

Chapter 2.1

Clustering Analysis and Algorithms

Xiangji Huang
York University, Canada

INTRODUCTION

Clustering is the process of grouping a collection of objects (usually represented as points in a multidimensional space) into classes of similar objects. Cluster analysis is a very important tool in data analysis. It is a set of methodologies for automatic classification of a collection of patterns into clusters based on similarity. Intuitively, patterns within the same cluster are more similar to each other than patterns belonging to a different cluster. It is important to understand the difference between clustering (unsupervised classification) and supervised classification.

Cluster analysis has wide applications in data mining, information retrieval, biology, medicine, marketing, and image segmentation. With the help of clustering algorithms, a user is able to understand natural clusters or structures underlying a data set. For example, clustering can help marketers discover distinct groups and characterize customer groups based on purchasing patterns in business. In biology, it can be used to derive plant and animal taxonomies, categorize genes

with similar functionality, and gain insight into structures inherent in populations.

Typical pattern clustering activity involves the following steps: (1) pattern representation (including feature extraction and/or selection), (2) definition of a pattern proximity measure appropriate to the data domain, (3) clustering, (4) data abstraction, and (5) assessment of output.

BACKGROUND

General references regarding clustering include Hartigan (1975), Jain and Dubes (1988), Kaufman and Rousseeuw (1990), Mirkin (1996), Jain, Murty, and Flynn (1999), and Ghosh (2002). A good introduction to contemporary data-mining clustering techniques can be found in Han and Kamber (2001). Early clustering methods before the '90s, such as k -means (Hartigan, 1975) and PAM and CLARA (Kaufman & Rousseeuw, 1990), are generally suitable for small data sets. CLARANS (Ng & Han, 1994) made improvements to CLARA in quality and scalability based

on randomized search. After CLARANS, many algorithms were proposed to deal with large data sets, such as BIRCH (Zhang, Ramakrishnan, & Livny, 1996), CURE (Guha, Rastogi, & Shim, 1998), Squashing (DuMouchel, Volinsky, Johnson, Cortes, & Pregibon, 1999) and Data Bubbles (Breuning, Kriegel, Kröger, & Sander, 2001).

MAIN THRUST

There exist a large number of clustering algorithms in the literature. In general, major clustering algorithms can be classified into the following categories.

Hierarchical Clustering

Hierarchical clustering builds a cluster hierarchy or a tree of clusters, also known as a *dendrogram*. Every cluster node contains child clusters; sibling clusters partition the points covered by their common parent. Such an approach allows the exploration of data on different levels of granularity. Hierarchical clustering can be further classified into *agglomerative* (bottom-up) and *divisive* (top-down) hierarchical clustering. An agglomerative clustering starts with one-point (singleton) clusters and recursively merges two or more most appropriate clusters. A divisive clustering starts with one cluster of all data points and recursively splits the most appropriate cluster. The process continues until a stopping criterion (for example, the requested number of k clusters) is achieved. Advantages of hierarchical clustering include (a) embedded flexibility regarding the level of granularity, (b) ease of handling of any forms of similarity or distance, and (c) applicability to any attribute types. Disadvantages of hierarchical clustering are (a) vagueness of termination criteria, and (b) the fact that most hierarchical algorithms do not revisit once-constructed clusters with the purpose of their improvement.

One of the most striking developments in hierarchical clustering is the algorithm BIRCH. BIRCH creates a height-balanced tree of nodes that summarize its zero, first, and second moments. Guha et al. (1998) introduced the hierarchical agglomerative clustering algorithm called CURE (Clustering Using Representatives). This algorithm has a number of novel features of general significance. It takes special care with outliers and with label assignment. Although CURE works with numerical attributes (particularly low-dimensional spatial data), the algorithm ROCK, developed by the same researchers (Guha, Rastogi, & Shim, 1999) targets hierarchical agglomerative clustering for categorical attributes.

Partitioning Clustering

Given a database of n objects and k , the number of clusters to form, a partitioning algorithm organizes the objects into k partitions ($k \leq n$), where each partition represents a cluster. The clusters are formed to optimize a partitioning criterion, often called a similarity function, such as distance, so that the objects within a cluster are similar, whereas the objects of different clusters are dissimilar in terms of the database attributes.

Partitioning clustering algorithms have advantages in applications involving large data sets for which the construction of a dendrogram is computationally prohibitive. A problem accompanying the use of a partitioning algorithm is the choice of the number of desired output clusters. A seminal paper (Dubes, 1987) provides guidance on this key design decision. The partitioning techniques usually produce clusters by optimizing a criterion function defined either locally (on a subset of the patterns) or globally (defined over all the patterns). Combinatorial search of the set of possible labelings for an optimum value of a criterion is clearly computationally prohibitive. In practice, the algorithm is typically run multiple times with different starting states, and the best

6 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/clustering-analysis-algorithms/24291

Related Content

Cross-Layer Distributed Attack Detection Model for the IoT

Hassan I. Ahmed, Abdurrahman A. Nasr, Salah M. Abdel-Mageid and Heba K. Aslan (2022). *International Journal of Ambient Computing and Intelligence* (pp. 1-17).

www.irma-international.org/article/cross-layer-distributed-attack-detection-model-for-the-iot/300794

Philosophical Foundations of Information Modeling

John M. Artz (2007). *International Journal of Intelligent Information Technologies* (pp. 59-74).

www.irma-international.org/article/philosophical-foundations-information-modeling/2423

Cost Efficiency Measures with Trapezoidal Fuzzy Numbers in Data Envelopment Analysis Based on Ranking Functions: Application in Insurance Organization and Hospital

Ali Ebrahimnejad (2012). *International Journal of Fuzzy System Applications* (pp. 51-68).

www.irma-international.org/article/cost-efficiency-measures-trapezoidal-fuzzy/68992

Exploring the Ethical Principles for the Implementation of Artificial Intelligence in Education: Towards a Future Agenda

Dilek enocak, Aras Bozkurt and Serpil Koçdar (2024). *Transforming Education With Generative AI: Prompt Engineering and Synthetic Content Creation* (pp. 200-213).

www.irma-international.org/chapter/exploring-the-ethical-principles-for-the-implementation-of-artificial-intelligence-in-education/338538

The Possibility of the Literary Work Generation by Computer

Akinori Abe (2016). *Computational and Cognitive Approaches to Narratology* (pp. 76-90).

www.irma-international.org/chapter/the-possibility-of-the-literary-work-generation-by-computer/159620