

Rough Sets, Kernel Set, and Spatiotemporal Outlier Detection

Alessia Albanese, *Member, IEEE*, Sankar K. Pal, *Fellow, IEEE*, and
Alfredo Petrosino, *Senior Member, IEEE*

Abstract—Nowadays, the high availability of data gathered from wireless sensor networks and telecommunication systems has drawn the attention of researchers on the problem of extracting knowledge from spatiotemporal data. Detecting outliers which are grossly different from or inconsistent with the remaining spatiotemporal data set is a major challenge in real-world knowledge discovery and data mining applications. In this paper, we deal with the outlier detection problem in spatiotemporal data and describe a rough set approach that finds the top outliers in an unlabeled spatiotemporal data set. The proposed method, called Rough Outlier Set Extraction (ROSE), relies on a rough set theoretic representation of the outlier set using the rough set approximations, i.e., lower and upper approximations. We have also introduced a new set, named Kernel Set, that is a subset of the original data set, which is able to describe the original data set both in terms of data structure and of obtained results. Experimental results on real-world data sets demonstrate the superiority of ROSE, both in terms of some quantitative indices and outliers detected, over those obtained by various rough fuzzy clustering algorithms and by the state-of-the-art outlier detection methods. It is also demonstrated that the kernel set is able to detect the same outliers set but with less computational time.

Index Terms—Spatiotemporal data, outlier detection, spatiotemporal uncertainty management, rough set and granular computing

1 INTRODUCTION

SPATIOTEMPORAL (ST) data mining is a growing research area dedicated to the discovery of hidden knowledge in large spatiotemporal databases, mainly through detecting periodic and/or frequent patterns and outliers. Particularly, outlier detection finds its applications in a broad spectrum of fields, such as fraud detection, intrusion detection in computer networking, and detecting motion or abnormal regions in image processing. The presence of outliers makes the modeling difficult due to the discordance the outliers introduce into the data; in this sense, the outlier detection task is attractive for two main reasons: the isolation of outliers, as a preventive step, can improve the performance of the predictive modeling by offering better data quality; on the contrary, the identification of outliers can be the main goal of the analysis as, for example, in fraud detection.

The most investigated approaches for outlier detection include:

1. distribution-based approaches that make use of standard statistical distribution to model the data declaring as outliers the objects that deviate from the model;
2. depth-based techniques that are based on computational geometry and compute different layers of

convex hulls declaring as outliers the objects belonging to the outer layers;

3. distance-based approaches that compute the proportion of database objects that are a specified distance from a target object; and
4. density-based approaches that assign a weight to each sample based on their local neighborhood density.

A different classification is based on the outlier detection output and divides into: labeling and scoring techniques. Labeling methods partition the data into two nonoverlapping sets (outliers and nonoutliers) and scoring methods offer a ranking list by assigning to each datum a factor reflecting its degree of outlierness. These former methods exploit a hard decision about the sets, the latter ones deal with a sort of soft decision about the membership of each datum to the set. The proposed method is the first rough method that improves and upgrades the “scoring methods,” proposing an effective soft granular computing-based solution exploiting the uncertainty region (boundary) to obtain more reliable results. Indeed, rough-set theory (RST) [41] is a paradigm to deal with uncertainty, vagueness, and incompleteness and it is proposed for indiscernibility in classification according to some similarity. Rough sets were extensively used for data mining but rarely for outlier detection in general-domain, the same for spatiotemporal specific-domain is hardly ever addressed and never for outlier detection in spatiotemporal data. In some sense, the few available outlier detection approaches interpret the rough set theory from the “operator-oriented point of view” [53]. In contrast, our method, called *Rough Outlier Set Extraction* (ROSE), exploits the set-oriented point of view of rough set theory to define the concept of outlier in terms of its lower and upper approximations (*rough outlier set*),

• A. Albanese and A. Petrosino are with the Computer Vision and Pattern Recognition Lab, Department of Applied Science, University of Naples Parthenope, Centro Direzionale Isola C4, I-80143 Naples, Italy. E-mail: {alessia.albanese, alfredo.petrosino}@uniparthenope.it.

• S.K. Pal is with the Indian Statistical Institute, 203, B.T. Road, Kolkata 700108, India. E-mail: sankar@isical.ac.in, skpal@ieee.org.

Manuscript received 25 Jan. 2012; revised 27 July 2012; accepted 11 Nov. 2012; published online 28 Nov. 2012.

Recommended for acceptance by J. Bailey.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-2012-01-0060. Digital Object Identifier no. 10.1109/TKDE.2012.234.

keeping into account those objects that can neither be ruled in nor ruled out as members of the target concept. Performance of ROSE in detecting outliers is found to be superior to best rough-fuzzy clustering algorithms in terms of various quantitative indices and to several state-of-the-art outlier detection methods.

Moreover, we introduce the concept of *kernel set*. Given a data set, the kernel set is a selected subset of elements able to describe the original data set in terms of data set structure. This paper includes two different versions of the ROSE algorithm on a test data set: one adopting, as input set, the entire set and the other adopting its kernel set. Experimental results show the advantages of considering the kernel set, in term of computational time, by comparing the *rough outlier set* extracted by the original data set with one extracted by the kernel set.

This paper is organized as follows: In Section 2, an overview on outlier detection approaches is given. Section 3 reports some preliminaries about rough set theory relevant to this work, indeed our approach is rough set based. Section 4 introduces the problem and reports the new rough set approach ROSE to extract the spatiotemporal rough outlier set. Section 5 introduces the new set kernel set. Sections 6.1, 6.2, and 6.3 present executed tests on three real-world (benchmark and test) data sets and the performance evaluation of the algorithm. Finally, conclusion remarks are given in Section 7 about ongoing and future work.

2 RELATED WORK

Most of the existing surveys on anomaly detection focus on a particular application domain or on a single research area, while the surveys, like [25], [14], [36] and two more recent brief surveys [44] and [49] are complete works that give the state of the art of anomaly detection techniques. They group anomaly detection into multiple categories and discuss techniques under each category. The discussed research issues include many topics to be taken into account to choose the appropriate outlier detection approach:

1. the detection method (parametric, i.e., distribution-based [7], depth-based [30], [29], [20]; graph-based methods [33], [48]; nonparametric, i.e., distance-based [31], [4], [43], [46]; density-based [12], [45], [54], [40], [55], [6]; clustering-based methods [24], [1], [21], [38]; and semiparametric, i.e., neural network-based, support vector machine-based techniques);
2. the nature of the detection algorithm, i.e., supervised, unsupervised, semi-supervised detection;
3. the nature of data, i.e., numerical, categorical, [11], [18] or mixed data [32], [37];
4. the size and the dimensionality of the data set, [2], [57], [47]; and
5. the nature of the target application [13], [22], [5].

This concerns the outlier detection methods in general domain. Concerning with specific spatiotemporal domain, only a few outlier detection methods have been proposed. Wu et al. [52] propose a spatiotemporal outlier detection algorithm called *Outstretch*, which discovers the sequences of spatial outliers over several time periods. Birant and Kut [9] present a ST-outlier detection approach based on

clustering concepts called *ST-DBSCAN* which is an improved version of the clustering technique *DBSCAN* [45] that supports also temporal aspects. Cheng and Li [17] further propose a four-step approach to detect spatiotemporal outliers, i.e., classification, aggregation, comparison, and verification. Wang et al. [50] also propose an approach to outlier detection in spatiotemporal domain. In a more recent work, Liu et al. [34] deal with the problem of detecting spatiotemporal outliers and causal relationships among them from traffic data streams.

Rough set theory has been recently introduced in the ST-domain literature for different aspects. In ST-domain, using the notion of rough sets, Bittner [10] defines approximations of ST-regions and relations between those approximations. Concerning outlier detection in general domain some works have been proposed: Nguyen [39] discusses a method for the detection and evaluation of outliers, as well as how to elicit the background domain knowledge from outliers using multilevel approximate reasoning schemes; Chen et al. [15] demonstrate an application of granular computing model using information tables for the outlier detection; Jiang et al. [27] propose a definition for outliers based on a rough outlier factor (ROF) as degree of outlierness for every object with respect to a given subset of universe. More recently, the same authors [28] propose a novel definition of outliers—sequence-based outliers—in information systems of rough set theory and an algorithm to find out such outliers. Concerning spatiotemporal outlier detection, no rough set theory-based approach has been proposed up to now.

3 ROUGH SET THEORY

Rough set theory, proposed by Pawlak [41], is a new and highly accepted paradigm used to deal with uncertainty, vagueness, and incompleteness. The main idea is based on the indiscernibility relation that describes indistinguishability of objects. Rough Set Theory can be approached as an extension of the Classical Set Theory, for use when representing incomplete knowledge. Concepts are represented by lower and upper approximations, according to which rough set methodology focuses on approximate representation of knowledge derivable from data [42].

3.1 Indiscernibility and Set Approximation

Let U be the universe of the discourse and A be the finite and nonempty set of attributes, then $S = \langle U, A \rangle$ is an information system. Let B a subset of A . With every subset of attributes $B \subseteq A$, an equivalence relation I_B on U can be easily associated:

$$I_B = \{(p, q) \in U \times U \mid \forall a \in B, a(p) = a(q)\}, \quad (1)$$

where I_B is called *B-indiscernibility relation*.

If $(p, q) \in I_B$, then objects p and q are indiscernible from each other by attributes B . The equivalence classes of the partition induced by the *B-indiscernibility relation* are denoted by $[p]_B$. These are also known as *granules*. We can approximate any subset X of U using only the information contained in B by constructing the lower and upper approximations of X . The sets $\{p \in U : [p]_B \subseteq X\}$ and $\{p \in U : [p]_B \cap X \neq \emptyset\}$, where $[p]_B$ denotes the equivalence

class of the object $p \in U$ relative to I_B , are called the B -lower and B -upper approximation of X in S and, respectively, denoted by $\underline{B}(X), \overline{B}(X)$. The objects in $\underline{B}(X)$ can be certainly classified as members of X on the basis of knowledge in B , while objects in $\overline{B}(X)$ can only be classified as possible members of X on the basis of B .

4 SPATIOTEMPORAL OUTLIER DETECTION

In this section, the spatiotemporal outlier detection problem is introduced by providing the problem formalization from a theoretical standpoint, together with its computational solution. A strict distinction between the spatial and temporal components is proposed in our definition of the problem. This may result useful in many contexts, for example, data sets which are characterized by only spatial information (we intend for spatial not only location information but also features detected at each location), where the temporal information is implicitly attached or is not present at all. In all such cases, the distinction allows us to consider just the spatial component, saving space, and time. In this way, time can be differently weighted for finding more efficiently temporal outlierness and for handling different scenarios, where spatial and temporal components get different importance in the data set. The proposed approach finds also spatiotemporal outliers.

4.1 Problem Definitions

Let us consider an information system $S = \langle U, A \rangle$ with U a spatiotemporal normalized data set and A its set of attributes. U can be written as follows:

$$U = \{p_i \equiv (z_{i1}, z_{i2}, \dots, z_{im}) \in [0, 1]^m, i = 1, \dots, N\},$$

where $p_i, i = 1, \dots, N$ is a m -dimensional feature vector and $A = \{a_1, a_2, a_3, \dots, a_m\}$ is the attribute set. In the following, we consider that at least three attributes must be present, i.e., the spatial attributes and the temporal one.

Given U , an integer $n > 0$ and a measure $d_{p_i}(U)$, defined over every $p_i \in U$, the general definition of the *Outlier Detection Problem* is as following:

Definition 1. The *Outlier Detection Problem* consists of finding $\bar{n} \geq n$ objects $p_1, p_2, \dots, p_n, p_{n+1}, \dots, p_{\bar{n}} \in U$ such that

$$\begin{aligned} d_{p_1}(U) &\geq d_{p_2}(U) \geq \dots \geq d_{p_n}(U) = d_{p_{n+1}}(U) \dots = d_{p_{\bar{n}}}(U) \\ &> d_{p_j}(U), \forall j = \bar{n} + 1, \dots, N. \end{aligned}$$

According to this definition, the concept of measure is used to determine the degree of dissimilarity of each object with respect to all others. Then, the n -Outlier Set can be formally defined as:

Definition 2. A n -Outlier Set $O \subseteq U$ is the set of $\bar{n} \geq n$ objects:

$$\begin{aligned} O = \{p_1, \dots, p_n, p_{n+1}, \dots, p_{\bar{n}} \in U \mid &d_{p_1}(U) \geq \dots \geq d_{p_n}(U) \\ &= d_{p_{n+1}}(U) \dots = d_{p_{\bar{n}}}(U) > d_{p_j}(U) \forall j = \bar{n} + 1, \dots, N\}, \end{aligned}$$

where $d_{p_i}(U), \forall i = 1, \dots, N$ is a measure defined and computed on U .

From Definition 2 it follows that $\tau = d_{p_n}(U)$ is the *outlierness threshold*, i.e., the minimum value among the

n maximum values of measures computed in U (associated with objects belonging to the n -Outlier Set), i.e.,

$$\tau = \inf\{\max_1(d_p(U), d_q(U)), \dots, \max_n(d_p(U), d_q(U))\}, \quad (2)$$

$$\forall p, q \in U.$$

Starting from the definition of spatial outlier and temporal outlier due to Birant and Kut [9] asserting: “a spatial outlier is a spatial referenced object whose nonspatial attribute values are significantly different from those of other spatially referenced objects in its spatial neighborhood,” and “a temporal outlier is an object whose nonspatial attribute value is significantly different from those of other objects in its temporal neighborhood,” we propose the following definitions applied only to spatiotemporal data:

Definition 3. A *Spatial Outlier (S-Outlier)* is an object whose spatial attribute value is significantly different from those of its closer objects (spatial neighborhood).

In this framework, the *Spatial Outlier* definition corresponds to:

Definition 4. Given U , an integer $n > 0$ and a measure on spatial component $d_{p_i}^s(U)$, defined over every $p_i \in U$, an object $p \in U$ is a *S-Outlier* iff $d_p^s(U) \geq \tau$ where τ is defined in (2).

Following Definition 4, it holds that:

Proposition 1. A *Spatial Outlier (S-Outlier)* is an object that belongs to the spatial n -Outlier Set indicated by O_s .

Similarly, we propose the following definition of *Temporal Outlier*, applied to only spatiotemporal data:

Definition 5. A *Temporal Outlier (T-Outlier)* is an object whose temporal attribute value is significantly different from those of its closer objects (temporal neighborhood).

In this framework, the *Temporal Outlier* definition corresponds to:

Definition 6. Given U , an integer $n > 0$ and a measure on temporal component $d_{p_i}^t(U)$, defined over every $p_i \in U$, an object $p \in U$ is a *T-Outlier* iff $d_p^t(U) \geq \tau$, where τ is defined in (2).

Equally, following Definition 6, it holds that:

Proposition 2. A *Temporal Outlier (T-Outlier)* is an object that belongs to the temporal n -Outlier Set indicated by O_t .

Definition 3 states that a spatial outlier has no objects or a small group of objects in its spatial neighborhood. The same is valid for a temporal outlier according to Definition 5. Following both definitions, the following holds:

Definition 7. A *Spatiotemporal Outlier (ST-Outlier)* is an object that respects both the definitions above.

To obtain a real degree of outlierness, an appropriate measure should be associated to each object; i.e., the euclidean distance computed between each object and all the other objects belonging to U . In real applications, characterized by an huge amount of data, this idea is unfeasible due to its high computational complexity ($O(N^2)$) where $N = |U|$.

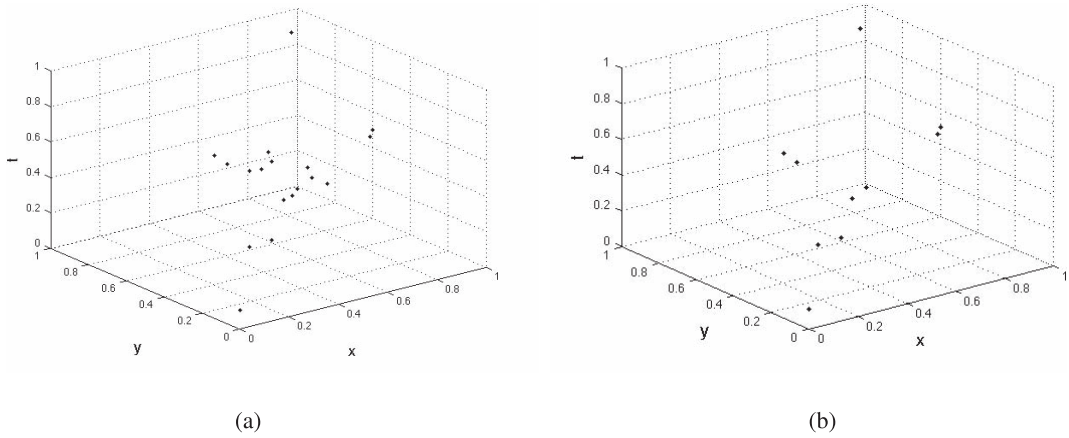


Fig. 1. (a) The example data set E and (b) its kernel set.

We preserve two aims: on one hand, we exploit the well-known outlier definition based on k -nearest neighbors [43], to associate to each object, a measure based on the distances among the object itself and its k -nearest neighbors rather than all N objects with $k \ll N$; on the other hand, we make use of a pruning strategy that discards objects that surely cannot belong to the n -Outlier Set, to address the problem of alleviating the computational cost.

In a Spatiotemporal context, the measure associated with each object is based upon the distances from its spatial k -nearest neighbors and its temporal k -nearest neighbors [3]. Precisely,

$$d_p^{s,t}(U) = \alpha \cdot d_p^s(U) + \beta \cdot d_p^t(U), \quad (3)$$

where

$$d_p^s(U) = \sum_{j=1}^k d^s(p, N^s(p, p_j)), \quad \forall p \in U, \quad (4)$$

$$d_p^t(U) = \sum_{j=1}^k d^t(p, N^t(p, p_j)), \quad \forall p \in U, \quad (5)$$

where $k > 0$ is the number of nearest neighbors to keep in account, $N^s(p, p_j)$ and $N^t(p, p_j)$ are, respectively, the j th spatial nearest neighbor and the j th temporal nearest neighbor of p , and α, β weight such that $\alpha + \beta = 1$. Definition 1, that introduces the Outlier Detection Problem, defines the Spatiotemporal Outlier Detection Problem, by selecting a measure as in (3).

To better illustrate the meanings of the previous and the following definitions, let us consider the example, a spatiotemporal data set $E = \{p_i \equiv (z_{i1}, z_{i2}, z_{i3}) \in [0, 1]^3, i = 1, \dots, 18\}$ where p_i is a three-dimensional feature vector and $A = \{a_1, a_2, a_3\}$ is the essential attribute set, i.e., a_1, a_2 are the spatial attributes and a_3 is the temporal attribute.

E is a labeled data set containing 18 elements as reported in Table 1 of Appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TKDE.2012.234>, and plotted in the Fig. 1. By fixing $k = 3$ and $n = 4$, the outlier sets (spatial, temporal outlier sets), on the basis of the previous definitions, are computed as follows: A 4-Spatial Outlier Set $O_s \subseteq E$ is the set of objects $p \in E$ that

significantly deviate from the rest of data with respect to the spatial component, i.e., $O_s = \{(0.95, 0.55, 0.50), (1, 0.60, 0.50), (0.01, 0.01, 0.1), (0.9, 0.9, 0.95)\}$. A 4-Temporal Outlier Set $O_t \subseteq E$ is the set of objects $p \in E$ that significantly deviate from the rest of data with respect to the temporal component, i.e., $O_t = \{(0.01, 0.01, 0.1), (0.20, 0.21, 0.3), (0.30, 0.22, 0.3), (0.9, 0.9, 0.95)\}$. If $n = 2$, a 2-Spatiotemporal Outlier Set $O_{s,t} \subseteq E$ is the set of objects $p \in E$ that significantly deviate from the rest of data with respect to the spatial and the temporal component, i.e., $O_{s,t} = \{(0.01, 0.01, 0.1), (0.9, 0.9, 0.95)\}$. O_s , O_t , and $O_{s,t}$ are shown in Fig. 2a as diamond and square, as triangle and square and only square, respectively. In Fig. 2b, a 2D projection has been reported to better visualize that the spatial outliers and spatiotemporal outliers are spatially far from the rest of data.

4.2 Rough Outlier Set Extraction

4.2.1 Theory

The goal of our approach is to exploit the rough set theory to define the Outlier Set such as a Rough Outlier Set.

Let $S = \langle U, A \rangle$ be an information system with U a spatiotemporal normalized data set and A its attribute set. If $n > 0$ is the required outlier number, we want to describe $O \subseteq U$ (n -Outlier Set) as

$$\langle \underline{B}(O), \overline{B}(O) \rangle (\text{Rough } n - \text{Outlier Set}), \quad (6)$$

where $\underline{B}(O)$ is the B -Lower approximation and $\overline{B}(O)$ is the B -Upper approximation of n -Outlier Set with respect to an attribute subset $B \subseteq A$.

The B -Lower approximation $\underline{B}(O)$ is defined as the set of objects that can be certainly classified as members of the set O on the basis of the knowledge in B , while the B -Upper approximation $\overline{B}(O)$ is defined as the set of possible members of O on the basis of the knowledge in B .

With this aim, let I_B be the B -indiscernibility relation on the universe U :

$$I_B = \{(p_i, p_j) \in U \times U : a(p_i) = a(p_j), \forall a \in B\}.$$

The equivalence classes $[p_j]_B$ or granules G_j of the partition induced by I_B on U are such that

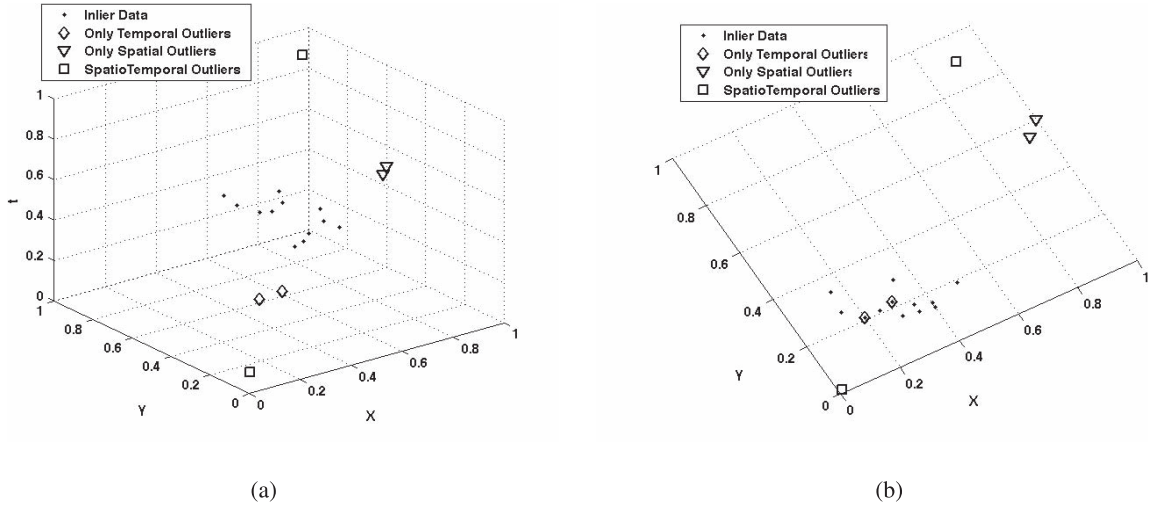


Fig. 2. Example data set: (a) Detected outlier sets (b) their xy -projection.

$$U = \bigcup_{j=1}^N G_j \text{ and } G_j \cap G_j = \emptyset, i \neq j.$$

The measure in (3) is used as a spatiotemporal weight $\bar{w}_{G_j}(s, t, i)$, to be assigned to every granule G_j , depending on space, s , and/or on time, t , and at iteration, i . The attribute subsets B include spatiotemporal attributes, or only spatial and only temporal attribute to define spatiotemporal outlier set, or only temporal set and only spatial outlier set, respectively. In this framework, the B -Lower and B -Upper approximations at iteration i can be defined as follows:

Definition 8. The B -Lower approximation $\underline{B}_i(O)$ of n -Outlier Set O , at iteration i , is

$$\underline{B}_i(O) = \{G_j \subseteq U : \bar{w}_{G_j} > \tau_i\},$$

where

$$\tau_i = \inf \{ \max_1^i (\bar{w}_{G_j}, \bar{w}_{G_k}), \dots, \max_n^i (\bar{w}_{G_j}, \bar{w}_{G_k}) \}, \quad (7)$$

$$\forall G_j, G_k \subseteq U.$$

Definition 9. The B -upper approximation $\bar{B}_i(O)$ of n -Outlier Set O , at iteration i , is

$$\bar{B}_i(O) = \{G_j \subseteq U : \bar{w}_{G_j} > \bar{\tau}_i\},$$

where

$$\bar{\tau}_i = \tau_{i-1}, \forall i \geq 2. \quad (8)$$

The threshold τ_1 is computed as the minimum value among the n higher values of weights assigned to the granules at first iteration, then, at second iteration, τ_2 will be the new minimum value among the new n higher values of weights reassigned to the granules at second iteration and $\bar{\tau}_2 = \tau_1$.

The iterative procedure will stop when the following convergence criterion will be satisfied:

Lemma 1. The construction of the lower approximation $\underline{B}(O)$ or the upper approximation $\bar{B}(O)$ of an n -Outlier Set O converges if it exists an index k such that the threshold does not vary anymore, i.e.,

$$\text{if } \bar{\tau}_k = \tau_k \text{ then } \underline{B}_k(O) = \bar{B}_k(O). \quad (9)$$

Proof. See Appendix, available in the online supplemental material. \square

Hence, the *Rough n -Outlier Set* is represented by

$$\langle \underline{B}_{k-1}(O), \bar{B}_{k-1}(O) \rangle. \quad (10)$$

In case of $B = A$ (every attribute is considered), the granules are

$$\forall p_j \in U : \{p_j\} \equiv G_j \forall j = 1, \dots, N, \quad (11)$$

so both spatial and temporal components are taken into account.

As instance, let us consider the labeled Example data set. In this case, the attribute set is $A = \{x, y, t\}$, i.e., x and y are cartesian coordinates and t is the temporal component.

Spatial Outliers In the case of spatial outliers, the reduction is made in terms of temporal component, i.e., $B = \{t\}$; so we have the following partition of the universe:

$$I_B = I_{\{t\}} = \{\{p_1, p_2\}, \{p_3, p_9\}, \{p_4\}, \{p_5\}, \{p_6\}, \{p_7, p_8\}, \{p_{10}\}, \{p_{11}\}, \{p_{12}\}, \{p_{13}\}, \{p_{14}\}, \{p_{15}\}, \{p_{16}\}, \{p_{17}\}, \{p_{18}\}\}.$$

The concept of *Spatial Outlier* can be appropriately defined on the basis of knowledge in $B = \{t\}$. Specifically, the B -lower approximation of the *Spatial Outlier Set* O_s is composed by the granules completely included into O_s , i.e., $\underline{B}(O_s) = \{\{p_7, p_8\}, \{p_{17}\}, \{p_{18}\}\}$ and the B -upper approximation is composed by the granules that have nontrivial intersection with O_s , i.e., $\bar{B}(O_s) = \{\{p_7, p_8\}, \{p_{17}\}, \{p_{18}\}\}$. In this case, the upper approximation does not give any additional information.

Temporal Outliers In the case of temporal outliers, the reduction is made by spatial components, i.e., $B = \{x, y\}$, getting

$$I_B = I_{\{x, y\}} = \{\{p_1, p_{12}\}, \{p_2, p_{13}\}, \{p_3\}, \{p_4\}, \{p_5\}, \{p_6\}, \{p_7\}, \{p_8\}, \{p_9\}, \{p_{10}\}, \{p_{11}\}, \{p_{14}\}, \{p_{15}\}, \{p_{16}\}, \{p_{17}\}, \{p_{18}\}\}.$$

The concept of *Temporal Outlier* can be equivalently get on the basis of knowledge in $B = \{x, y\}$. The B -lower

approximation of the Temporal Outlier Set O_t is composed by the granules completely included into O_t , i.e., $\underline{B}(O_t) = \{\{p_{17}\}, \{p_{18}\}\}$ and the B -upper approximation is composed by the granules that have a nontrivial intersection with O_t , i.e., $\overline{B}(O_t) = \{\{p_1, p_{12}\}, \{p_2, p_{13}\}, \{p_{17}\}, \{p_{18}\}\}$. In this case, the notion of rough set arises; indeed, the upper approximation gives additional information.

4.2.2 ROSE Algorithm

The *Rough Outlier Set Extraction Algorithm* is designed to receive as input the universe U , the number k of nearest neighbors, and the number n of outliers to detect. The output of the (iterative) procedure is the *Rough Outlier Set* (Upper, Lower Approximation, and Negative Region). The algorithm selects, at each step, a small subset of objects, called *WorkingSet*, from the overall data set U . To this aim, *ExtractElements* extracts a number of elements equal to a fixed percentage of the cardinality of U that has to be greater than k . The following main steps are computed. For all selected objects, the procedure computes the euclidean distances among the objects in the *WorkingSet* and all the objects of U , considering the spatial components, the temporal components or both of them (general case $B = A$) depending upon the chosen attribute subset B with respect to the Rough Outlier Set has been computed. Algorithm ROSE related to the general case has been shown. *UpdateUpperApprox* and *UpdateLowerApprox* at first iteration create the same set of n top outliers at that step, i.e., the n objects that have an associated measure higher than the others. Then, at next iterations, *UpdateUpperApprox* and *UpdateLowerApprox* compute the Lower and Upper approximation of *Rough Outlier Set*, using the τ (computed by *LowerWeight*) and τ_{prev} thresholds as, respectively, defined in (7) and (8). At each iteration i , the pruning strategy selects objects from U that have their measure under the computed threshold to build the Negative Region. The *LowerWeight* function computes the τ threshold (and consequently τ_{prev} is the saved value of τ before to be updated). At each iteration, the thresholds have been computed as the weight minimum value among the weight maximum n values, as defined in (7)). The difference set between the Universe set and the Negative Region is the Kernel Set.

Algorithm 1. ROSE - Rough Outlier Set Extraction.

```

beginROSEExtraction( $U, n, k$ )
  LowerOutlierSet = null; UpperOutlierSet = null
   $w_{s,t,k}(q) = 0$ 
   $\tau_{prev} = 0; \tau = 0$ 
  WorkingSet = ExtractElements( $U$ )
  while (WorkingSet! = null) do
    for  $p \in U$  do
      for  $q \in WorkingSet$  do
        if (LowerOutlierSet == null and
            UpperOutlierSet == null)
          or ( $w_{s,t,k}(q) \geq \tau_{prev}$ ) then
             $d_s(p, q) = CalculateSpDistance(p, q)$ 
             $d_t(p, q) = CalculateTempDistance(p, q)$ 
            BuildTreeKNN( $p, q, d_s, d_t, k$ )
        else
          AddNegativeRegion( $p$ )

```

```

      end if
    end for
  end for
  for  $q \in WorkingSet$  do
     $w_{s,t,k}(q) = CalculateWeight(q)$ 
    UpperOutlierSet = UpdateUpperApprox
      ( $\tau_{prev}, n, w_{s,t,k}(q)$ )
    LowerOutlierSet = UpdateLowerApprox( $\tau, n, w_{s,t,k}(q)$ )
  end for
   $\tau = LowerWeight(UpperOutlierSet)$ 
  if ( $\tau! = 0$ ) then
     $\tau_{prev} = \tau$ 
  end if
   $U = U - WorkingSet$ 
  WorkingSet = ExtractElements( $U$ )
end while
end ROSEExtraction()

```

4.2.3 ROSE Algorithm—Time Complexity

The ROSE algorithm has worst-case time complexity $O(|U|^2)$, but practical complexity $O(|U|^{1+d})$, with $d < 1$ and U the universe.

5 THE KERNEL SET: RELEVANCE TO OUTLIER DETECTION

The present section introduces a new set, called kernel set, and states that it is a relevant set for outlier detection. Given a data set U , the kernel set is a subset, of lower cardinality, that can be used instead of U , to detect the same outlier set. The time complexity reduction of the use of kernel set is quantified by measuring kernel set dimensionality over that of U .

5.1 Definition

Let us now define a new set, called *Kernel Set*, $K \subseteq U$, as a selected subset of the universe U that characterizes the overall data set. Intuitively, this set is a subset of objects of U that maintains the general structure of the universe U . The Kernel Set is built by construction, in an iterative way, adding each object having specific properties.

Definition 10. Given U and two integers $n > 0, k > 0$ (number of nearest neighbors), $d(U)$ a measure defined on U , the Kernel Set K is built by adding each object $p \in U$ such that one of the following properties holds:

1. $d_p(U) \geq \tau$,
2. if $d_p(U) < \tau$, then $\exists q \in U$ such that $p \in NN^k(q)$ and $d_q(U) < \tau$ and $d_q(K - \{p\}) \geq \tau$,

where $NN^k(q)$ is the set of k -nearest neighbors of q and $d(K)$ is the restriction of $d(U)$ on $K \subseteq U$.

The Definition 10 states that the objects that belong to the Kernel Set are:

1. object p for which $d_p(U) \geq \tau$ and, hence, belongs to n -Outlier Set.
2. object p that, even if $d_p(U) < \tau$, is one of the nearest neighbors of an object q for which $d_q(U) < \tau$ and $d_q(K - \{p\}) \geq \tau$.

The second property states that once these objects p have been added to K , the measure of the object q becomes less than τ also in K as in U . Otherwise, the global structure of the data set should be altered.

Also, the *Kernel Set* is built for the Example data set like:

$$K = \{(0.01, 0.01, 0.1), (0.9, 0.9, 0.95), (0.95, 0.55, 0.5), \\ (1.0, 0.6, 0.5), (0.2, 0.21, 0.3), (0.3, 0.22, 0.3), \\ (0.3, 0.16, 0.55), (0.35, 0.15, 0.6), (0.15, 0.26, 0.76), \\ (0.16, 0.34, 0.77)\}.$$

This set is also reported in Fig. 1b. The *Kernel Set* contains all elements of the *Outlier Set*.

5.2 Properties

Let us start to prove the following propositions related to the new set.

Proposition 3. *The measure computed in K is an upper bound of the measure computed in U such that*

$$d_p(U) \leq d_p(K), \forall p \in U,$$

where $d_p(U) = \sum_{j=1}^k d(p, N(p, p_j))$ and $N(p, p_j)$ is the j th nearest neighbor of p .

Proof. See Appendix, available in the online supplemental material. \square

The following proposition is valid:

Proposition 4. *A Kernel Set contains the n -Outlier Set: $K \supseteq O$.*

Proof. $\forall p \in O : d_p(U) > \tau \Rightarrow p \in K$.

The proof clearly follows from definition of K . \square

Proposition 5. *The Outlier Set O_K computed starting from Kernel Set K is a superset of O computed from U :*

$$O_K \supseteq O.$$

Proof. See Appendix, available in the online supplemental material. \square

5.3 Significance to Outlier Detection

The *kernel Set* is a meaningful subset of the universe U with the following properties:

- *Kernel Set* is a subset with lower cardinality than U ,
- the “same results” in terms of *rough outlier set* are obtained using *Kernel Set* instead of U .
- *Kernel Set* can be considered as the model learned during a training phase.

In the following, we propose the comparison between the obtained results, in terms of *rough outlier set*, executing ROSE algorithm, once using, as input, the entire universe U and another time computed using, as input, the *kernel Set* K .

5.4 Computational Benefits

Let us consider the two versions (or runnings) of ROSE algorithm, to appreciate the computational benefits. At the first run, ROSE algorithm receives, as input, the entire data set U , while at the second run, ROSE receives the kernel set K of U that is a subset of U . A computational benefit, coming from using kernel set instead of the entire universe,

is derived. Indeed, $O(|U|^{1+d}) < O(|K|^{1+d})$, being $K \subset U$. To quantify the computational benefits coming from the use of the kernel set, we evaluate the dimensionality of kernel set K with respect to U . The experimental results have been provided in the following Section 6.4.

6 EXPERIMENTAL RESULTS AND DISCUSSION

Our outlier detection method is based on rough set theory and is specific for spatiotemporal data. At the best of our knowledge, there is no rough approach to outlier detection for spatiotemporal data to compare with. Hence, three different experimental tests have been executed. The first test is oriented to demonstrate the ability of the outlier detection algorithm and the role of the kernel set working on a real-world spatiotemporal data set; the comparisons on this data set are made using rough-fuzzy clustering methods. The second test is intended to compare our results with other outlier detection methods (also rough-oriented) for general domain on a UCI repository data set. The third test is oriented to compare our performance with outlier detection methods (not rough approach) tailored for spatiotemporal domain, on a spatiotemporal data set. Two Sections 6.4 and 6.5 end this section: one concerning an experimental evaluation of the dimension reduction percentage of the kernel set with respect to its starting data set U and one concerning a sensitivity analysis about the parameters k and n of the algorithm.

6.1 School Buses Data Set

For the first test, we make tests on a real-world data set, named School Buses [19]. The data set is publicly available and consists of 145 trajectories (about 69,000 entries) of two school buses collecting and delivering students around Athens metropolitan area in Greece for 108 distinct days. The structure of each record is as follows: $\{obj_id, traj_id, date, time, lat, lon, x, y\}$ where obj_id is the school bus identification, $traj_id$ is the unique trajectory identification, the date and time are the sampling time stamps every 30 seconds (date in $dd/mm/yyyy$ format and time in $hh:mm:ss$ format), the (lat, lon) and (x, y) are the bus location, in WGS84 and in GGRS87 reference systems, respectively. In our case, the obj_id and $traj_id$ are not considered, date and time fields are converted in just one field t consisting of a time string corresponding to the elements year, month, day, hour, minute, and second. Moreover, the lat and lon are redundant and are not considered, because x and y give the same information. Hence, the normalized representation of the data set is illustrated in Fig. 3b: in a 3D cartesian reference system, x and y are the spatial coordinates and the third dimension is time t . In Fig. 3a, the trajectory map of school buses is shown. In the following Fig. 3c, the testing data set consisting of half of the original data set (about 30,000 entries) with some added temporal outliers is shown.

6.1.1 Rough Outlier Set Extraction—Spatial Rough Outlier Set Extraction from U

Let U denote the spatiotemporal normalized School Buses data set:

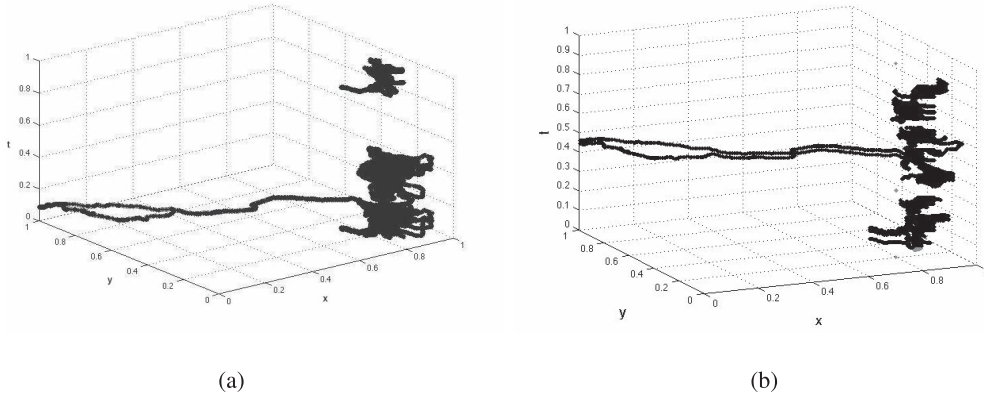


Fig. 3. School Buses data set: (a) Normalized data set, (b) testing subset with added temporal outliers highlighted in gray color.

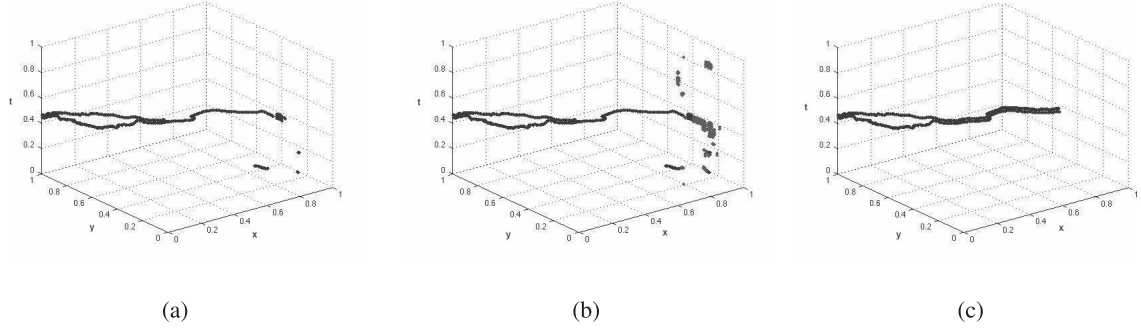


Fig. 4. (a) Intermediate step: lower approx, (b) Intermediate step: Lower Approx U boundary, (c) Last-1 Step: Lower approx.

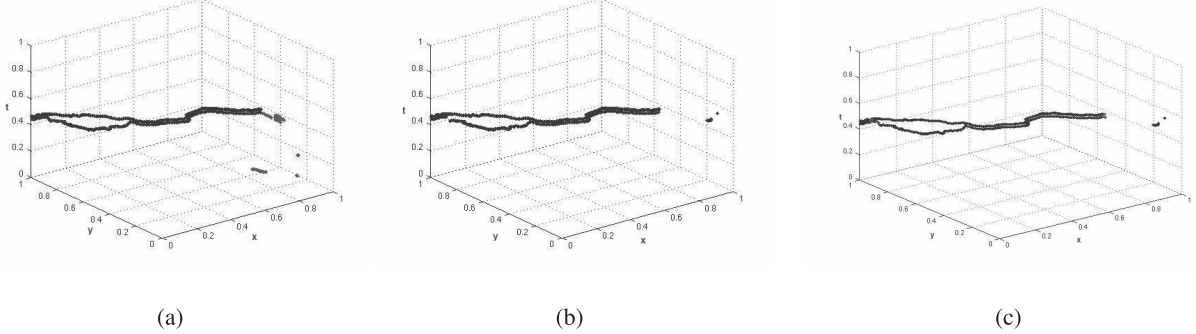


Fig. 5. (a) Last-1 step: lower Approx U boundary, (b) last step: lower approx, (c) last step: lower Approx U boundary.

$$U = \{p_i \equiv (z_{i,1}, z_{i,2}, z_{i,3}) \in [0, 1]^3, i = 1, \dots, N\},$$

where $(z_{i,1}, z_{i,2})$ are cartesian coordinates of the i th object, $z_{i,3}$ is the relative time stamp. Let $\langle U, A \rangle$ be the information system, with the attribute set $A = \{x, y, t\}$, i.e., x and y are the spatial components and t is the temporal component.

We want to describe $O \subseteq U$ (*Outlier Subset*) as the *rough outlier subset* $\langle B(O), \bar{B}(O) \rangle$, where $B \subseteq A$ is constituted by the spatial attributes, (x, y) . Selecting only spatial components, the results of selected iterations, an *intermediate* step, the *last-1* and the *last* one have been shown. Specifically, the lower, upper approximation (lower and boundary) at an intermediate step of *Spatial Rough Outlier Set* are represented and shown in Figs. 4a and 4b, where boundaries are reported in gray color.

Figs. 4c and 5a show the lower, upper approximation (lower and boundary) at *last-1* step, while Figs. 5b and 5c show the same approximations at *last* step. In the last figure, we can see the advantages of keeping into account

the boundary. Otherwise, many interesting objects (belonging to the boundary) should be missed.

6.1.2 Rough Outlier Set Extraction—Spatiotemporal Rough Outlier Set Extraction from U

Let $\langle U, A \rangle$ be the information system, with the attribute set $A = \{x, y, t\}$, i.e., x and y are the spatial components and t is the temporal component. Now we are considering $B = A$, so we are looking for *spatiotemporal Rough Outlier Set*.

The spatiotemporal outliers will be more relevant than spatial and temporal outliers (see temporal outliers injected in the Fig. 3b). Hence, the lower approximation includes the most part of spatial and temporal outliers, while the upper approximation includes the remaining part of temporal outliers and some other spatial outliers have been detected. In this section, we show the lower, lower approximation with boundary at last step. Fig. 6a shows the lower approximation, while Fig. 6b shows the lower approximation with boundaries in gray color.

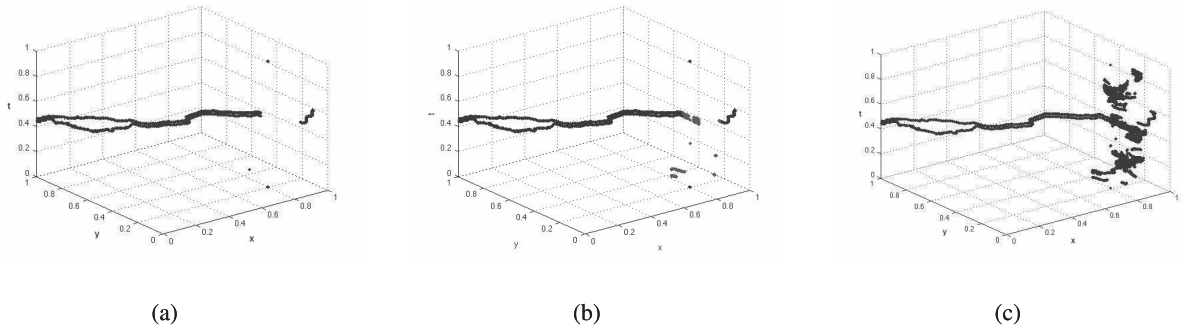


Fig. 6. (a) Last step: lower approx, (b) Last step: lower Approx U boundary, (c) School Buses data set: its Kernel set.

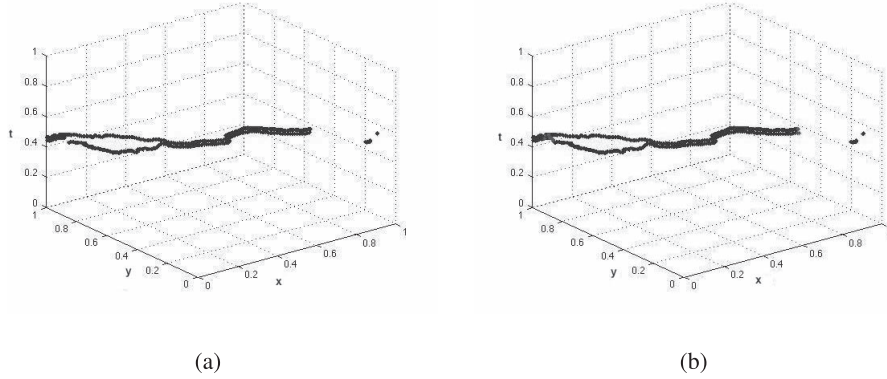


Fig. 7. ROSE results from Kernel set of School Buses data set—last step: (a) lower approx, (b) lower Approx U boundary.

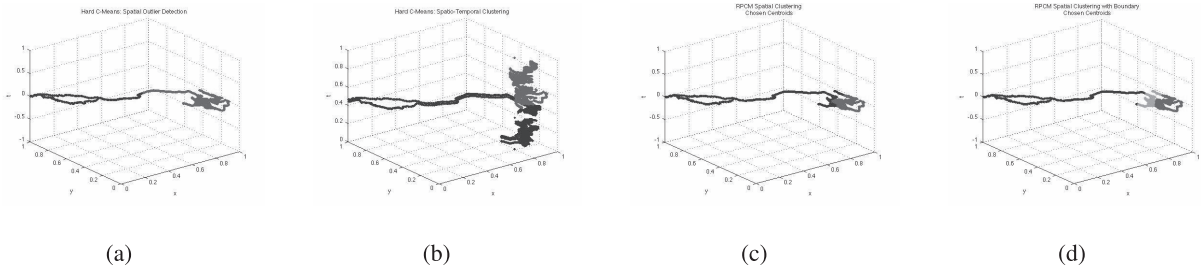


Fig. 8. Hard C-Means Clusters Results: (a) spatial outlier detection, (b) spatiotemporal outlier detection—spatial outlier detection: (c) RPCM clusters results, (d) RPCM clusters results with boundary.

6.1.3 Rough Outlier Set Extraction—Spatial Rough Outlier Set Extraction from the Kernel Set

The section reports the tests aimed to demonstrate the use of the Kernel Set. This set is a selected subset, able to describe the original data set both in terms of data structure and in terms of obtained results. In particular, we want to show the advantages of using this set and the benefits of considering it. To this aim, we show the rough outlier set extracted by the universe U and the rough outlier set extracted by the kernel set. The results show the advantages of considering this set. Fig. 6c shows the Kernel set of School Buses data set. Starting from the *Kernel Set*, the rough outlier set is built by our approach ROSE. Let be $B \subseteq A$ constituted by the spatial attributes, i.e., (x, y) . Selecting only spatial components, the results of last iteration of the test of spatial rough outlier set extraction from the Kernel set is reported. Fig. 7a shows the lower approximation at the last iteration, while Fig. 7b shows the lower approximation with boundaries in gray color. Thus, we compare these results with the last test of rough outlier set extraction from the entire Universe U , shown in Fig. 5c.

Comparing Figs. 5c and 7b, we can appreciate that the results are quite similar with an interesting computational benefit coming from considering the *Kernel set* instead of the entire universe U .

6.1.4 Quantitative Measures and Indices

In this section, we use performance indices as introduced by Maji and Pal in [35] such as α index, ρ index, and γ index, to evaluate the performance of our algorithm compared with *Hard C-Means* and with other *rough-fuzzy* clustering algorithms, incorporating the concepts of rough sets. So, the algorithms adopted for comparison are: Hard C-Means, RFCM—Rough Fuzzy C-Means, RPCM—Rough Possibilistic C-Means, RFPCM—Rough Fuzzy Possibilistic C-Means. To analyze the performance of our proposed algorithm, tests have been performed on the School Buses data set. Figs. 8a and 8b show the clusters computed by Hard C-Means clustering algorithm (number of clusters set to 2) in spatial and spatiotemporal outlier detection, respectively. Figs. 8c and 8d and Fig. 9 show the results of each rough-fuzzy algorithm in spatial outlier detection.

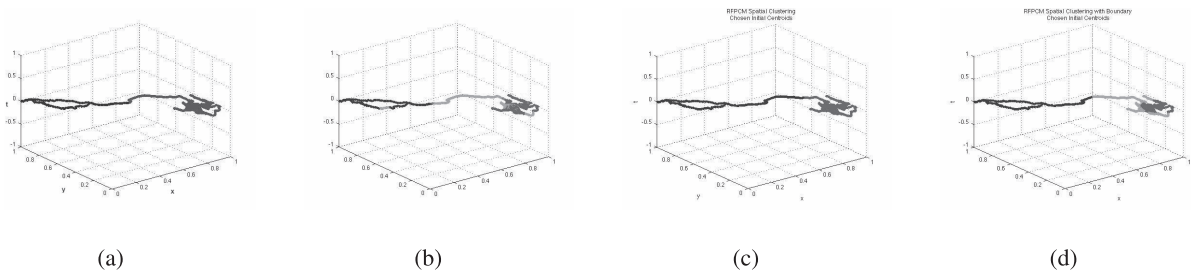


Fig. 9. Spatial outlier detection: (a) RFCM clusters results, (b) RFCM clusters results with boundary, (c) RFPCM clusters results, (d) RFPCM clusters results with boundary.

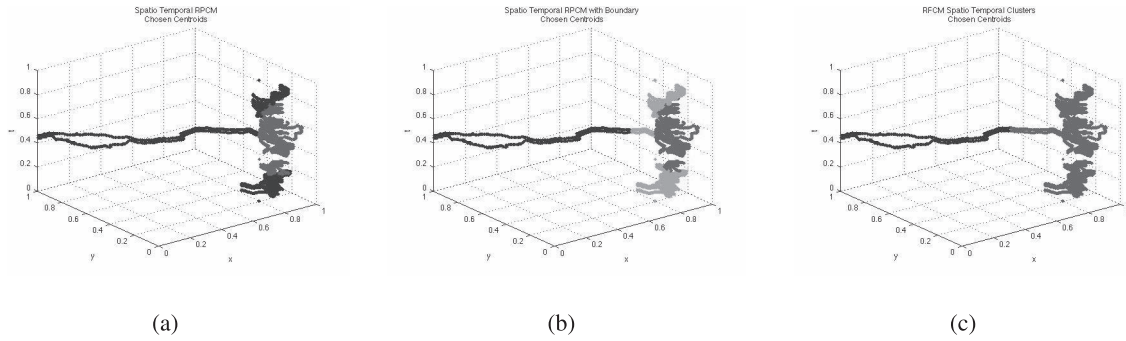


Fig. 10. ST outlier detection: (a) RPCM clusters results, (b) RPCM clusters results with boundary, (c) RFPCM clusters results.

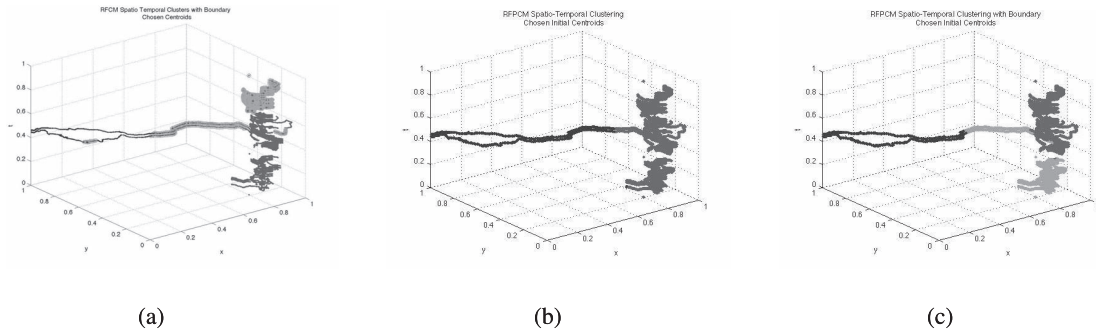


Fig. 11. ST outlier detection: (a) RFCM clusters results with boundary, (b) RFPCM clusters results, (c) RFPCM clusters results with boundary.

TABLE 1
Spatial Outlier Detection—Quantitative Evaluation of Algorithms—Chosen Initial Centroids

Methods	α Index	ρ Index	γ Index
ROSE	0.9836	0.0164	0.9987
RFCM	0.5448	0.4551	0.9250
RPCM	0.4725	0.5274	0.7919
RFPCM	0.5645	0.4354	0.9007

Legenda:

ROSE = Rough Outlier Set Extraction
RFCM = Rough Fuzzy C-Means
RPCM = Rough Possibilistic C-Means
RFPCM = Rough Fuzzy Possibilistic C-Means

In Figs. 9a and 9c, the two clusters are drawn with gray and black colors after the assignment of the boundary to clusters, while in the Figs. 9b and 9d the boundaries (before the assignment) are drawn with light gray color.

Figs. 10 and 11 show the results of rough-fuzzy algorithms in spatiotemporal outlier detection. The parameters have been set as follows: $c = 2$ (inlier and outlier cluster), ω and $\tilde{\omega}$ are equal to 0.5 to give the same importance to the lower approximation and to the boundary. Several runs have been done with different initializations and different parameters, related to initial centroid choice. These parameters have been maintained constant across all runs. The tests show that the best results are obtained for particular choices of initial centroids rather than for random choices of initial centroids. So, we report

only the final prototypes of the best solution. Tables 1 and 2 report the best results obtained using different algorithms for $c = 2$ in case of the same choice of initial centroids for HCM, RFCM, RPCM, and RFPCM. Tables 1 and 2 compare the performance of these different rough-fuzzy clustering algorithms with respect to α , ρ , γ in spatial and spatiotemporal outlier detection, respectively. The results reported in Tables 1 and 2 establish the fact that although the hybridization versions of c -means algorithm were not designed as outlier detectors, they generate good prototypes for $c = 2$. In spatial outlier detection, the RFPCM provides the best results as shown in Fig. 9; the results of other two versions of rough clustering are quite similar to that of the RFPCM, while in spatiotemporal outlier detection, the RPCM outperforms them as shown in

TABLE 4
ROSE Results: Comparison on Grand St. Bernard Data Set—Spatial and Temporal Outliers

Methods	Running Average		Mahalanonis Dist.		Density	
	DR(%)	FPR(%)	DR(%)	FPR(%)	DR(%)	FPR(%)
TOD:	72.3	10.5	100	15.0	100	15.1
ROSE _T Low:	80	1.2	96	1.0	100	0
ROSE _T Upp:	87.5	1.7	100	1.2	100	0
SOD:	24.5	3.3	100	4.3	100	4.4
POD:	29.8	1.8	80	3.7	75	3.8
ROSE _S Low:	96.2	1.0	92	1.0	92	0.3
ROSE _S Upp:	98.1	1.2	96	1.1	100	1.1

TABLE 5
Kernel Set Dimension Computation on Different Data Sets

Dataset Name	Universe cardinality	Kernel set cardinality	Reduction %
School Buses	30414	17101	13313 (44%)
Wisconsin Breast Cancer - Original	699	486	213 (30%)
Wisconsin Breast Cancer - Unbalanced	483	308	175 (36%)
Grand St. Bernard Dataset	2101	1535	566 (27%)

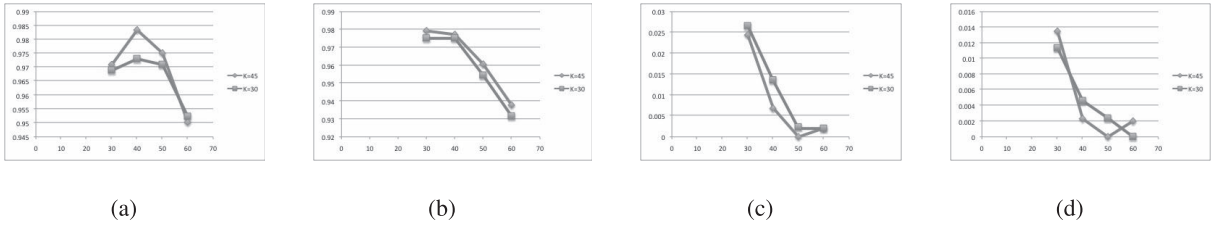


Fig. 12. Wisconsin Breast Cancer data set—for two fixed k values: (a) Accuracy: lower approximation, (b) Accuracy: upper approximation, (c) False Alarm Probability: lower approximation, (d) False Alarm Probability: upper approximation.

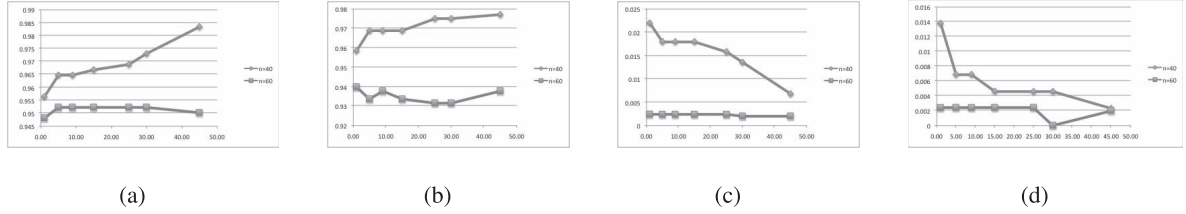


Fig. 13. Wisconsin Breast Cancer data set—for two fixed n values: (a) Accuracy: lower approximation, (b) Accuracy: upper approximation, (c) False Alarm Probability: lower approximation, (d) False Alarm Probability: upper approximation.

positives on running average technique and in a bit lower or comparable way than SOD on the other two labeling techniques. Globally, the achieved ROSE results outperform the compared state-of-the-art techniques on this spatiotemporal data set.

6.4 Kernel Set Dimension

Experimental computations, about the dimension reduction between four analyzed data sets and their kernel sets, have been widely executed. The kernel sets dimensions, reported in the Table 5 are the average dimensions on 10 executions, varying the input parameters of the ROSE algorithm. The computed data provide an average value of reduction percentage equal to 46 percent. The analyzed data sets are the following: School Buses, Wisconsin Breast Cancer (original version), Wisconsin Breast Cancer (unbalanced version), and Grand St. Bernard. This reduction significantly drops down the algorithm time complexity and, hence, its computational cost.

6.5 Sensitivity Analysis of Input Parameters

This section ends this evaluation section and is intended to conduct a sensitivity analysis about the input parameters k

and n of the algorithm to evaluate the algorithm behavior. The comparison have been done on the Wisconsin Breast Cancer Data Set doing several different combinations of k and n parameters. In particular, the n and k parameters have been chosen in the following way: 1) keeping the value of n fixed (at 40, at 60) the value of k was varying at 1, 5, 9, 15, 25, 30, 45 and 2) keeping the value of k fixed (at 30, at 45) the value of n was varying at 30, 40, 50, and 60. The results have been shown in the following figures: the first couple of Figs 12a and 12b shows the accuracy curves for $k = 45$ and $k = 30$ varying n ; the second couple of Figs. 13a and 13b shows the accuracy curves for $n = 40$ and $n = 60$ varying k . Then, the Figs. 12c and 12d and the Figs. 13c and 13d show the false alarm probability curves for $k = 45$ and $k = 30$ varying n and those for $n = 40$ and $n = 60$ varying k , respectively. For $n = 50$ and $k = 45$, a reversal trend between lower and upper approximation as for $n = 40$ and a bit lower accuracy for $n = 50$ respect to $n = 40$ clearly appear. Hence, increasing too much the number n of outliers to be searched not surely improve the results. A zero false alarm probability has been reported in both cases for $k = 45$.

7 CONCLUSIONS

The manuscript extends outlier detection using a new rough set approach to spatiotemporal data. Specifically, the rough set-based outlier detection method has been theoretically grounded based on a definition of outlier set as rough set. A remarkable note should be made for the definition of a new set, called kernel set, that has been demonstrated to be able to generate the “same” output results in terms of rough outlier set with time computational benefits. The experimental results on three real-world data sets prove that the performance of ROSE in detecting outliers are superior when compared to several other methods. On the real-world School Buses data set, ROSE has been compared with C-Means clustering algorithm and other rough-fuzzy clustering algorithms (Rough Fuzzy C-Means, Rough Possibilistic C-Means, Rough Fuzzy Possibilistic C-Means), incorporating the concepts of rough sets, producing reasonable results both in terms of quantitative and qualitative standpoints. On the benchmark Wisconsin Breast Cancer data set, ROSE has been also compared with several state-of-the-art outlier detection methods, also rough-oriented, for general domain (SEQ, DIS, NED, KNN, RNN), demonstrating higher, and just sometimes comparable, performance. Another comparison has been made on the WSN Grand ST. Bernard data set with spatiotemporal methods (Zhang’s TOD, SOD, POD) that use the same data set, demonstrating the ROSE superiority even in this case. The approach is computationally less intensive compared with these approaches. The ROSE algorithm appear to consistently outperform other rough and not rough approaches in medium to large problem settings, showing to be able to do well also on data sets of varying sizes. Since spatiotemporal outlier detection might turn out to be useful in many different research fields, we hope that this work will spark further interest in such problems that are challenging and relatively unexplored.

ACKNOWLEDGMENTS

S.K. Pal acknowledges the J.C. Bose Fellowship of the Government of India.

REFERENCES

- [1] C.C. Aggarwal and P. Yu, “Finding Generalized Projected Clusters in High Dimensional Spaces,” *Proc. ACM SIGMOD Int’l Conf. Management Data*, pp. 70-81, 2000.
- [2] C.C. Aggarwal and P.S. Yu, “An Effective and Efficient Algorithm for High-Dimensional Outlier Detection,” *VLDB J.*, vol. 14, pp. 211-221, 2005.
- [3] A. Albanese and A. Petrosino, “A Non Parametric Approach to the Outlier Detection in Spatio-Temporal Data Analysis,” *Information Technology and Innovation Trends in Organizations*, D’Atri, et al., eds., pp. 101-108, Springer Verlag, 2011.
- [4] F. Angiulli and C. Pizzuti, “Outlier Mining in Large High-Dimensional Data Sets,” *IEEE Trans. Knowledge and Data Eng.*, vol. 17, no. 2, pp. 203-215, Feb. 2005.
- [5] F. Angiulli and F. Fasseti, “Distance-Based Outlier Queries in Data Streams: The Novel Task and Algorithms,” *J. Data Mining and Knowledge Discovery*, vol. 20, no. 2, pp. 290-324, 2010.
- [6] M. Ankerst, M.M. Breunig, H.-P. Kriegel, and J. Sander, “Optics: Ordering Points To Identify The Clustering Structure,” *Proc. ACM SIGMOD Int’l Conf. Management Data (SIGMOD ’99)*, pp. 49-60, 1999.
- [7] V. Barnett and T. Lewis, *Outliers in Statistical Data*. John Wiley & Sons, 1994.
- [8] S.D. Bay, “The UCI KDD Repository,” <http://kdd.ics.uci.edu>, 1999.
- [9] D. Birant and A. Kut, “Spatio-Temporal Outlier Detection in Large Databases,” *J. Computing and Information Technology*, vol. 14, no. 4, pp. 291-297, 2006.
- [10] T. Bittner, “Rough Sets in Spatio-Temporal Data Mining,” *Proc. First Int’l Workshop Temporal, Spatial, and Spatio-Temporal Data Mining-Revised Papers (TSDM ’00)*, pp. 89-104, 2000.
- [11] S. Boriah, V. Chandola, and V. Kumar, “Similarity Measures for Categorical Data: A Comparative Evaluation,” *Proc. Eighth SIAM Int’l Conf. Data Mining*, pp. 243-254, 2008.
- [12] M.M. Breunig, H.-P. Kriegel, R.T. Ng, and J. Sander, “LOF: Identifying Density Based Local Outliers,” *Proc. ACM SIGMOD Int’l Conf. Management of Data*, pp. 93-104, 2000.
- [13] A. Ceglar, J.F. Roddick, and D.M.W. Powers, “CURIO: A Fast Outlier and Outlier Cluster Detection Algorithm for Large Data Sets,” *Proc. Second Int’l Workshop Integrating Artificial Intelligence and Data Mining*, pp. 37-45, 2007.
- [14] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly Detection: A Survey,” *ACM Computing Surveys*, vol. 41, no. 3, pp. 15:1-15:58, 2009.
- [15] Y. Chen, D. Miao, and R. Wang, “Outlier Detection Based on Granular Computing,” *Proc. Sixth Int’l Conf. Rough Sets and Current Trends Computing*, pp. 283-292, 2008.
- [16] Y. Chen, D. Miao, and H. Zhang, “Neighborhood Outlier Detection,” *Expert Systems with Applications*, vol. 37, no. 12, pp. 8745-8749, 2010.
- [17] T. Cheng and Z. Li, “A Multiscale Approach to Detect Spatio-Temporal Outliers,” *Trans. GIS*, vol. 10, no. 2, pp. 253-263, 2006.
- [18] K. Das and J. Schneider, “Detecting Anomalous Records in Categorical Data Sets,” *Proc. 13th ACM SIGKDD Int’l Conf. Knowledge Discovery and Data Mining*, pp. 220-229, 2007.
- [19] E. Frentzos, K. Gratsias, N. Pelekis, and Y. Theodoridis, “Nearest Neighbor Search on Moving Object Trajectories,” *Proc. Ninth Int’l Symp. Spatial and Temporal Databases (SSTD ’05)*, pp. 328-345, 2005.
- [20] A.K. Ghosh and P. Chaudhuri, “On Maximum Depth Classifiers,” *Scandinavian J. Statistics*, vol. 32, no. 2, pp. 327-350, 2005.
- [21] S. Guha, R. Rastogi, and K. Shim, “CURE: An Efficient Clustering Algorithm for Large Databases,” *Proc. ACM SIGMOD Int’l Conf. Management Data*, vol. 27, no. 2, pp. 73-84, 1998.
- [22] J.M.P. Gutierrez and J.F. Gregori, *Clustering Techniques Applied to Outlier Detection of Financial Market Series Using a Moving Window Filtering Algorithm*, Unpublished working paper series, no. 948, European Central Bank, pp. 1-45, 2008.
- [23] S. Harkins, H.X. He, G.J. Williams, and R.A. Baxter, “Outlier Detection Using Replicator Neural Networks,” *Proc. Fourth Int’l Conf. Data Warehousing and Knowledge Discovery*, pp. 170-180, 2002.
- [24] Z. He, X. Xu, and S. Deng, “Discovering Cluster-Based Local Outliers,” *J. Pattern Recognition Letters*, vol. 24, pp. 1641-1650, 2003.
- [25] V. Hodge and J. Austin, “A Survey of Outlier Detection Methodologies,” *J. Artificial Intelligence Rev.*, vol. 22, no. 2, pp. 85-126, 2004.
- [26] F. Ingelrest, G. Barrenetxea, G. Schaefer, M. Vetterli, O. Couach, and M. Parlange, “SensorScope: Application-Specific Sensor Network for Environmental Monitoring,” *J. ACM Trans. Sensor Networks*, vol. 6, no. 2, pp. 1-32, 2010.
- [27] F. Jiang, Y. Sui, and C. Cao, “Outlier Detection Based on Rough Membership Function,” *Proc. Fifth Int’l Conf. Rough Sets and Current Trends Computing (RSCTC ’06)*, pp. 388-397, 2006.
- [28] F. Jiang, Y. Sui, and C. Cao, “Some Issues about Outlier Detection in Rough Set Theory,” *Expert Systems with Applications*, vol. 36, no. 3, pp. 4680-4687, 2009.
- [29] R. Jörnsten, “Clustering and Classification Based on the L1 Data Depth,” *J. Multivariate Analysis*, vol. 90, no. 1, pp. 67-89, 2004.
- [30] T. Johnson, I. Kwok, and R.T. Ng, “Fast Computation of 2-Dimensional Depth Contours,” *Proc. Fourth Int’l Conf. Knowledge Discovery and Data Mining*, pp. 224-228, 1998.
- [31] E. Knorr and R. Ng, “Algorithms for Mining Distance-Based Outliers in Large Data Sets,” *Proc. 24th Int’l Conf. Very Large Data Bases (VLDB ’98)*, pp. 392-403, 1998.
- [32] A. Koufakou and M. Georgiopoulos, “A Fast Outlier Detection Strategy for Distributed High-Dimensional Data Sets with Mixed Attributes,” *Data Mining and Knowledge Discovery*, vol. 20, no. 2, pp. 259-289, 2010.
- [33] J. Laurikkala, M. Juhola, and E. Kentala, “Informal Identification of Outliers in Medical Data,” *Proc. Fifth Workshop Intelligent Data Analysis Medicine Pharmacology (IDAMAP)*, pp. 20-24, 2000.

- [34] W. Liu, Y. Zheng, S. Chawla, J. Yuan, and X. Xie, "Discovering Spatio-Temporal Causal Interactions in Traffic Data Streams," *Proc. 17th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 1010-1018, 2011.
- [35] P. Maji and S.K. Pal, "Rough Set Based Generalized Fuzzy C-Means Algorithm and Quantitative Indices," *IEEE Trans. Systems, Man, and Cybernetics-Part B: Cybernetics*, vol. 37, no. 6, pp. 1529-1540, Dec. 2007.
- [36] M. Markos and S. Sameer, "Novelty Detection: A Review Part 1: Statistical Approaches," *Signal Processing*, vol. 83, no. 12, pp. 2481-2497, 2003.
- [37] E. Muller, I. Assent, U. Steinhausen, and T. Seidl, "OutRank: Ranking Outliers in High Dimensional Data," *Proc. IEEE 24th Int'l Conf. Data Eng. Workshop*, pp. 600-603, 2008.
- [38] R.T. Ng and J. Han, "CLARANS: A Method for Clustering Objects for Spatial Data Mining," *IEEE Trans. Knowledge and Data Eng.*, vol. 14, no. 5, pp. 1003-1016, Sept./Oct. 2002.
- [39] T.T. Nguyen, "Outlier Detection: An Approximate Reasoning Approach," *Proc. Int'l Conf. Rough Sets and Intelligent Systems Paradigms (RSEISP '07)*, pp. 495-504, 2007.
- [40] S. Papadimitriou, H. Kitagawa, P.B. Gibbons, and C. Faloutsos, "LOCI: Fast Outlier Detection Using the Local Correlation Integral," *Proc. 19th Int'l Conf. Data Eng. (ICDE '03)*, pp. 315-326, 2003.
- [41] Z. Pawlak, *Rough Sets: Theoretical Aspects of Reasoning About Data*. Kluwer, 1991.
- [42] Z. Pawlak and A. Skowron, "A Rough Set Approach for Decision Rules Generation," *Proc. Workshop W12: The Management Uncertainty in AI at 13th IJCAI*, 1993.
- [43] S. Ramaswamy, R. Rastogi, and K. Shim, "Efficient Algorithms for Mining Outliers from Large Data Sets," *Proc. ACM SIGMOD Int'l Conf. Management Data*, pp. 427-438, 2000.
- [44] N.N.R. Ranga Suri, N. Murty, and G. Athithan, "Data Mining Techniques for Outlier Detection," *Visual Analytics and Interactive Technologies: Data, Text and Web Mining Applications*, chapter 2, pp. 22-38, IGI Global Snippet, 2010.
- [45] J. Sander, M. Ester, H.-P. Kriegel, and X. Xu, "Density-Based Clustering in Spatial Databases: The Algorithm GDBSCAN and Its Applications," *Data Mining and Knowledge Discovery*, vol. 2, no. 2, pp. 169-194, 1998.
- [46] P. Sun and S. Chawla, "On Local Spatial Outliers," *Proc. IEEE Fourth Int'l Conf. Data Mining*, pp. 209-216, Nov. 2004.
- [47] Y. Tao, X. Xiao, and S. Zhou, "Mining Distance-Based Outliers from Large Databases in Any Metric Space," *Proc. 12th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 394-403, 2006.
- [48] P.M. Valero Mora, F.W. Young, and M. Friendly, "Visualizing Categorical Data in ViSta," *Computational Statistics Data Analysis*, vol. 43, no. 4, pp. 495-508, 2003.
- [49] K. Venkateswara Rao, A. Govardhan, and K.V. Chalapati Rao, "Spatio Temporal Data Mining: Issues, Task and Applications," *Int'l J. Computer Science Eng. Survey*, vol. 3, no. 1, pp. 39-52, 2012.
- [50] X.R. Wang, J.T. Lizier, O. Obst, M. Prokopenko, and P. Wang, "Spatiotemporal Anomaly Detection in Gasmonitoring Sensor Networks," *Proc. European Conf. Wireless Sensor Networks (EWSN)*, pp. 90-105, 2008.
- [51] G.J. Williams, R.A. Baxter, H.X. He, S. Harkins, and L.F. Gu, "A Comparative Study of RNN for Outlier Detection in Data Mining," *Proc. IEEE Int'l Conf. Data Mining (ICDM '03)*, pp. 709-712, 2002.
- [52] E. Wu, W. Liu, and S. Chawla, "Spatio-Temporal Outlier Detection in Precipitation Data," *Proc. Second Int'l Conf. Knowledge Discovery from Sensor Data*, pp. 115-133, 2008.
- [53] Y.Y. Yao, "Two Views of the Theory of Rough Sets in Finite Universes," *Int'l J. Approximate Reasoning*, vol. 15, pp. 291-317, 1996.
- [54] T. Zhang, R. Ramakrishnan, and M. Livny, "BIRCH: An Efficient Data Clustering Method for Very Large Databases," *Proc. ACM SIGMOD Int'l Conf. Management Data*, vol. 25, no. 2, pp. 103-114, 1996.
- [55] Y. Zhang, S. Yang, and Y. Wang, "LDBOD: A Novel Local Distribution Based Outlier Detector," *Pattern Recognition Letters*, vol. 29, no. 7, pp. 967-976, 2008.
- [56] Y. Zhang, N.A.S. Hamm, N. Meratnia, A. Stein, M. van de Voort, and P.J.M. Havinga, "Statistics-Based Outlier Detection for Wireless Sensor Networks," *Int'l J. Geographical Information Science*, vol. 26, no. 8, pp. 1373-1392, 2012.
- [57] C. Zhu, H. Kitagawa, and C. Faloutsos, "Example-Based Robust Outlier Detection in High Dimensional Data Sets," *Proc. IEEE Fifth Int'l. Conf. Data Mining (ICDM '05)*, pp. 829-832, Nov. 2005.



member of the IEEE Computer Society and a member of the IEEE.



Sankar K. Pal received the PhD degrees from Calcutta University and Imperial College, London. He joined the Indian Statistical Institute in 1975 as a CSIR senior research fellow where he became a full professor in 1987, a distinguished scientist in 1998, and the director in 2005. He is a distinguished scientist of the Institute and its former director. He is also a J.C. Bose fellow of the Government of India. He founded the Machine Intelligence Unit and the Center for Soft Computing Research at the Institute in Calcutta which are enjoying international recognition. He was at UC Berkeley and UMD, College Park, the NASA JSC, Houston, Texas, and the US Naval Research Lab, Washington DC. He has been serving as a distinguished visitor of the IEEE Computer Society since 1987 and held several visiting positions in Italy, Poland, Hong Kong, and Australian Universities. He is a fellow of the TWAS, IAPR, IFSA, and all four National Academies for science/engineering in India. He is a coauthor of 17 books and more than 300 research publications in the areas of pattern recognition and machine learning, image processing, data mining, web intelligence, soft computing, and bioinformatics. He is/was on the editorial boards of 20 journals including IEEE Transactions. He has received several national and international awards including the most coveted S.S. Bhatnagar Prize in India in 1990 and Padma Shri in 2013. He is a fellow of the IEEE.



Alfredo Petrosino is a professor of computer science at the University of Naples Parthenope, where he heads the research laboratory CVPRLab at the University of Naples Parthenope (cvprlab.uniparthenope.it). He held positions at the University of Salerno, International Institute of Advanced Scientific Studies (IIASS), National Institute for the Physics of Matter (INFN), and lastly as a researcher and a senior researcher at the National Research Council (CNR). He taught at the Universities of Salerno, Siena, Naples Federico II, and Naples Parthenope. He is a member of the International Association for Pattern Recognition, US and the International Neural Networks Society, US. He has coedited six books and more than 100 research publications in the areas of computer vision, image, and video analysis, pattern recognition, neural networks, fuzzy and rough sets, data mining. He is/was an associate editor of the *Pattern Recognition Journal*, member of the editorial board of the *Pattern Recognition Letters*, the *International Journal of Knowledge Engineering and Soft Data Paradigms*, and editor of the *IEEE Transactions SMC-Part A, Fuzzy Sets and Systems, Image and Vision Computing, Parallel Computing*. He is a senior member of the IEEE.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.