

## Cluster Analysis

The process of grouping a set of physical or abstract objects into classes of similar objects is called clustering.

A cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters.

### Applications

- market research → clustering can help marketers discover distinct groups in their customers based on their purchasing patterns.
- pattern recognition
- data analysis
- image processing
- in biology → categorize genes with similar functionality

clustering is also called data segmentation. because clustering partitions large data sets into groups according to their similarity.

→ Clustering can also be used for outlier detection. Applications of outlier detection include the detection of credit card fraud and monitoring of criminal activities in e-commerce.

For example, exceptional cases in credit card transactions, such as very expensive and frequent purchases, may be of interest as possible fraudulent activity.

The following are typical requirements of clustering in data mining:

- (1) Scalability → Many clustering algorithms work well on small data sets, ~~however~~, clustering on a sample of a given large data set may lead to biased results. Highly scalable clustering algorithms are needed.
- (2) Ability to deal with different types of attributes → Many algorithms are designed to cluster interval based (numerical) data. However, applications may require clustering other types of data, such as binary, categorical (nominal) and ordinal data.
- (3) Discovery of clusters with arbitrary shape. Many clustering algorithms determine clusters based on Euclidean distance. Algorithms based on such distance find spherical clusters with similarity, size and density. It is important to develop algorithms that can detect clusters of arbitrary shape.
- (4) Minimal requirements for domain knowledge to determine input parameters. → Many clustering algorithms require users to input certain parameters in cluster analysis (such as the no. of desired clusters)

The clustering results can be quite sensitive to input parameters. Parameters are often difficult to determine for high-dimensional objects.

### (5) Ability to deal with noisy data

Most real world databases contain outliers or missing, unknown or erroneous data. Some clustering algorithms are sensitive to such data and may lead to clusters of poor quality.

### (6) Incremental clustering and insensitivity to the order of input records :-

Some clustering algorithms can not incorporate newly inserted data onto existing clustering structures and determine a new clustering from scratch.

Some clustering algorithms are sensitive to the order of input data.

It is important to develop incremental clustering algorithms that are insensitive to the order of input.

### (7) High dimensionality :-

A database or a data warehouse can contain several dimensions or attributes.

Many clustering algorithms are good at handling two or three dimensions.

Finding clusters of data objects in high dimensional space is challenging.

### (8) constraint based clustering :-

Real world applications may need to perform clustering under various kinds of constraints.

A challenging task is to find groups of data with good clustering behavior that satisfy specified constraints.

Interpretability and usability.

Users expect clustering results to be interpretable, comprehensible and usable.

## II Types of data in cluster analysis

Data matrix  $\rightarrow$  (object-by-variable structure)

This represents  $n$  objects, with  $p$  variables.

The structure is  $\rightarrow n \times p$  matrix. Also called two-mode matrix

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

Dissimilarity matrix (object by object structure)

Represented by  $n \times n$  table.

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

where  $d(i,j)$  is the measured difference between objects  $i$  and  $j$ . Also called one-mode matrix.

$$d(i,j) = d(j,i) \text{ and } d(i,i) = 0.$$

(3)

Interval Scaled Variables  $\rightarrow$  continuous measurements of a roughly linear scale.

Ex - weight and height, weather temp. etc.

The measurement unit used can affect the clustering analysis. For example, changing

measurements units from meters to inches for height or from kg to pounds for weight, may lead to different clustering structures. so we have to standardize the data.

(1) calculate the mean absolute deviation.

$$S_f = \frac{1}{n} [ |x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f| ]$$

$m_f \rightarrow$  mean value of f. More robust to outliers.

(2) calculate the Z-score.

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

(3) Euclidean distance.

$$d(c_i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{in} - x_{jn})^2}$$

Manhattan or city block distance.

$$d(c_i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{in} - x_{jn}|$$

Both the distances satisfy the following.

1.  $d(c_i, j) \geq 0$  = distance is non-negative no.

2.  $d(c_i, i) = 0$ , the distance of an object to itself.

3.  $d(c_i, j) = d(j, i)$ , distance is symmetric.

4.  $d(c_i, j) \leq d(i, k) + d(k, j)$

Minkowski distance.

$$d(c_i, j) = \left[ (x_{i1} - x_{j1})^p + (x_{i2} - x_{j2})^p + \dots + (x_{in} - x_{jn})^p \right]^{\frac{1}{p}}$$

p = the integer.

Binary variables  $\rightarrow$  it has only two states: 0 or 1  
 How can we compute the dissimilarity between two binary variables?

Ans  $\rightarrow$  By computing a dissimilarity matrix from the given binary data.

		object j	
		1	0
object i	1	$q$	$r$
	0	$s$	$t$
	sum	$q+s$	$r+t$

where  $q =$  no. of variables that equal 1 for both i & j objects.  
 $r =$  " " " that equal 1 for i and 0 for j objects.  
 $s =$  " " " that equal 0 for i and 1 for j =  
 $t =$  " " " that equal 0 for i and 0 for j objects.  
 $p =$  total no. of variables =  $q+r+s+t$ .

Binary variable is symmetric  $\rightarrow$  if both of its states are equally valuable and carry the same weight. Ex: gender - male or female. symmetric binary dissimilarity (or distance) can be assess between objects i and j by

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

A binary variable is asymmetric; if the outcomes of the states are not equally important. Ex  $\rightarrow$  +ve or -ve outcomes of a disease test. The asymmetric binary dissimilarity, where the no. of -ve matches, t is considered unimportant and thus is ignored is calculated as

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

Eg suppose that a patient record table containing the attributes name, gender, fever, cough, test-1, test-2, test-3 and test-4, where name is an object identifier, gender is a symmetric attribute and remaining attributes are asymmetric binary.

Let the values yes (y) and positive (p+) set to 1.

and the value N (no or -ve) set to 0.  
suppose that the distance between objects is computed based only on the asymmetric variables. The distance between each pair of these patients is.

$$d(\text{Jack, Mary}) = \frac{r+s}{q+r+s} = \frac{0+1}{2+0+1} = 0.33$$

$$d(\text{Jack, Jim}) = \frac{1+t}{t+t+1} = 0.67.$$

$$d(\text{Mary, Jim}) = \frac{1+2}{1+1+2} = 0.75$$

name gender fever cough +1 +2 +3 +4

			y(1)	N(0)	p(1)	N(0)	N(0)	N(0)
(i) Jack	M		y(1)	N(0)	p(1)	N(0)	N(0)	N(0)
(ii) Mary	F		y(1)	N(0)	p(1)	N(0)	p(1)	N(0)
Jim	M		y(1)	y(1)	N(0)	N(0)	N(0)	N(0)

$q = \text{no. of variables that equal 1 for both } i, j$

$r = \text{no. of variables that equal 1 for } i \text{ and 0 for } j$

$s = \text{no. of variables that equal 0 for } i \text{ and 1 for } j$

## Categorical variables

A categorical variable is a generalization of the binary variable that can take more than two states.

Ex map\_color is a categorical variable

That have, five states, red, yellow, green, pink and blue.

Let the no. of states of a categorical variable be M. The states can be denoted by letters, symbols <sup>OK</sup>, <sub>subset of integers.</sub> The dissimilarity between two objects i and j can be computed based on the ratio of mismatched -

$$d(i, j) = \frac{P-m}{P} \quad \text{where } m = \text{no. of matches} \\ (\text{i.e. the no. of variables for which } i \text{ and } j \text{ are in the same state}) \\ P = \text{total no. of variables.}$$

<u>Ex</u>	Object identifier	<u>test</u>
1	code-A	
2	code-B	
3	code-C	
4	code-A	

Dissimilarity matrix =

$$= \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & & & \\ d(2,1) & 0 & & \\ d(3,1) & d(3,2) & 0 & \\ d(4,1) & d(4,2) & d(4,3) & 0 \end{bmatrix}$$

Here, we have one categorical variable, test-1, so,  $P=1$ .

$d(i, j) = 0$  if objects i and j match and 1 if the objects differ.

(1)

## Ordinal variables

A discrete ordinal variable, resembles a categorical variable, except that the states of the ordinal value are ordered in a meaningful sequence.

Ex Assistant, associate and professor.

## continuous ordinal variables

A set of continuous data for an unknown scale, i.e. the relative ordering of the values is essential but their actual magnitude is not.

Ex relative ranking in sports (e.g. Gold, Silver, bronze)

~~ordinal variables may~~  
Suppose,  $f$  is a variable from a set of ordinal variables describing  $n$  objects.

The discriminant computation involves the following:

1. The value of  $f$  for the  $i$ th object is defined if has  $M_f$  ordered states, representing the ranking  $1, \dots, M_f$ . Replace each  $x_{if}$  by its corresponding rank,  $r_{if} \in \{1, \dots, M_f\}$
2. since each ordinal variable can have a different no. of states, it is often necessary to map the range of each variable onto  $\{0, 1\}$ . This can be done by replacing the rank  $r_{if}$  of the  $i$ th object in the  $f$ th variable by

$$Z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

3. Dissimilarity can then be computed using any of the distance measures, using  $Z_{if}$  to represent the  $f$  value for the  $i$ th object

<u>Ex</u>	Object identifier	test-2	Rank	Normalized value
	1	Excellent	3	1
	2	Fair	1	0
	3	Good	2	0.5
	4	Excellent	3	1

$M_f = 3$ , three states, Ex, Fair, Good.

① Replace each value for test-2 by its rank.

$$\text{Ex} = 3, \text{Fair} = 1, \text{Good} = 2$$

② Normalize the ranking by mapping rank 1 to 0.0, rank 2 to 0.5 and rank 3 to 1.0.

③ Use Euclidean distance, which results in the following dissimilarity matrix.

$$\begin{bmatrix} 0 & & & \\ d(2,1) & 0 & & \\ d(3,1) & d(3,2) & 0 & \\ d(4,1) & d(4,2) & d(4,3) & 0 \end{bmatrix}$$

$$\sqrt{(1-0)^2 + (0-0)^2 + (0-1)^2} = \sqrt{2}$$

$$= \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1 & 0.5 & 0 \end{bmatrix}$$

## (2)

### Ratio scaled variables

A ratio scaled variable makes a +ve measurement on a non-linear scale, such as an exponential scale.

$Ae^{Bt}$  or  $Ae^{-Bt}$ , where A and B are +ve constants, t is time.

Ex → growth of a bacteria population.

There are three methods to handle ratio-scaled variables:-

(1) Treat ratio-scaled variables like interval scaled variables. But it's not a good choice.

(2) Apply logarithmic transformation:-

A ratio scaled variable of having value  $x_{it}$  for object i, ~~is~~ if  $y_{it} = \log(x_{it})$ . Then treat them as interval-valued.

(3) Treat  $x_{it}$  as continuous ordinal data and treat their ranks as interval-valued.

<u>Ex</u>	Object	$t_{it}^{B+3}$	$\log(t_{it}^{B+3})$	Then use
		445	2.65	the Euclidean
1	22	1.34		distance or
2	164	2.21		the transforming
3	120	3.08		value, we get
4				

The dissimilarity matrix =

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & & 0 & & \\ d(4,1) & & & 0 & \\ & d(4,2) & & d(4,3) & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & & & & \\ 1.31 & 0 & & & \\ 0.44 & 0.87 & 0 & & \\ 0.43 & 1.74 & 0.87 & 0 & \end{bmatrix}$$

1  
2  
3  
4

## variables of mixed types

Suppose that the data set contains variables of mixed type. The dissimilarity  $d(i, j)$  between objects  $i$  and  $j$  is defined as

$$d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}} \quad (1)$$

where the indicator  $\delta_{ij}^{(f)} = 1$  if either  $x_{if}$  or  $x_{jf}$  is missing - (OR)

(2)  $x_{if} = x_{jf} = 0$  and variable  $f$  is asymmetric binary  
~~so~~  $\delta_{ij}^{(f)} = 1$  otherwise  $\delta_{ij}^{(f)} = 0$

The dissimilarity between  $i$  and  $j$  is computed as follows -

(1) If  $f$  is interval-scaled:  $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{\max_{h \in H} x_{hf} - \min_{h \in H} x_{hf}}$  where  $h$  runs over all non-missing objects for variable  $f$ .

(2) If  $f$  is binary or categorical:  $d_{ij}^{(f)} = 1 - \delta_{ij}^{(f)}$ , if  $x_{if} = x_{jf}$ , otherwise  $d_{ij}^{(f)} = 1$

(3) If  $f$  is ordinal, compute the ranks  $r_{if}$  and  $r_{jf}$  and treat  $r_{if}$  as interval-scaled:  $r_{if} = \frac{r_{if}-1}{M_f-1}$

(4) If  $f$  is ratio-scaled: either perform logarithmic transformation and treat the transformed data as interval-scaled.

(OR) treat  $f$  as continuous ordinal data, compute  $r_{if}$  and  $r_{jf}$ , and then treat  $r_{if}$  as interval-scaled.

Ex	object	$r_{if-1}$ (categorical)	$r_{if-2}$ (ordinal)	$r_{if-3}$ (ratio-scaled)
	1	code-A	excellent	445
	2	code-B	fair	28
	3	code-C	good	164
	4	code-A	excellent	1210

(3)

The procedure we followed for test-1 and test-2 are the same.

Let us consider for test-3. After applying logarithmic transformation we get.

$$2.65 \quad \text{Let } \max_n x_n = 3.08$$

$$1.34 \quad \min_n x_n = 1.34$$

$$2.21$$

$$3.08$$

Then divide the values in the dissimilarity matrix by  $3.08 - 1.34 = 1.74$ .

$$\begin{bmatrix} 0 \\ 1.34 \\ 0.49/1.74 & 0 \\ 0.43/1.74 & 1.74/1.74 & 0 \\ 0.43/1.74 & 1.74/1.74 & 0.87/1.74 & 0 \end{bmatrix}$$

$$d_3 = \begin{bmatrix} 0 \\ 0.25 \\ 0.25 \\ 0.25 & 0 \\ 0.25 & 0.50 & 0 \\ 0.25 & 1 & 0.5 & 0 \end{bmatrix}$$

Now are the dissimilarity matrices for the 3-varieties using eqn C(1)

$$d_1 = \begin{bmatrix} 0 \\ 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$$\text{Now. } d_{(2,1)} = \frac{1(C1) + 1(C1) + 1(0.25)}{3} = 0.92$$

$$d_{(3,1)} = \frac{1(C1) + 1(0.5) + 1(0.25)}{3} = 0.58$$

$$d_2 = \begin{bmatrix} 0 \\ 1 & 0 \\ 0.5 & 0.5 & 0 \\ 0 & 1.0 & 0.50 \end{bmatrix}$$

The final result is

$$\begin{bmatrix} 0 \\ 0.92 & 0 \\ 0.58 & 0.67 & 0 \\ 0.09 & 1 & 0.67 \end{bmatrix}$$

3] 1.25 1.50  
1.3 8.5  
8.5 21

## vector objects

To measure distance between complex objects.

similarity function,  $S(x, y)$ , to compare two vectors  $x$  and  $y$ . calculate cosine measure

$$S(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$x^t$  = transposition of vector  $x$

$\|x\|$  = Euclidean norm of vector  $x$

$\|y\|$  = " " " of "  $y$

$s$  = cosine angle between vectors  $x$  and  $y$ .

Euclidean norm of vector  $x$  =  $\sqrt{x_1^2 + x_2^2 + \dots + x_p^2}$

Ex Given  $x = (1, 1, 0, 0)$   
" " " "  
 $y = (0, 1, 1, 0)$

$$S(x, y) = \frac{0 + 1 + 0 + 0}{\sqrt{2} \sqrt{2}} = 0.5$$

Tanimoto coefficient =  $S(x, y) = \frac{x^t \cdot y}{x^t \cdot x + y^t \cdot y - x^t \cdot y}$

## Types of clustering Method

partitioning methods - K-means, K-medoids

Hierarchical methods  $\rightarrow$  agglomerative or divisive

Density based methods

Grid based methods

Model based methods