

Clustering Analysis

Day 1 : February 25, 2021

What is clustering?

What is clustering?

- The organization of unlabeled data into similarity groups called clusters.
- A cluster is a collection of data items which are "similar" between them, and "dissimilar" to data items in other clusters.



History of clustering

- John Snow, a London physician plotted the location of cholera deaths on a map during an outbreak in the 1850s.
- The locations indicated that cases were clustered around certain intersections where there were polluted wells -- thus exposing both the problem and the solution.





From: Nina Mishra HP Labs

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)
- Typical ways to use/apply cluster analysis
 - As a stand-alone tool to get insight into data distribution, or
 - As a preprocessing (or intermediate) step for other algorithms

Reminder!



Elements of clustering



For example, by optimizing the criterion function

Similarity



Cluster evaluation

- Intra-cluster cohesion (compactness):
 - Cohesion measures how near the data points in a cluster are to the cluster centroid.
 - Sum of squared error (SSE) is a commonly used measure.
- Inter-cluster separation (isolation):
 - Separation means that different cluster centroids should be far away from one another.
- In most applications, expert judgments are still the key

Number of clustering



- Possible approaches
 - 1. fix the number of clusters to k
 - find the best clustering according to the criterion function (number of clusters may vary)

Cluster analysis applications

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
 - Dynamic trend detection
 - Clustering stream data and detecting trends and patterns
 - Multimedia data analysis, biological data analysis and social network analysis
 - Ex. Clustering images or video/audio clips, gene/protein sequences, etc.

Elements and Types of Clustering

Requirements and Challenges

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. non-exclusive (e.g., one document may belong to more than one class)

Similarity measure

Distance-based (e.g., Euclidean, road network, vector) vs. connectivitybased (e.g., density or contiguity)

Clustering space

Full space (often when low dimensional) vs. subspaces (often in highdimensional clustering)

Requirements and Challenges

Quality

- Ability to deal with different types of attributes: Numerical, categorical, text, multimedia, networks, and mixture of multiple types
- Discovery of clusters with arbitrary shape
- Ability to deal with noisy data

Scalability

- Clustering all the data instead of only on samples
- High dimensionality
- Incremental or stream clustering and insensitivity to input order

Constraint-based clustering

- User-given preferences or constraints; domain knowledge; user queries
- Interpretability and usability



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, timeseries data, sequences, stream data, networked data, uncertain data

Additional Insight-Centered

Visual insights, semi-supervised, ensemble-based, validation-based

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

Typical Clustering Methodologies

High Dimensional Clustering

Subspace clustering: clustering on various subspaces
 Bottom-up, top-down, correlation-based methods vs. δ-cluster method

Dimensionality reduction: a vertical form of clustering

Probabilistic latent semantic indexing (PLSI) then LDA: topic modeling of text data

□ Non-negative matrix factorization (NMF) clustering (an example of co-clustering)

Spectral clustering: use the spectrum of the similarity matrix of the data to perform dimensionality reduction for clustering in fewer dimensions

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

Clustering Different Types of Data

Time series data: Sensor data, stock markets, temporal tracking, forecasting, etc.

Sequence data: Weblogs, biological sequences, system command sequences

Stream data: Real-time data

User Insights and Clustering

Visual insights: One picture is worth a thousand words

- Human eyes: High-speed processor linking with a rich knowledge-base
- A human can provide intuitive insights; HD-eye: visualizing HD clusters

Semi-supervised insights: Passing user's insights or intention to system

- User-seeding: A user provides a number of labeled examples, approximately representing categories of interest
- Multi-view and ensemble-based insights
 - Multi-view clustering: Multiple clusterings represent different perspectives
 - Multiple clustering results can be ensembled to provide a more robust solution

Validation-based insights: Evaluation of the quality of clusters generated

May use case studies, specific measures, or pre-existing labels

Clustering techniques



Recommended readings

Major Reference Books on Cluster Analysis

- Jiawei Han, Micheline Kamber, and Jian Pei. Data Mining: Concepts and Techniques. Morgan Kaufmann, 3rd ed., 2011 (Chapters 10 & 11)
- Charu Aggarwal and Chandran K. Reddy (eds.). Data Clustering: Algorithms and Applications. CRC Press, 2014
- Mohammed J. Zaki and Wagner Meira, Jr.. Data Mining and Analysis: Fundamental Concepts and Algorithms. Cambridge University Press, 2014
- Reference paper for this lecture
 - Charu Aggarwal. An Introduction to Clustering Analysis. in Aggarwal and Reddy (eds.). Data Clustering: Algorithms and Applications (Chapter 1). CRC Press, 2014

Similarity Measures for Clustering

Good clustering

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
- Similarity measure is critical for cluster analysis

Similarity, Dissimilarity, and Proximity

Similarity measure or similarity function

- A real-valued function that quantifies the similarity between two objects
- Measure how two data objects are alike: The higher value, the more alike
- Often falls in the range [0,1]: 0: no similarity; 1: completely similar

Dissimilarity (or distance) measure

- Numerical measure of how different two data objects are
- In some sense, the inverse of similarity: The lower, the more alike
- Minimum dissimilarity is often 0 (i.e., completely similar)
- □ Range [0, 1] or $[0, \infty)$, depending on the definition
- Proximity usually refers to either similarity or dissimilarity

Data and Dissimilarity Matrices

🗖 Data matrix

- A data matrix of n data points with I 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

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nl} \end{pmatrix}$$
$$\begin{pmatrix} 0 & & & \\ d(2,1) & 0 & & \\ & \ddots & \ddots & \ddots & \ddots \end{pmatrix}$$

 $(\mathbf{r} + \mathbf{r})$

 $: : \cdot \\ d(n,1) d(n,2) \dots 0$

Example



Data Matrix

point	attribute1	attribute2
x1	1	2
x2	3	5
x3	2	0
<i>x4</i>	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

Numeric Data: Minkowski Distance

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 $i = (x_{i1}, x_{i2}, ..., x_{il})$ and $j = (x_{j1}, x_{j2}, ..., x_{jl})$ 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)

□ $d(i, j) \le d(i, k) + d(k, j)$ (Triangle Inequality)

- A distance that satisfies these properties is a metric
- Note: There are nonmetric dissimilarities, e.g., set differences

Minkowski Distance

- $\square p = 1$: (L₁ norm) Manhattan (or city block) distance
 - □ E.g., the Hamming distance: the number of bits that are different between two binary vectors $d(i, j) = |x_{i1} - x_{i1}| + |x_{i2} - x_{i2}| + \dots + |x_{il} - x_{il}|$

 \square *p* = 2: (L₂ norm) Euclidean distance

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

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

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

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

Example



Man	hattan	(L.)
intern	naccan	(-1/

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

Euclidean (L₂)

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_{∞})

L∞	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

Proximity Measure for Binary Attributes

A contingency table for binary data

		Ob	ject j	
-		1	0	sum
Object i	1	q	r	q + r
Object /	0	8	t	s+t
	sum	q + s	r+t	p

Distance measure for symmetric binary variables:

Distance measure for asymmetric binary variables:

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

- □ Jaccard coefficient (*similarity* measure for *asymmetric* binary variables): $sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$
- □ Note: Jaccard coefficient is the same as "coherence": (a concept discussed in Pattern Discovery)

$$coherence(i, j) = \frac{sup(i, j)}{sup(i) + sup(j) - sup(i, j)} = \frac{q}{(q+r) + (q+s) - q}$$

Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	М	Y	N	Р	N	N	N
Mary	F	Y	N	Р	N	P	N
Jim	Μ	Y	Р	N	Ν	N	N

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

 $d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$
 $d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$
 $d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$



1

0

Jim

Proximity Measure for Categorical Attributes

Categorical data, also called nominal attributes

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

Method 1: Simple matching

D *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

Ordinal Variables

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: $r_{if} \in \{1, ..., M_f\}$
- 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(freshman, senior) = 1, d(junior, senior) = 1/3

Compute the dissimilarity using methods for interval-scaled variables

Attributes of Mixed Type

A dataset may contain all attribute types

Nominal, symmetric binary, asymmetric binary, numeric, and ordinal

One may use a weighted formula to combine their effects:

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

1

If f is numeric: Use the normalized distance

□ If f is binary or nominal: $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$; or $d_{ij}^{(f)} = 1$ otherwise

If f is ordinal

Compute ranks
$$z_{if}$$
 (where $z_{if} = \frac{r_{if} - 1}{M_f - 1}$)

Treat z_{if} as interval-scaled

Cosine Similarity of Two Vectors

A document can be represented by a bag of terms or a long vector, with each attribute recording the *frequency* of a particular term (such as word, keyword, or phrase) in the document

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

Other vector objects: Gene features in micro-arrays

Applications: Information retrieval, biologic taxonomy, gene feature mapping, etc.

Cosine measure: If d_1 and d_2 are two vectors (e.g., term-frequency vectors), then

$$cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$$

where \bullet indicates vector dot product, ||d||: the length of vector d

Recommended Readings

L. Kaufman and P. J. Rousseeuw, Finding Groups in Data: An Introduction to Cluster Analysis, John Wiley & Sons, 1990

Mohammed J. Zaki and Wagner Meira, Jr. Data Mining and Analysis: Fundamental Concepts and Algorithms. Cambridge University Press, 2014

Jiawei Han, Micheline Kamber, and Jian Pei. Data Mining: Concepts and Techniques. Morgan Kaufmann, 3rd ed., 2011

Charu Aggarwal and Chandran K. Reddy (eds.). Data Clustering: Algorithms and Applications. CRC Press, 2014

Partitioning-Based Clustering Methods

Partitioning Algorithms: Basic Concepts

- Partitioning method: Discovering the groupings in the data by optimizing a specific objective function and iteratively improving the quality of partitions
- K-partitioning method: Partitioning a dataset D of n objects into a set of K clusters so that an objective function is optimized (e.g., the sum of squared distances is minimized, where c_k is the centroid or medoid of cluster C_k)
 - A typical objective function: Sum of Squared Errors (SSE)

$$SSE(C) = \sum_{k=1}^{\infty} \sum_{x_{i \in C_k}} ||x_i - c_k||^2$$

- Problem definition: Given K, find a partition of K clusters that optimizes the chosen partitioning criterion
 - Global optimal: Needs to exhaustively enumerate all partitions
 - Heuristic methods (i.e., greedy algorithms): K-Means, K-Medians, K-Medoids, etc.

K-Means Clustering

- K-Means (MacQueen'67, Lloyd'57/'82)
 - Each cluster is represented by the center of the cluster
- Given K, the number of clusters, the K-Means clustering algorithm is outlined as follows
 - Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Re-compute the centroids (i.e., mean point) of each cluster
- Until convergence criterion is satisfied
- Different kinds of measures can be used
 - Manhattan distance (L₁ norm), Euclidean distance (L₂ norm), Cosine similarity

Example



K-Means Discussion

- Efficiency: O(tKn) where n: # of objects, K: # of clusters, and t: # of iterations
 - Normally, K, t << n; thus, an efficient method</p>
- K-means clustering often terminates at a local optimal
 - Initialization can be important to find high-quality clusters
- Need to specify K, the number of clusters, in advance
 - There are ways to automatically determine the "best" K
 - In practice, one often runs a range of values and selected the "best" K value
- Sensitive to noisy data and outliers
 - Variations: Using K-medians, K-medoids, etc.
- K-means is applicable only to objects in a continuous n-dimensional space
 - Using the K-modes for categorical data
- Not suitable to discover clusters with non-convex shapes
- □ Using density-based clustering, kernel K-means, etc.

Variations of K-Means Clustering

There are many variants of the K-Means method, varying in different aspects



Initialization of K-Means

- Different initializations may generate rather different clustering results (some could be far from optimal)
- Original proposal (MacQueen'67): Select k seeds randomly
 - Need to run the algorithm multiple times using different seeds
- □ There are many methods proposed for better initialization of *k* seeds
 - **K-Means++** (Arthur & Vassilvitskii'07):
 - The first centroid is selected at random
 - The next centroid selected is the one that is farthest from the currently selected (selection is based on a weighted probability score)
 - The selection continues until k centroids are obtained



Example



000

×

2

*

-1

0

-2

□ This run of k-Means generates a poor quality clustering

Outliers: From K-Means to K-Medoids

- The K-Means algorithm is sensitive to outliers!—since an object with an extremely large value may substantially distort the distribution of the data
- K-Medoids: Instead of taking the mean value of the object in a cluster as a reference point, medoids can be used, which is the most centrally located object in a cluster
- The K-Medoids clustering algorithm:
 - Select K points as the initial representative objects (i.e., as initial K-medoids)

Repeat

- Assigning each point to the cluster with the closest medoid
- Randomly select a non-representative object o_i
- Compute the total cost S of swapping the medoid m with o_i
- □ If S < 0, then swap m with o_i to form the new set of medoids
- Until convergence criterion is satisfied

Partitioning Around Medoids (PAM)



K-Medoids Discussion

K-Medoids Clustering: Find representative objects (medoids) in clusters

PAM (Partitioning Around Medoids: Kaufmann & Rousseeuw 1987)

- Starts from an initial set of medoids, and
- Iteratively replaces one of the medoids by one of the non-medoids if it improves the total sum of the squared errors (SSE) of the resulting clustering
- PAM works effectively for small data sets but does not scale well for large data sets (due to the computational complexity)
- □ Computational complexity: PAM: O(K(n K)²) (quite expensive!)

Efficiency improvements on PAM

CLARA (Kaufmann & Rousseeuw, 1990):

PAM on samples; O(Ks² + K(n – K)), s is the sample size

CLARANS (Ng & Han, 1994): Randomized re-sampling, ensuring efficiency + quality

K-Medians & Outliers

- Medians are less sensitive to outliers than means
 - Think of the median salary vs. mean salary of a large firm when adding a few top executives!
- K-Medians: Instead of taking the mean value of the object in a cluster as a reference point, medians are used (L₁-norm as the distance measure)
- □ The criterion function for the *K*-*Medians* algorithm:
- □ The *K*-*Medians* clustering algorithm:
 - Select K points as the initial representative objects (i.e., as initial K medians)

Repeat

- Assign every point to its nearest median
- Re-compute the median using the median of each individual feature
- Until convergence criterion is satisfied

$$S = \sum_{k=1}^{K} \sum_{x_{i \in C_k}} |x_{ij} - med_{kj}|$$

K-Modes: Clustering Categorical Data

- K-Means cannot handle non-numerical (categorical) data
 - Mapping categorical value to 1/0 cannot generate quality clusters for highdimensional data
- **K-Modes:** An extension to K-Means by replacing means of clusters with modes
- Dissimilarity measure between object X and the center of a cluster Z
 - $\Box \quad \Phi(x_j, z_j) = 1 n_j^r / n_j \text{ when } x_j = z_j \text{ ; } 1 \text{ when } x_j \neq z_j$
 - where z_j is the categorical value of attribute j in Z_j, n_j is the number of objects in cluster l, and n_j^r is the number of objects whose attribute value is r
- This dissimilarity measure (distance function) is frequency-based
- Algorithm is still based on iterative *object cluster assignment* and *centroid update*
- A fuzzy K-Modes method is proposed to calculate a fuzzy cluster membership value for each object to each cluster
- A mixture of categorical and numerical data: Using a *K-Prototype* method

Kernel K-Means Clustering

- Kernel K-Means can be used to detect non-convex clusters
 - Content of the second s
- Idea: Project data onto the high-dimensional kernel space, and then perform K-Means clustering



- Map data points in the input space onto a high-dimensional feature space using the kernel function
- Perform K-Means on the mapped feature space
- Computational complexity is higher than K-Means
 - Need to compute and store n x n kernel matrix generated from the kernel function on the original data
- The widely studied spectral clustering can be considered as a variant of <u>Kernel K-Means</u> clustering

Kernel Functions and Kernel K-Means

Typical kernel functions:

- □ Polynomial kernel of degree h: $K(X_i, X_j) = (X_i \cdot X_j + 1)^h$
- Gaussian radial basis function (RBF) kernel: $K(X_i, X_i) = e^{-||X_i X_j||^2/2\sigma^2}$
- Gigmoid kernel: $K(X_i, X_j) = tanh(\kappa X_i \cdot X_j \delta)$

□ The formula for kernel matrix K for any two points $x_i, x_j \in C_k$ is $K_{x_ix_j} = \phi(x_i) \bullet \phi(x_j)$

The SSE criterion of *kernel K-means*:

 $SSE(C) = \sum_{i=1}^{K} \sum_{i=1}^{K} ||\phi(x_i) - c_k||^2$ The formula for the cluster centroid:

$$c_{k} = \frac{\sum_{x_{i \in C_{k}}} \phi(x_{i})}{\mid C_{k} \mid}$$

 \Box Clustering can be performed without the actual individual projections $\phi(x_i)$ and $\phi(x_i)$ for the data points x_i , $x_i \in C_k$

Example

Gaussian radial basis function (RBF) kernel: $K(X_i, X_j) = e^{-||X_i - X_j||^2/2\sigma^2}$

Suppose there are 5 original 2-dimensional points:

 $\square x_1(0, 0), x_2(4, 4), x_3(-4, 4), x_4(-4, -4), x_5(4, -4)$

 \Box If we set σ to 4, we will have the following points in the kernel space

C E.g.,
$$||x_1 - x_2||^2 = (0 - 4)^2 + (0 - 4)^2 = 32$$
, therefore, $K(x_1, x_2) = e^{-\frac{32}{2 \cdot 4^2}} = e^{-1}$

22

Origi	nal S	pace		RBF Kernel Space ($\sigma = 4$)				
	x	y	$K(x_l, x_1)$	$K(x_l, x_2)$	$K(x_i, x_3)$	$K(x_i, x_4)$	$K(x_i, x_5)$	
<i>x</i> ₁	0	0	0	$e^{-\frac{4^2+4^2}{2\cdot4^2}} - e^{-1}$	e^{-1}	e ⁻¹	e ⁻¹	
<i>x</i> ₂	4	4	e^{-1}	0	e^{-2}	e^{-4}	e^{-2}	
<i>x</i> ₃	-4	4	e ⁻¹	e ⁻²	0	e ⁻²	e^{-4}	
<i>x</i> ₄	-4	-4	e^{-1}	e^{-4}	e^{-2}	0	e^{-2}	
<i>x</i> ₅	4	-4	e ⁻¹	e ⁻²	e^{-4}	e^{-2}	0	

Example: Kernel K-Means Clustering



- The above data set cannot generate quality clusters by K-Means since it contains noncovex clusters
- Gaussian RBF Kernel transformation maps data to a kernel matrix K for any two points $x_{i}, x_{j}: K_{x_{i}x_{j}} = \phi(x_{i}) \bullet \phi(x_{j})$ and Gaussian kernel: $K(X_{i}, X_{j}) = e^{-||X_{i} X_{j}||^{2}/2\sigma^{2}}$
- K-Means clustering is conducted on the mapped data, generating quality clusters

Recommended Readings

- J. MacQueen. Some Methods for Classification and Analysis of Multivariate Observations. In Proc. of the 5th Berkeley Symp. on Mathematical Statistics and Probability, 1967
- S. Lloyd. Least Squares Quantization in PCM. IEEE Trans. on Information Theory, 28(2), 1982
- A. K. Jain and R. C. Dubes. Algorithms for Clustering Data. Prentice Hall, 1988
- L. Kaufman and P. J. Rousseeuw. Finding Groups in Data: An Introduction to Cluster Analysis. John Wiley & Sons, 1990
- R. Ng and J. Han. Efficient and Effective Clustering Method for Spatial Data Mining. VLDB'94
- B. Schölkopf, A. Smola, and K. R. Müller. Nonlinear Component Analysis as a Kernel Eigenvalue Problem. Neural computation, 10(5):1299–1319, 1998
- I. S. Dhillon, Y. Guan, and B. Kulis. Kernel K-Means: Spectral Clustering and Normalized Cuts. KDD'04
- D. Arthur and S. Vassilvitskii. K-means++: The Advantages of Careful Seeding. SODA'07
- C. K. Reddy and B. Vinzamuri. A Survey of Partitional and Hierarchical Clustering Algorithms, in (Chap. 4) Aggarwal and Reddy (eds.), Data Clustering: Algorithms and Applications. CRC Press, 2014
- M. J. Zaki and W. Meira, Jr. Data Mining and Analysis: Fundamental Concepts and Algorithms. Cambridge Univ. Press, 2014

Next session!

Hierarchical Clustering

 Agglomerative: all cases start separately, then similar cases are joined

Divisive: all cases start in one category, then split

Dendrogram shows all levels of splits



DBSCAN

DBSCAN

- Density-based spatial clustering of applications with noise
- Works on local density
- Finds nonconvex and nonlinearly separable clusters
- Points can also be classified as noise



United in Leaning