

Preface

Classification is an integral part of any pattern recognition system. Depending on whether a set of labeled training samples is available or not, classification can be either supervised or unsupervised. Clustering is an important unsupervised classification technique where a number of data points are grouped into clusters such that points belonging to the same cluster are similar in some sense and points belonging to different clusters are dissimilar in the same sense. Cluster analysis is a complex problem as a variety of similarity and/or dissimilarity measures exist in the literature without any universal definition. In a crisp clustering technique, each pattern is assigned to exactly one cluster, whereas in the case of fuzzy clustering, each pattern is given a membership degree to each class. Fuzzy clustering is inherently more suitable for handling imprecise and noisy data with overlapping clusters.

For partitioning a data set, one has to define a measure of similarity or proximity based on which cluster assignments are done. The measure of similarity is usually data dependent. It may be noted that, in general, one of the fundamental features of shapes and objects is symmetry, which is considered to be important for enhancing their recognition. Examples of symmetry abound in real life, such as the human face, human body, flowers, leaves, and jellyfish. As symmetry is so common, it may be interesting to exploit this property while clustering a data set. Based on this observation, in recent years a large number of symmetry-based similarity measures have been proposed. This book is focused on different issues related to clustering, with particular emphasis on symmetry-based and metaheuristic approaches.

The aim of a clustering technique is to find a suitable grouping of the input data set so that some criteria are optimized. Hence, the problem of clustering can be posed as an optimization problem. The objective to be optimized may represent different characteristics of the clusters, such as compactness, symmetrical compactness, separation between clusters, connectivity within a cluster, etc. A straightforward way to pose clustering as an optimization problem is to optimize some cluster validity index that reflects the goodness of the clustering solutions. All possible values of the chosen optimization criterion (validity index) define the complete search space. Traditional partitional clustering techniques, such as K -means and fuzzy C -means, employ a greedy search technique over the search space in order to optimize

the compactness of the clusters. Although these algorithms are computationally efficient, they suffer from some drawbacks. They often get stuck at some local optima depending on the choice of the initial cluster centers. They are not able to determine the appropriate number of clusters from data sets and/or are capable of detecting clusters of some specific shape only.

To overcome the problem of getting stuck at local optima, several metaheuristic optimization tools, such as genetic algorithms (GAs), simulated annealing (SA), differential evolution (DE), etc., have been widely used to reach the global optimum value of the chosen validity measure. These techniques perform multimodal search in complex landscapes and provide near-optimal solutions for the objective or fitness function of an optimization problem. They have applications in fields as diverse as pattern recognition, image processing, neural networks, machine learning, job shop scheduling, and very large-scale integration (VLSI design), to mention just a few.

The two fundamental questions that need to be addressed in any typical clustering scenario are: (i) how many clusters are actually present in the data and (ii) how real or good is the clustering itself. That is, whatever the clustering technique, one has to determine the number of clusters and also the validity of the clusters formed. The measure of validity of clusters should be such that it will be able to impose an ordering of the clusters in terms of their goodness. Several cluster validity indices have been proposed in the literature, e.g., the Davies-Bouldin (DB) index, Dunn's index, Xie-Beni (XB) index, I -index, CS-index, etc., to name just a few. In recent years, several symmetry-based cluster validity indices have also been developed which are able to detect any kind of symmetric cluster from data sets. In this book, we discuss in detail several existing well-known cluster validity indices as well as some recently proposed symmetry-based versions.

Conventional GA-based clustering techniques, as well as related approaches, use some validity measure for optimization. Most of the existing clustering algorithms are able to determine hyperspherical/ellipsoidal-shaped clusters depending on the distance norm used. In recent years, some symmetry-based clustering techniques have been developed which can determine the appropriate number of clusters and the appropriate partitioning from data sets having any type of clusters, irrespective of their geometrical shape and overlapping nature, as long as they possess the characteristic of symmetry. A major focus of this book is on GA-based clustering techniques, which use symmetry as a similarity measure. Some of these clustering techniques are also able to detect the number of clusters automatically.

In many real-life situations one may need to optimize several objectives simultaneously. These problems are known as multiobjective optimization problems (MOOPs). In this regard, a multitude of metaheuristic single-objective optimization techniques such as genetic algorithms, simulated annealing, differential evolution, and their multiobjective versions have been developed. In this book we discuss in detail some existing single- and multiobjective optimization techniques. Moreover, a newly developed multiobjective simulated annealing-based technique is elaborately described and its effectiveness for solving several benchmark test problems is shown.

In the realm of clustering, simultaneous optimization of multiple validity indices, capturing different characteristics of the data, is expected to provide improved robustness to different data properties. Hence, it is useful to apply some multiobjective optimization (MOO) technique to solve the clustering problem. In MOO, search is performed over a number of objective functions. In contrast to single-objective optimization, which yields a single best solution, in MOO the final solution set contains a number of Pareto-optimal solutions, none of which can be further improved with regard to any one objective without degrading another. In a part of this book some recently developed multiobjective clustering techniques are discussed in detail.

A major focus of this book is on several real-life applications of clustering techniques in domains such as remote sensing satellite images, magnetic resonance (MR) brain images, and face detection. Analysis of remote sensing satellite images has significant utility in different areas such as climate studies, assessment of forest resources, examining marine environments, etc. An important task in remote sensing applications is the classification of pixels in the images into homogeneous regions, each of which corresponds to some particular land cover type. This problem has often been modeled as a segmentation problem, and clustering methods have been used to solve it. However, since it is difficult to have a priori information about the number of clusters in satellite images, the clustering algorithms should be able to automatically determine this value. Moreover, in satellite images it is often the case that some regions occupy only a few pixels, while the neighboring regions are significantly large. Thus, automatically detecting regions or clusters of such widely varying sizes presents a challenge in designing segmentation algorithms. In this book, we explore the applications of some symmetry-based clustering algorithms to classify remote sensing imagery in order to demonstrate their effectiveness.

Automatically classifying brain tissues from magnetic resonance images (MRI) has a major role in research and clinical study of neurological pathology. Additionally, regional volume calculations may provide even more useful diagnostic information. Automatic segmentation of brain MR images is a complex problem. Clustering approaches have been widely used for segmentation of MR brain images. A part of this book is dedicated to the applications of some symmetry-based clustering algorithms to classify MR brain images in order to demonstrate their effectiveness.

In recent years the problem of human face recognition has gained huge popularity. There are wide applications of face recognition systems including secure access control and financial transactions. The first important step of fully automatic human face recognition is human face detection. Face detection is the technique to automatically determine the locations and sizes of faces in a given input image. In this book a procedure to detect human faces based on symmetry is also described in detail.

This book aims to provide a treatise on clustering in a metaheuristic framework, with extensive applications to real-life data sets. Chapter 1 provides an introduction to pattern recognition, machine learning, classification, clustering, and related areas, along with different applications of pattern recognition. Chapter 2 describes in detail existing metaheuristic-based single- and multiobjective optimization techniques. An elaborate discussion on a recently developed multiobjective simulated annealing-based technique, archived multiobjective simulated annealing (AMOSA), is also

provided in this chapter. The utility of this technique to solve several benchmark test problems is shown. Chapter 3 mainly describes the different types of similarity measures developed in the literature for handling binary, categorical, ordinal, and quantitative variables. It also contains a discussion on different normalization techniques. Chapter 4 gives a broad overview of the existing clustering techniques and their relative advantages and disadvantages. It also provides a detailed discussion of several recently developed single- and multiobjective metaheuristic clustering techniques. Chapter 5 presents symmetry-based distances and a genetic algorithm-based clustering technique that uses this symmetry-based distance for assignment of points to different clusters and for fitness computation. In Chap. 6, some symmetry-based cluster validity indices are described. Elaborate experimental results are also presented in this chapter. Application of these symmetry-based cluster validity indices to segment remote sensing satellite images is also presented. In Chap. 7, some automatic clustering techniques based on symmetry are presented. These techniques are able to determine the appropriate number of clusters and the appropriate partitioning from data sets having symmetric clusters. The effectiveness of these clustering techniques is shown for many artificial and real-life data sets, including MR brain image segmentation. Chapter 8 deals with an extension of the concept of point symmetry to line symmetry-based distances, and genetic algorithm-based clustering techniques using these distances. Results are presented for some artificial and real-life data sets. A technique using the line symmetry-based distance for face detection from images is also discussed in detail. Some multiobjective clustering techniques based on symmetry are described in detail in Chap. 9. Three different clustering techniques are discussed here. The first one assumes the number of clusters a priori. The second and third clustering techniques are able to detect the appropriate number of clusters from data sets automatically. The third one, apart from using the concept of symmetry, also uses the concept of connectivity. A method of measuring connectivity between two points in a cluster is described for this purpose. A connectivity-based cluster validity index is also discussed in this chapter. Extensive experimental results illustrating the greater effectiveness of the three multiobjective clustering techniques over the single-objective approaches are presented for several artificial and real-life data sets.

This book contains an in-depth discussion on clustering and its various facets. In particular, it concentrates on metaheuristic clustering using symmetry as a similarity measure with extensive real-life applications in data mining, satellite remote sensing, MR brain imaging, gene expression data, and face detection. It is, in this sense, a complete book that will be equally useful to the layman and beginner as to an advanced researcher in clustering, being valuable for several reasons.

It includes discussions on traditional as well as some recent symmetry-based similarity measurements. Existing well-known clustering techniques along with metaheuristic approaches are elaborately described. Moreover, some recent clustering techniques based on symmetry are described in detail. Multiobjective clustering is another emerging topic in unsupervised classification. A multiobjective data clustering technique is described elaborately in this book along with extensive experimental results. A chapter of this book is wholly devoted to discussing some existing and

new multiobjective clustering techniques. Extensive real-life applications in remote sensing satellite imaging, MR brain imaging, and face detection are also provided in the book. Theoretical analyses of some recent symmetry-based clustering techniques are included.

The book will be useful to graduate students and researchers in computer science, electrical engineering, system science, and information technology as both text and reference book for some parts of the curriculum. Theoretical and experimental researchers will benefit from the discussions in this book. Researchers and practitioners in industry and R & D laboratories working in fields such as pattern recognition, data mining, soft computing, remote sensing, and biomedical imaging will also benefit.

The authors gratefully acknowledge the initiative and support for the project provided by Mr. Ronan Nugent of Springer. Sanghamitra Bandyopadhyay acknowledges the support provided by the Swarnajayanti Project Grant (No. DST/SJF/ET-02/2006-07) of the Department of Science and Technology, Government of India. The authors acknowledge the help of Mr. Malay Bhattacharyya, senior research fellow of ISI Kolkata, for drawing some of the figures of Chap. 5. Sriparna Saha also acknowledges the help of her students Mridula, Abhay, and Utpal for proofreading the chapters.

Kolkata, India

Sanghamitra Bandyopadhyay
Sriparna Saha

Unsupervised Classification

Similarity Measures, Classical and Metaheuristic
Approaches, and Applications

Bandyopadhyay, S.; Saha, S.

2013, XVIII, 262 p., Hardcover

ISBN: 978-3-642-32450-5