$\Sigma_{\text{col}}$		3	3	6

- ❑ Gender is a symmetric attribute (not counted in)

- ❑ The remaining attributes are asymmetric binary

- ❑ Let the values Y and P be 1, and the value N be 0

❑ Distance:  $d(i, j) = \frac{r + s}{q + r + s}$

$$d(jack, mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(jack, jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(jim, mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

		Mary		
		1	0	$\Sigma_{\text{row}}$
Jack		1	1	2
		0	1	4
$\Sigma_{\text{col}}$		2	4	6

## 2.4 Distance between Categorical Attributes Ordinal Attributes and Mixed Types

- ❑ Categorical data, also called nominal attributes

- ❑ Example: Color (red, yellow, blue, green), profession, etc.

- ❑ Method 1: Simple matching

- ❑  $m$ : # of matches,  $p$ : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- ❑ Method 2: Use a large number of binary attributes

- ❑ Creating a new binary attribute for each of the  $M$  nominal states

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)
- Can be treated like interval-scaled
  - Replace *an ordinal variable value* by its rank:  $r_{if} \in \{1, \dots, M_f\}$
  - Map the range of each variable onto [0, 1] by replacing *i*-th object in the *f*-th variable by
 
$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
  - Example: freshman: 0; sophomore: 1/3; junior: 2/3; senior 1
    - Then distance:  $d(\text{freshman, senior}) = 1$ ,  $d(\text{junior, senior}) = 1/3$
  - Compute the dissimilarity using methods for interval-scaled variables

## 2.5 Proximity Measure between Two Vectors Cosine Similarity

### Example: Calculating Cosine Similarity

- Calculating Cosine Similarity:
 
$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$$

where  $\bullet$  indicates vector dot product,  $\|d\|$ : the length of vector  $d$
- Ex: Find the **similarity** between documents 1 and 2.
- $d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$        $d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$
- First, calculate vector dot product  

$$d_1 \bullet d_2 = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 1 + 2 \times 1 + 0 \times 0 + 0 \times 1 = 25$$
- Then, calculate  $\|d_1\|$  and  $\|d_2\|$   

$$\|d_1\| = \sqrt{5 \times 5 + 0 \times 0 + 3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0} = 6.481$$

$$\|d_2\| = \sqrt{3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 1 \times 1 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1} = 4.12$$
- Calculate cosine similarity:  $\cos(d_1, d_2) = 25 / (6.481 \times 4.12) = 0.94$

# Covariance Matrix

- The variance and covariance information for the two variables  $X_1$  and  $X_2$  can be summarized as  $2 \times 2$  covariance matrix as

$$\begin{aligned}\Sigma &= E[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^T] = E[\begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \end{pmatrix} (X_1 - \mu_1 \quad X_2 - \mu_2)] \\ &= \begin{pmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] \end{pmatrix} \\ &= \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}\end{aligned}$$

- Generalizing it to  $d$  dimensions, we have,

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{d1} & x_{d2} & \dots & x_{dd} \end{pmatrix} \quad \Sigma = E[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^T] = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_d^2 \end{pmatrix}$$

## Week1-Quiz

2021年1月21日 11:12

1. Answer : 4.6

1. The following real dataset contains information about two different flowers: Iris setosa and Iris versicolor.

Species	Sepal length	Sepal width	Petal length	Petal width
Iris setosa	4.9	3.0	1.4	0.2
Iris versicolor	5.6	2.5	3.9	1.1

**What is the Manhattan distance between these two objects?**

2. Answer: 1/4

2. The following real dataset contains two samples from the dataset for Prediction of Molecular Bioactivity for Drug Design - Binding to Thrombin, with sampled features. For each activity (F1, F2, ..., F10), the class value (0/1) indicates if the activity is active or inactive.

Cases	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10
1	0	1	1	0	0	0	1	1	1	1
2	0	1	0	0	1	0	1	0	0	1

**Assume all the activities are *asymmetric* binary variables. What is the Jaccard coefficient between case 1 and case 2?**

3. Answer: 8/21

3. The following real world dataset contains two samples from Car Evaluation Database, which was derived from a simple hierarchical decision model originally developed for the demonstration of DEX ( Bohanec, M., & Rajkovic, V. (1990). Expert system for decision making. *Sistemica* 1(1), 145-157.). The model evaluates cars according to the following concept structure:

1 / 1 point

CAR	car acceptability
.PRICE	overall price
.. buying	buying price
.. maint	price of the maintenance
. TECH	technical characteristics
.. COMFORT	comfort
... doors	number of doors
... persons	capacity in terms of persons to carry
... lug_boot	the size of luggage boot
... safety	estimated

The attribute values are as follows:

Attribute	Values (categorical)
buying	v-high, high, med, low
maint	v-high, high, med, low
doors	2, 3, 4, 5 - more
persons	2, 4, more
lug_boot	small, med, big
safety	low, med, high

Case	buying	maint	doors	persons	lug_boot	safety
Car 1	med	v-high	3	more	small	med
Car 2	high	v-high	4	4	big	med

To calculate the distance between objects with categorical attributes, we use a set of binary attributes to represent each categorical attribute. Assume all the binary attributes are **symmetric**. What is the distance between Car 1 and Car 2?

5. Answer : 0.4

5. Given the following two short texts with punctuation removed, calculate the cosine similarity between them based on the **1 / 1 point** bag of words model.

Text1: all grown-ups were once children but only few of them remember it

Text2: all children should be very understanding of grown-ups

6. Answer 0.882

6. With regard to the species of Iris versicolor, we have sampled data on the features of sepal length and sepal width, as follows. **1 / 1 point**

Feature	Sepal length	Petal length
Case 1	7.0	3.2
Case 2	6.4	3.2
Case 3	6.9	3.1
Case 4	5.5	2.3
Case 5	6.5	2.8

What is the correlation coefficient between sepal length and sepal width?