

Clustering Techniques for Big Data Analysis

Stelios Zimeras

 <http://orcid.org/0000-0002-9940-0041>

University of the Aegean, Greece

ABSTRACT

Clustering is the process by which data is classified into semantically consistent clusters based on some measure of similarity. Typically, clustering is an unsupervised machine learning problem, meaning that the structure of the data must be detected without any label being available as to which category it belongs to. Various clustering techniques have been developed, which aim to find coherent groups among a large number of data registered in large databases. We could say that the clustering technique is directly related to the optimization technique and thus its applications multiply in finding homogeneous groups of elements. This work deals with clustering algorithms and their application to big data. First, the clustering concept, objectives, and techniques are studied. Then, the main clustering algorithms are analyzed, their positive and negative characteristics, the steps to be followed for their application, their mathematical formulas, and a small application for each one on a small data set.

INTRODUCTION TO CLUSTERING ANALYSIS

Clustering is the process by which data is classified into semantically consistent clusters based on some measure of similarity (Jain & Dubes, 1988). Typically, clustering is an unsupervised machine learning problem, meaning that the structure of the data must be detected without any label being available as to which category it belongs to. The goal in clustering is to create groups, each of which will gather homogeneous elements and each of these groups maintains a center, the most central (Singh & Singh; 2024).

Various clustering techniques have been developed, aimed at finding coherent groups among a large amount of data. The clustering technique is directly related to the optimization technique. The mathematical formulation is given by the following analysis. Let be the data set $X = \{x_i \in R^k, k=1,2,\dots, i=1,2,\dots,n\}$ and let be a set of clusters $C = \{c_j, j=1,2,\dots,m\}$ and $m > 1$. The clustering problem is to find a matching $f: X \rightarrow C$ so that each observation x_i from the data set is associated with a cluster $C_j, 1 \leq j \leq m$, and for each observation the similarity between it and any other observation from the same cluster is greater than the similarity between it and any observation from other clusters (Singh & Singh, 2024; Zhou et. al., 2024;

DOI: 10.4018/407568

Moujahid & Dornaika, 2025). To solve a data clustering problem, the clustering stages are described based on the following steps:

- Representation of the characteristics
- Defining a metric
- Data clustering technique
- Evaluation of the final result.

Distance metrics

A metric or distance on a set X is a non-negative function $d: X \times X \rightarrow R$ for which the following three axioms hold:

$$d(x, y) = 0 \leftrightarrow x = y$$

$$d(x, y) = d(y, x)$$

$$d(x, z) \leq d(x, y) + d(y, z)$$

To cluster observations, we need to define a specified distance between points and between sets of points. Let consider dimension space n (in R^n) and two points $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$; the most common distances between the two points are (Figure 1):

Euclidean: The Euclidean distance between two points is the length of a line segment between two points

$$d(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^2 \right)^{1/2}$$

Manhattan: The Manhattan distance is a measurement of distance between two points in an N -dimensional vector space

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

18 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-techniques-for-big-data-analysis/407568

Related Content

Prediction of User Interests for Providing Relevant Information Using Relevance Feedback and Re-ranking

L. Sai Ramesh, S. Ganapathy, R. Bhuvaneshwari, K. Kulothungan, V. Pandiyaraju and A. Kannan (2015). *International Journal of Intelligent Information Technologies* (pp. 55-71).

www.irma-international.org/article/prediction-of-user-interests-for-providing-relevant-information-using-relevance-feedback-and-re-ranking/139740

Weight-Aware Multidimensional Advertising for TV Programs

Jianmin Wang, Yi Liu, Ting Xie and Yuchu Zuo (2013). *International Journal of Ambient Computing and Intelligence* (pp. 1-11).

www.irma-international.org/article/weight-aware-multidimensional-advertising-for-tv-programs/104157

Transformative Impact of AI in Green Finance: A Catalyst for Sustainable Development in India

Reenu Kumari, Komal Sharma, Rajesh Kumar and Ahu Cokun Özer (2024). *AI-Driven Decentralized Finance and the Future of Finance* (pp. 309-324).

www.irma-international.org/chapter/transformative-impact-of-ai-in-green-finance/355313

Institutionalization of Business Intelligence for the Decision-Making Iteration

Shaheb Ali, Rafiqul Islam and Ferdausur Rahman (2019). *International Journal of Intelligent Information Technologies* (pp. 101-118).

www.irma-international.org/article/institutionalization-of-business-intelligence-for-the-decision-making-iteration/221355

Cyber Security Patterns Students Behavior and Their Participation in Loyalty Programs

Witold Chmielarz and Oskar Szumski (2018). *International Journal of Ambient Computing and Intelligence* (pp. 16-31).

www.irma-international.org/article/cyber-security-patterns-students-behavior-and-their-participation-in-loyalty-programs/205573