

A parallel matrix-based method for computing approximations in incomplete information systems

Junbo Zhang, Jian-Syuan Wong, Yi Pan, *Senior Member, IEEE*, and Tianrui Li*, *Senior Member, IEEE*

Abstract—As the volume of data grows at an unprecedented rate, large-scale data mining and knowledge discovery present a tremendous challenge. Rough set theory, which has been used successfully in solving problems in pattern recognition, machine learning, and data mining, centers around the idea that a set of distinct objects may be approximated via a lower and upper bound. In order to obtain the benefits that rough sets can provide for data mining and related tasks, efficient computation of these approximations is vital. The recently introduced cloud computing model, MapReduce, has gained a lot of attention from the scientific community for its applicability to large-scale data analysis. In previous research, we proposed a MapReduce-based method for computing approximations in parallel, which can efficiently process complete data but fails in the case of missing (incomplete) data. To address this shortcoming, three different parallel matrix-based methods are introduced to process large-scale, incomplete data. All of them are built on MapReduce and implemented on Twister that is a lightweight MapReduce runtime system. The proposed parallel methods are then experimentally shown to be efficient for processing large-scale data.

Index Terms—Rough sets, Data mining, MapReduce, Matrix, Incomplete information systems

1 INTRODUCTION

With the development of information technology, the data volume is growing at an unprecedented rate. Large-scale data mining and knowledge discovery has been a tremendous challenge. The parallel program model MapReduce supports large distributed data sets on clusters of computers [1], [2], which allows for an efficient analysis of large amounts of data. As one of the most important and effective cloud computing techniques, MapReduce has become a popular computing model for cloud platforms. Many cloud platforms support the MapReduce framework. For instance, Amazon EC2 and

Microsoft Azure both support open source MapReduce Hadoop [3] and Twister [4] runtime systems such as Amazon Elastic MapReduce [5] and Twister4Azure [6].

Rough set theory is a powerful mathematical tool that can be used to process inconsistent information in decision-making situations [7], [8], [9], [10]. It plays an important role in the fields of pattern recognition, machine learning and data mining during last decades [11], [12], [13], [14], [15], [16], [17], [18], [19], [20]. For example, Hu et al. proposed a novel feature evaluation measure, called the neighborhood decision error rate, and then presented a forward greedy strategy for searching the feature subset based on neighborhood rough sets [14]. Qian et al. introduced a theoretic framework based on rough set theory, called positive approximation, which can be used to accelerate a heuristic process for feature selection from incomplete data [18].

To our knowledge, the most current methods based on rough sets are sequential and can only run on a single computer to deal with small data sets. For example, ROSETTA is such a toolkit that can be used for data analysis, which can be downloaded <http://www.lcb.uu.se/tools/rosetta/>. To expand the applications of rough sets to the field of large scale data mining and knowledge discovery, we proposed a parallel method for computing approximations based on MapReduce in our previous work [21], [22]. However, existing studies on parallel methods are mainly based on complete information systems. They do not consider the missing data in real databases, which frequently occur in decision-making problems. This motivates the need for an investigation of parallel data mining methods based on rough sets in incomplete information systems (IIS).

Our contributions in this paper include:

- 1) Traditional methods for computing approximations in IIS are based on crisp set view, which are difficult to parallelize and fail to process large-scale incomplete data with rough set theory. To achieve this goal, a new matrix view of the rough set theory in IIS is presented and the relation matrix is then introduced for computing approximations.
- 2) A MapReduce-based parallel method to construct

Junbo Zhang is with the School of Information Science and Technology, Southwest Jiaotong University, Chengdu 610031, China & the Department of Computer Science, Georgia State University, Atlanta, GA 30303, USA (e-mail: JunboZhang86@163.com, jbzhang@cs.gsu.edu);

Jian-Syuan Wong and Yi Pan are with the Department of Computer Science, Georgia State University, Atlanta, GA 30303, USA (e-mail: jwong9@student.gsu.edu, pan@cs.gsu.edu);

Tianrui Li is with the School of Information Science and Technology, Southwest Jiaotong University, Chengdu 610031, China (e-mail: trli@swjtu.edu.cn).

*Corresponding author.

the relation matrix is designed for fast computing approximations. The key feature of this method is a Sub-Merge operation, which reduces the space requirement and accelerate the process of merging the relation matrices.

- 3) An incremental method is applied to the process of merging the relation matrices. With this feature, the relation matrix is updated in parallel and incrementally to efficiently accelerate the computational process.
- 4) A sparse matrix method is employed to optimize the proposed matrix-based method and further improve the performance of the algorithm.

The rest of the paper is organized as follows: we provide the basic concepts of the MapReduce model and rough sets under IIS in Section 2. We introduce the Matrix representation of the approximations in IIS in Section 3. Three different parallel matrix-based methods are proposed in the incomplete decision table (IDT) in Section 4. Experimental analysis is given in Section 5. The paper ends with conclusions and further research topics in Section 6.

2 PRELIMINARIES

The basic concepts, notations and results of the MapReduce model [1], [2] and rough sets as well as their extensions are briefly reviewed [8], [23], [24], [25], [26], [27].

2.1 The MapReduce Model and Twister

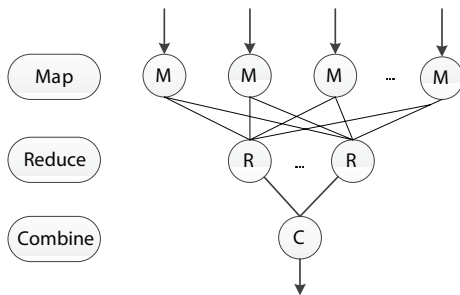


Fig. 1. The MapReduce framework in Twister

This section reviews some basic ideas of MapReduce and relative features of Twister. The MapReduce model, by Google, is designed to handle large-scale data sets in a distributed environment and has become a well-known parallel model in cloud computing [2]. It provides a convenient way to parallelize data analysis. The advantages of MapReduce include convenience, robustness, and scalability. Figure 1 shows the MapReduce framework in Twister. The computation in MapReduce model includes two major phases, the Map phase and the Reduce phase. Each phase utilizes key/value pairs as input and output. The programmer can choose types of keys and values depending on the requirements. In addition, the Map

and Reduce functions are also required for the Map phase and the Reduce phase, respectively. Input data is taken by the master node and divided into sub-program for the Map phase. A Map function can use a key/value pair $(K1, V1)$, and also generate the output in another key/value pair $(K2, V2)$ format. The intermediate results then are collected by the runtime system and used to produce the key/list values $(K2, \text{a list of } (V2))$ pairs as the input of a Reduce function. Once the Reduce function finishes its process, another key/list values $(K3, \text{a list of } (V3))$ pair is produced as the output.

After Google's work, various MapReduce platforms were developed. (1) Apache Hadoop [3] is one of the most popular run-time systems based on the MapReduce model; (2) Twister [4] is another run-time system especially designed for iterative MapReduce applications; (3) Phoenix [28] is a shared-memory implementation of MapReduce model for data-intensive processing tasks; (4) Mars [29] is a MapReduce runtime system for graphic processors (GPUs). In this paper, the proposed algorithms are all implemented with Twister, which is reviewed as follows.

Twister, a lightweight MapReduce runtime system, is designed for facilitating iterative applications based on MapReduce model. In [4], Twister also demonstrates that it is a good fit for several iterative scientific applications such as Blast. Besides the support of iterative MapReduce, Twister also has other attractive features that are different to typical non-iterative MapReduce runtime systems: I. Intermediated results are stored in the distributed memory of worker role which is more efficient than writing and reading intermediate data in the local disk; II. Input file can be partitioned into several files as the input data of the Map phase; III. The publish/subscribe messaging infrastructure is applied as the network broker to make communications become more efficient. IV. The *Combine* phase in Twister is used to collect results from all the reduce outputs and there is only one result generated in the *Combine* phase.

2.2 Rough sets

Definition 1: [8] An information system is defined as a pair $IS = (U, A)$ of non-empty, finite set U and A , where U is a *universe of objects*, and A is a set consisting of *attributes*, i.e., functions $a : U \rightarrow V_a$, where V_a is the set of values of attribute a , called the *domain* of a . Sometimes we distinguish in an information system $IS = (U, A)$ a partition of A into two classes $C, D \subseteq A$ of attributes, called *condition* and *decision* attributes, respectively. $DT = (U, C \cup D)$ is called a decision table.

Definition 2: [8] Let $B \subseteq A$ be a subset of attributes. The indiscernibility relation, denoted by I_B , is an equivalence relation defined as:

$$\forall x, y, I_B(x, y) \Leftrightarrow \forall a \in B : a(x) = a(y) \quad (1)$$

where $a(x)$ denotes the value of attribute a of object x .

The classical rough set analysis depends on the indiscernibility relation that describes indistinguishability of objects. It is based on the equivalence relation. The objects in an equivalence class satisfy reflexive, symmetric and transitive. In addition, it is not able to handle missing data, and requires the information system should be complete. However, missing data appears frequently in real applications [24]. The classical rough set approach, based on complete information systems, cannot be directly applied in information systems with missing attribute values. To process missing data, some extended rough set models [24], [26], [27] were presented.

Definition 3: [24] An information system (U, A) is called as an incomplete information system (IIS) if there exists a in A and x in U that satisfy that the value $a(x)$ is unknown, denoted as “*”, which is considered as an “everything is possible” value.

Under this definition of IIS, the toleration and similarity relation are proposed respectively to deal with unknown data in [24], [26].

Definition 4: [24] Let $B \subseteq A$ be a subset of attributes. The tolerance relation, denoted by T_B , is defined as:

$$T_B = \{(x, y) | \forall a \in B, (a(x) = a(y)) \vee (a(y) = *) \vee (a(x) = *)\}. \quad (2)$$

Obviously, the tolerance relation T is reflexive and symmetric, but not transitive.

Definition 5: [26] Let $B \subseteq A$ be a subset of attributes. The similarity relation, denoted by S_B , is defined as:

$$S_B = \{(x, y) | \forall a \in B, (a(x) \neq *) \wedge (a(x) = a(y))\}. \quad (3)$$

Obviously, the similarity relation S is reflexive and transitive, but not symmetric.

In many cases, the information system is incomplete due to the following two reasons [27]. One is that the value of an attribute is lost for a specific case. As an example, the value is currently unavailable although it was known originally owing to a variety of reasons, e.g., it was recorded but erased later. Another reason is the value of an attribute is irrelevant or unimportant. For example, it is feasible to diagnose a patient in spite of the fact that some clinical test results are not taken. Such missing attribute values do not matter for the final outcome and are called as “do not care” conditions [27]. Under this assumption that missing values in IIS are lost or “do not care” condition, a new version of the IIS is given as follows:

Definition 6: [27] $IIS = (U, A)$ is an incomplete information system if there exists a in A and x in U that satisfy that the value $a(x)$ is missing. All the missing values are denoted by “?” or “*”, where the lost value is denoted by “?”, “do not care” condition is denoted by “*”.

Under this new definition of IIS, Grzymala-Busse integrated the tolerance relation [24] and the similarity relation [26] and put forward a characteristic relation [27].

Definition 7: [27] In the incomplete information system $IIS = (U, A)$, for $B \subseteq A$, the characteristic relation K_B is defined as follows:

$$K_B = \{(x, y) | \forall a \in B, (a(x) \neq ?) \wedge (a(x) = a(y) \vee a(x) = * \vee a(y) = *)\}. \quad (4)$$

where “?” and “*” are missing values which mean the lost value and “do not care”, respectively.

When $(x, y) \in K_B$, we call x and y are indiscernible w.r.t. B . Let $K_B(x) = \{y | y \in U, (x, y) \in K_B\}$. We call $K_B(x)$ the characteristic class for x w.r.t. K_B .

The characteristic relation is reflexive, but not symmetric and transitive. It is a generalization of tolerance relation and similarity relation.

Table 1 shows a description of above binary relations.

TABLE 1
A description of relations

Relation	Reflexive	Symmetric	Transitive
Equivalence	✓	✓	✓
Tolerance	✓	✓	×
Similarity	✓	×	✓
Characteristic	✓	×	×

Definition 8: [23] Given an incomplete information system $IIS = (U, A)$, $\forall X \subseteq U$, $B \subseteq A$, the lower and upper approximations of X in terms of the characteristic relation K_B are defined as

$$\underline{K_B}(X) = \{x \in U | K_B(x) \subseteq X\} \quad (5)$$

$$\overline{K_B}(X) = \{x \in U | K_B(x) \cap X \neq \emptyset\} \quad (6)$$

Here, these two approximations divide the universe U into three disjoint regions: the positive region $POS_{K_B}(X)$, the boundary region $BND_{K_B}(X)$ and the negative region $NEG_{K_B}(X)$, respectively.

$$\begin{cases} POS_{K_B}(X) = \underline{K_B}(X) \\ BND_{K_B}(X) = \overline{K_B}(X) - \underline{K_B}(X) \\ NEG_{K_B}(X) = U - \overline{K_B}(X) \end{cases} \quad (7)$$

Definition 9: [23] Given an incomplete decision table $IDT = (U, C \cup D)$, $B \subseteq C$. Let $U/D = \{D_1, D_2, \dots, D_r\}$ be a partition over the decision D . Then the lower and upper approximations of the decision D with respect to attributes B are defined as

$$\underline{K_B}(D) = \bigcup_{j=1}^r \underline{K_B}(D_j) \quad (8)$$

$$\overline{K_B}(D) = \bigcup_{j=1}^r \overline{K_B}(D_j) \quad (9)$$

where

$$\underline{K_B}(D_j) = \{x \in U | K_B(x) \subseteq D_j\} \quad (10)$$

$$\overline{K_B}(D_j) = \{x \in U | K_B(x) \cap D_j \neq \emptyset\} \quad (11)$$

The positive region $POS_{K_B}(D)$, the boundary region $BND_{K_B}(D)$ and the negative region $NEG_{K_B}(D)$ of D with respect to the attribute subset B are defined as

$$\begin{cases} POS_{K_B}(D) = \underline{K_B}(D) \\ BND_{K_B}(D) = \overline{K_B}(D) - \underline{K_B}(D) \\ NEG_{K_B}(D) = U - \overline{K_B}(D) \end{cases} \quad (12)$$

3 MATRIX REPRESENTATION OF THE APPROXIMATIONS IN INCOMPLETE INFORMATION SYSTEMS

We here introduce the matrix representation of the lower and upper approximations in IIS.

From the matrix point of view, a set of axioms were constructed to characterize classical rough set upper approximation [30]. In our previous work, we defined a basic vector $H(X)$, which was induced from the relation matrix. Then four cut matrices of $H(X)$, denoted by $H^{[\mu,\nu]}(X)$, $H^{(\mu,\nu]}(X)$, $H^{[\mu,\nu)}(X)$ and $H^{(\mu,\nu)}(X)$, were derived for the approximations, positive, boundary and negative regions intuitively in set-valued information systems [31]. In this paper, to process large-scale incomplete data effectively, we also introduce the matrix-based methods.

First, we review the vector representation of the subset X of the domain U .

Definition 10: [31] Let $U = \{x_1, x_2, \dots, x_n\}$, and X be a subset of U . The characteristic function $G(X) = (g_1, g_2, \dots, g_n)^T$ (T denotes the transpose operation) is defined as

$$g_i = \begin{cases} 1, & x_i \in X \\ 0, & x_i \notin X \end{cases} \quad (13)$$

where $G(X)$ assigns 1 to an element that belongs to X and 0 to an element that does not belong to X . For example, if $U = \{x_1, x_2, x_3\}$, $X = \{x_1, x_3\}$, then we have $G(X) = \{1, 0, 1\}$.

Therefore, we use the characteristic relation to construct the relation matrix.

Definition 11: [32] Given an incomplete information system $IIS = (U, A)$. Let $B \subseteq A$ and K_B be a characteristic relation on U , $RM_{n \times n}^{K_B} = (m_{ij})_{n \times n}$ be an $n \times n$ matrix representing K_B , called the relation matrix w.r.t. B . Then

$$m_{ij} = \begin{cases} 1, & (x_i, x_j) \in K_B \\ 0, & (x_i, x_j) \notin K_B \end{cases} \quad (14)$$

According the above definitions, we obtain the following corollary easily.

Corollary 1: Let $RM_{n \times n}^{K_B} = (m_{ij})_{n \times n}$ and K_B be a characteristic relation on U . Then $m_{ii} = 1$, $1 \leq i, j \leq n$.

The induced diagonal matrix of the relation matrix is given for the later calculation of the basic vector.

Definition 12: [31] Let $B \subseteq A$, K_B be a characteristic relation on U , and $\Lambda_{n \times n}^{K_B}$ be an induced diagonal matrix of $RM = (m_{ij})_{n \times n}$. Then

$$\begin{aligned} \Lambda_{n \times n}^{K_B} &= \text{diag}\left(\frac{1}{\lambda_1}, \frac{1}{\lambda_2}, \dots, \frac{1}{\lambda_n}\right) \\ &= \text{diag}\left(\frac{1}{\sum_{j=1}^n m_{1j}}, \frac{1}{\sum_{j=1}^n m_{2j}}, \dots, \frac{1}{\sum_{j=1}^n m_{nj}}\right) \end{aligned} \quad (15)$$

where $\lambda_i = \sum_{j=1}^n m_{ij}$, $1 \leq i \leq n$.

Corollary 2: $\Lambda_{n \times n}^{K_B} = \text{diag}\left(\frac{1}{|K_B(x_1)|}, \frac{1}{|K_B(x_2)|}, \dots, \frac{1}{|K_B(x_n)|}\right)$ and $1 \leq |K_B(x_i)| \leq n$, $1 \leq i \leq n$.

To give the vector representation of rough set approximations, the positive region, the boundary region and the negative region, we need to introduce the basic vector as well as its four cut matrices.

Definition 13: [31] The n -column vector called a basic vector, denoted by $H(X)$, is defined as:

$$H(X) = \Lambda_{n \times n}^{K_B} \bullet M_{n \times n}^{K_B} \bullet G(X) \quad (16)$$

where \bullet is the dot product of matrices.

Definition 14: [31] Let $0 \leq \mu \leq \nu \leq 1$. Four cut matrices of $H(X)$, denoted by $H^{[\mu,\nu]}(X)$, $H^{(\mu,\nu]}(X)$, $H^{[\mu,\nu)}(X)$ and $H^{(\mu,\nu)}(X)$, are defined as follows.

(1) $H^{[\mu,\nu]}(X) = (h'_i)_{n \times 1}$

$$h'_i = \begin{cases} 1, & \mu \leq h_i \leq \nu \\ 0, & \text{else} \end{cases} \quad (17)$$

(2) $H^{(\mu,\nu]}(X) = (h'_i)_{n \times 1}$

$$h'_i = \begin{cases} 1, & \mu < h_i \leq \nu \\ 0, & \text{else} \end{cases} \quad (18)$$

(3) $H^{[\mu,\nu)}(X) = (h'_i)_{n \times 1}$

$$h'_i = \begin{cases} 1, & \mu \leq h_i < \nu \\ 0, & \text{else} \end{cases} \quad (19)$$

(4) $H^{(\mu,\nu)}(X) = (h'_i)_{n \times 1}$

$$h'_i = \begin{cases} 1, & \mu < h_i < \nu \\ 0, & \text{else} \end{cases} \quad (20)$$

Remark 1: These four cut matrices are Boolean matrices.

According to the cut matrices of the basic vector, we present the vector representation of rough set approximations as well as the positive region, the boundary region and the negative region.

Corollary 3: Given any subset $X \subseteq U$ in an incomplete information system $IIS = (U, A)$, where $U = \{x_1, x_2, \dots, x_n\}$. $B \subseteq A$ and K_B is a characteristic relation on U . $H(X) = (h_1, h_2, \dots, h_n)^T$ is the basic vector. Then the lower and upper approximations of X in IIS can be computed from the cut matrix of $H(X)$ as follows.

(1) The n -column boolean vector $G(\underline{K_B}(X))$ of the lower approximation $\underline{K_B}(X)$:

$$G(\underline{K_B}(X)) = H^{[1,1]}(X) \quad (21)$$

(2) The n -column boolean vector $G(\overline{K_B}(X))$ of the upper approximation $\overline{K_B}(X)$:

$$G(\overline{K_B}(X)) = H^{(0,1]}(X) \quad (22)$$

Corollary 4: The positive region $POS_{K_B}(X)$, the boundary region $BND_{K_B}(X)$, and the negative region $NEG_{K_B}(X)$ can also be generated from the cut matrix of $H(X)$, respectively as follows.

(1) The n -column boolean vector $G(POS_{K_B}(X))$ of the positive region:

$$G(POS_{K_B}(X)) = H^{[1,1]}(X) \quad (23)$$

(2) The n -column boolean vector $G(BND_{K_B}(X))$ of the boundary region:

$$G(BND_{K_B}(X)) = H^{(0,1)}(X) \quad (24)$$

(3) The n -column boolean vector $G(NEG_{K_B}(X))$ of the negative region:

$$G(NEG_{K_B}(X)) = H^{[0,0]}(X) \quad (25)$$

3.1 Matrix-based method in the incomplete decision table

We have introduced the vector representation of rough set approximations of the subset X of the domain U . However, in IDT, the domain U is partitioned by the decision D . Therefore, we need to give the representation of rough set approximations of all subsets of the domain U . That is the matrix representation of rough set approximations.

Definition 15: Given an incomplete decision table $IDT = (U, C \cup D)$, $B \subseteq C$, let $U/D = \{D_1, D_2, \dots, D_r\}$ be a partition over the decision D . $\forall D_j \in U/D$, $G(D_j)$ is an n -column boolean vector of D_j . Let $DM_{n \times r} = (G(D_1), G(D_2), \dots, G(D_r))$ be an $n \times r$ boolean matrix, called a decision matrix.

Similarly, we combine all the basic vectors as the basic matrix and give its computational method as follows.

Corollary 5: Given an incomplete decision table $IDT = (U, C \cup D)$, $B \subseteq C$. Let $U/D = \{D_1, D_2, \dots, D_r\}$ be a partition over the decision D . $\forall D_j \in U/D$, $H(D_j)$ is the basic vector. Let $HD = (H(D_1), H(D_2), \dots, H(D_r))$ be an $n \times r$ matrix, called a basic matrix. Then, HD can be computed as follows.

$$HD = \Lambda_{n \times n}^{K_B} \bullet RM_{n \times n}^{K_B} \bullet DM_{n \times r} \quad (26)$$

where \bullet is the dot product of matrices.

Therefore, we present the vector representation of rough set approximations of each subset as well as the positive region, the boundary region and the negative region of the decision D .

Corollary 6: Given an incomplete decision table $IDT = (U, C \cup D)$, $B \subseteq C$. Let $U/D = \{D_1, D_2, \dots, D_r\}$ be a partition over the decision D . $HD = (H(D_1), H(D_2), \dots, H(D_r))$ is a basic matrix. $\forall j = 1, 2, \dots, r$, the upper and lower approximations of D_j in IIS can be computed from the cut matrix of HD as follows.

(1) The n -column boolean vector $G(\underline{K_B}(D_j))$ of the lower approximation $\underline{K_B}(D_j)$:

$$G(\underline{K_B}(D_j)) = H^{[1,1]}(D_j) \quad (27)$$

(2) The n -column boolean vector $G(\overline{K_B}(D_j))$ of the upper approximation $\overline{K_B}(D_j)$:

$$G(\overline{K_B}(D_j)) = H^{(0,1)}(D_j) \quad (28)$$

The positive region $POS_{K_B}(D)$, the boundary region $BND_{K_B}(D)$, and the negative region $NEG_{K_B}(D)$ can also be generated from the cut matrix of HD , respectively as follows.

(3) The n -column boolean vector $G(POS_{K_B}(D))$ of the positive region:

$$G(POS_{K_B}(D)) = \sum_{j=1}^r H^{[1,1]}(D_j) \quad (29)$$

where \sum is the boolean addition (logical operation OR) of boolean vectors.

(4) The n -column boolean vector $G(BND_{K_B}(D))$ of the boundary region:

$$\begin{aligned} G(BND_{K_B}(D)) &= \sum_{j=1}^r H^{(0,1)}(D_j) - \sum_{j=1}^r H^{[1,1]}(D_j) \\ &= \sum_{j=1}^r H^{(0,1)}(D_j) \end{aligned} \quad (30)$$

(5) The n -column boolean vector $G(NEG_{K_B}(D))$ of the negative region:

$$\begin{aligned} G(NEG_{K_B}(D)) &= J - \sum_{j=1}^r H^{(0,1)}(D_j) \\ &= \left(\sum_{j=1}^r H^{(0,1)}(D_j) \right)^{[0,0]} \end{aligned} \quad (31)$$

where $J = (1, 1, \dots, 1)^T$ is an $n \times 1$ matrix of ones.

3.2 Sequential Matrix-based Algorithm (SMA) in incomplete decision table

From the view of the matrix, according to Definitions 11, 15 and Corollaries 5, 6, we present a sequential method for computing rough set approximations in IDT and its corresponding naive Sequential Matrix-based Algorithm (SMA) (Figure 2 and Algorithm 1).

Algorithm SMA (see Algorithm 1) is a sequential matrix-based algorithm for computing approximations of the decision, the positive region, the boundary region and the negative region. Step 2 is to construct the relation matrix according to Definition 11, whose time complexity is $O(|U|^2|C|)$. By Definition 12, Step 3 is to compute the induced diagonal matrix. Its time complexity is $O(|U|^2)$. Step 4 is to construct the decision matrix according to Definition 15, which costs $O(|U||D|)$. Step 5 is to compute the basic matrix HD by Corollary 5 with time complexity $O(|U|^2r)$. Finally, Step 6 is to compute and output the approximations of the decision, the positive region, the boundary region and the negative region according to Corollary 6 with time complexity $O(|U|r)$. Therefore, the total complexity of Algorithm 1 is $O(|U|^2(|C| + r))$. In many real-applications, e.g., gene selection [33] which includes thousands of condition attributes, there is only one decision attribute and a few decision classes. It means $|C| \gg |D|$ and $|C| \gg r$. Thus the most intensive computation to occur is the construction of the relation matrix. In following sections, we will discuss how to construct the relation matrix in parallel.

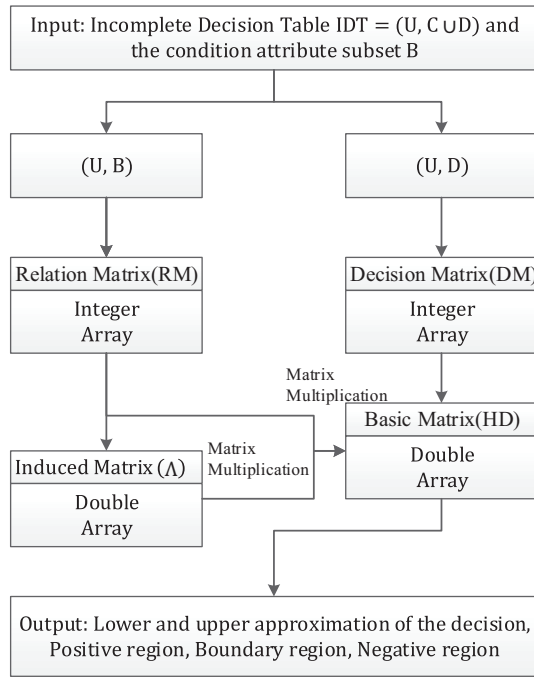


Fig. 2. The sequential method for computing rough set approximations in IDT. In order to facilitate the operations of matrix multiplication, we use integer arrays to construct the relation matrix and the decision matrix in our implementation.

Algorithm 1: A naive Sequential Matrix-based Algorithm (SMA) for computation of rough set approximations

Input:

The incomplete decision table $IDT = (U, C \cup D)$ and the condition attribute subset $B \subseteq C$.

Output:

The approximations of the decision: $\underline{K}_B(D_i)$ and $\overline{K}_B(D_i)$, $i = 1, 2, \dots, r$; the positive region $POS_{K_B}(D)$, the boundary region $BND_{K_B}(D)$, and the negative region $NEG_{K_B}(D)$.

1 **begin**

2 **Construct** the relation matrix $RM_{n \times n}^{K_B} = (m_{ij})_{n \times n}$ in terms of the condition attributes B by using the characteristic relation K_B ;

3 **Compute** the induced diagonal matrix

$$\Lambda_{n \times n}^{K_B} = \text{diag} \left(\frac{1}{\sum_{j=1}^n m_{1j}}, \frac{1}{\sum_{j=1}^n m_{2j}}, \dots, \frac{1}{\sum_{j=1}^n m_{nj}} \right);$$

4 **Construct** the decision matrix

$DM_{n \times r} = (G(D_1), G(D_2), \dots, G(D_r))$ in terms of the decision attribute D ;

5 **Compute** the basic matrix

$$HD = \Lambda_{n \times n}^{K_B} \bullet RM_{n \times n}^{K_B} \bullet DM_{n \times r};$$

6 **Compute and Output** the approximations of the decision, the positive region, the boundary region and the negative region by using the cut matrices of HD .

7 **end**

4 PARALLEL MATRIX-BASED METHODS IN INCOMPLETE DECISION TABLE

As mentioned above, the most intensive computation to occur is the construction of the relation matrix. To accelerate the process of the construction of the relation matrix, we use the MapReduce framework to parallelize this process.

4.1 S1: A parallel strategy based on MapReduce

In this subsection, we give a simple parallel strategy based on MapReduce, called S1.

We first give the definition of the partition of IDT with regard to the attribute.

Definition 16: Given an incomplete decision table $IDT = (U, C \cup D)$. Let $IDT = \bigcup_{t=1}^{|C|} S_t \cup S_D$, where $S_t = (U, c_t)$ and $S_D = (U, D)$. It satisfies $C = \bigcup_{t=1}^{|C|} c_t$. It means that IDT is divided into $|C| + 1$ sub-decision tables.

Remark 2: The sub-decision table S_D is used to construct the decision matrix in sequential according to Definition 15.

In order to show the merge operation of the matrices, we need to introduce the logical conjunction of two matrices firstly.

Definition 17: Given two $\mu \times \nu$ boolean matrices $Y = (y_{ij})_{\mu \times \nu}$ and $Z = (z_{ij})_{\mu \times \nu}$, the logical conjunction \wedge of two matrices is defined as

$$Y \wedge Z = (y_{ij} \wedge z_{ij})_{\mu \times \nu} \quad (32)$$

Actually, we partition IDT with regard to the attribute. Each attribute in IDT can generate a relation matrix. Here, we give the following lemma to show how to merge two relation matrices.

Lemma 1: Let $b_1, b_2 \in C$ and $K_{\{b_1, b_2\}}$, K_{b_1} , K_{b_2} be the characteristic relation on U . Then

$$RM_{n \times n}^{K_{\{b_1, b_2\}}} = RM_{n \times n}^{K_{b_1}} \wedge RM_{n \times n}^{K_{b_2}} \quad (33)$$

Proof. Suppose $RM_{n \times n}^{K_{\{b_1, b_2\}}} = (m_{ij})_{n \times n}$, $RM_{n \times n}^{K_{b_1}} = (\varphi_{ij}^{b_1})_{n \times n}$ and $RM_{n \times n}^{K_{b_2}} = (\varphi_{ij}^{b_2})_{n \times n}$.

(a) If $(x_i, x_j) \in K_{\{b_1, b_2\}}$, then $m_{ij} = 1$, that is, x_i and x_j are indiscernible w.r.t. $\{b_1, b_2\}$. By Definition 7, we have $x_i K_{b_1} x_j$ and $x_i K_{b_2} x_j$. Therefore, $\varphi_{ij}^{b_1} = 1$ and $\varphi_{ij}^{b_2} = 1$. Hence, $m_{ij} = \varphi_{ij}^{b_1} \wedge \varphi_{ij}^{b_2} = 1$.

(b) If $(x_i, x_j) \notin K_{\{b_1, b_2\}}$, then $m_{ij} = 0$, that is, x_i and x_j are discernible w.r.t. $\{b_1, b_2\}$. By Definition 7, we have $(x_i, x_j) \notin K_{b_1}$ or $(x_i, x_j) \notin K_{b_2}$. Therefore, $\varphi_{ij}^{b_1} = 0$ or $\varphi_{ij}^{b_2} = 0$. Hence, $m_{ij} = \varphi_{ij}^{b_1} \wedge \varphi_{ij}^{b_2} = 0$.

To sum up, $\forall i, j \in \{1, 2, \dots, n\}$, $m_{ij} = \varphi_{ij}^{b_1} \wedge \varphi_{ij}^{b_2}$, namely, $m_{n \times n}^{K_{\{b_1, b_2\}}} = m_{n \times n}^{K_{b_1}} \wedge m_{n \times n}^{K_{b_2}}$. \square

Similarly, multiple relation matrices generated by different attributes can be merged according to the following corollary.

Corollary 7: Let $B \subseteq C$, K_B be the characteristic relation on U . Then

$$RM_{n \times n}^{K_B} = \bigwedge_{b \in B} RM_{n \times n}^{K_b} \quad (34)$$

According to Lemma 1, each condition attribute will generate a sub-relation matrix. For $IDT = (U, C \cup D)$, it will cost the space $|U|^2|C|$ to construct the relation matrix. As an example, there exists a data set with 10000 objects and 1000 attributes. It means the required memory space for storing an intermediate matrix of an attribute is about 95 MB ($10000^2/(2^{20}) \approx 95.36$). If we apply it to analyze this data, then about 93 GB memory space ($1000 \text{ attributes} \times 95 \text{ MB}$) is required to save all the intermediate results. To overcome this weakness, we introduce a local merge operation, called Sub-Merge, which can be viewed as a local Reduce. We continue the above example. Assume that there are 8 compute nodes, then there are $1000/8 = 125$ attributes in each node. Each node processes these 125 attributes one by one. Before passing intermediate matrices to the Reduce phase, the proposed Sub-Merge operation merges these intermediate matrices. It means that each node only generates one intermediate matrix with the Sub-Merge, but generates 125 intermediate matrices without the Sub-Merge. With the feature of the Sub-Merge, only 760 MB memory space ($8 \times 95 \text{ MB}$) is required to save all the intermediate results. Moreover, for single node, it only needs 95 MB to store. Without the Sub-Merge operation, each node needs 11875 MB memory space ($125 \times 95 \text{ MB}$), which is always greater than the memory of the node. To verify this viewpoint, we have conducted the experiment to test the algorithm without the Sub-Merge operation. Unfortunately, we always receive the error message “Out of Memory” because it needs more memory space to store. Besides, the relation matrix is asymmetric only when processing the missing data “?” according to Definition 7. It means that the missing data “?” can be processed separately when there are only a few “?”. In this situation, it almost can reduce half of space. In our later implementation, we only use the (a)symmetry property to construct the relation matrix but still store a full matrix because its performance improvement may be less significant when the memory space is enough and there are many missing data “?”.

The Sub-Merge operation is based on the merge of relation matrices of different attribute sets. Therefore, we give a new definition of the partition of IDT with regard to the attribute set.

Definition 18: Given an incomplete decision table $IDT = (U, C \cup D)$. Let $IDT = \bigcup_{t=1}^f S_t \cup S_D$, where $S_t = (U, C_t)$ and $S_D = (U, D)$. It satisfies: (1) $C = \bigcup_{t=1}^f C_t$; (2) $C_j \cap C_k = \emptyset, \forall j, k \in \{1, 2, \dots, f\}$ and $j \neq k$. It means that IDT is divided into $f + 1$ sub-decision tables.

Be different from Definition 16, this definition shows that each sub-decision table includes an attribute set rather than a single condition attribute.

Lemma 2: Let $B_1, B_2 \subseteq C, B_1 \cap B_2 = \emptyset$ and $K_{\{B_1 \cup B_2\}}$, where K_{B_1}, K_{B_2} are two characteristic relations on U . Then

$$RM_{n \times n}^{K_{\{B_1 \cup B_2\}}} = RM_{n \times n}^{K_{B_1}} \wedge RM_{n \times n}^{K_{B_2}} \quad (35)$$

Proof. The proof is similar to that of Lemma 1. \square

Corollary 8: Let $B = \bigcup_{t=1}^f B_t \subseteq C, K_B$ be the characteristic relation on U . Then

$$RM_{n \times n}^{K_B} = \bigwedge_{t=1}^f RM_{n \times n}^{K_{B_t}} \quad (36)$$

According to Lemma 2, each condition attribute sub-set can generate a sub-relation matrix independently. And, the relation matrix can be constructed through the logical conjunction of sub-relation matrices according to Corollary 8. It is easy to know that the element m_{ij} in the relation matrix are independent. Thus, they can be computed in parallel. Figure 3 gives the parallel strategy based on MapReduce. Intuitively, we also present the parallel algorithm (see Algorithm 2). Assume that the input data are split and stored in different nodes. In each Map, the sub-relation matrix is constructed. Its size is $n \times n$, where $n = |U|$. Simply, the row number of sub-relation matrix is set as the key of output in the Map phase. The key domain is $\{1, 2, \dots, n\}$. It means that there are n tasks in the Reduce phase. Preferably, we can use a parameter η to control the number of tasks in the Reduce phase. Namely, the key domain is $\{1, 2, \dots, \eta\}$. Specifically, (1) If $\eta = 1$, it is the coarsest level of granularity. Thus there is only one task in the Reduce phase, which processes $n \times n$ elements. (2) If $\eta = n \times n$, it is the finest level of granularity. Thus there are $n \times n$ tasks in the Reduce phase, each of which processes only one element. The Map and Reduce functions in S1 are outlined in Algorithms 3 and 4, respectively. Specially, the process Sub-Merge in the Map function is used to merge sub-relation matrices. After the Reduce phase, we use *Combine* to merge the relation vectors and construct the relation matrix.

Algorithm 2: A Parallel Matrix-based Algorithm (PMA) for computation of rough set approximations

- 1 **Call** the MapReduce module to construct the relation matrix;
 - 2 **Construct** the decision matrix;
 - 3 **Split** the relation matrix into a set of row blocks and the decision matrix into a set of column blocks, and **Send** the data blocks to different nodes;
 - 4 **Call** the MapReduce based matrix operations module to compute the basic matrix;
 - 5 **Output** the approximations of the decision, the positive region, the boundary region and the negative region;
-

The only sequential step is the construction of the decision matrix, which is same to that decision matrix of Algorithm 1. As for the remaining matrix multiplications, we employ the MapReduce based matrix multiplication module that has been implemented in the runtime system Twister [4]. Before it, we need to split the relation matrix and the decision matrix. It is based on

Algorithm 3: S1-Map(*key*, *value*)

Input:

//*key*: The condition attribute B_t

//*value*: The sub-decision table $S_t = (U, B_t)$

Output:

key': Row number in sub-relation matrix

value': A byte array

```

1 begin
2   Let byte  $RM_{n \times n}^{K_{B_t}} = \{1\}$ ; //  $n \times n$  matrix
3   for each  $b \in B_t$  do
4     Construct the sub-relation matrix  $RM_{n \times n}^{K_b}$ ;
5      $RM_{n \times n}^{K_{B_t}} = RM_{n \times n}^{K_{B_t}} \wedge RM_{n \times n}^{K_b}$ ;
6   end
7   for  $i = 1$  to  $n$  do
8      $key' = i$ ;
9      $value = RM_{n \times n}^{K_{B_t}}(i)$ ; // the  $i$ th row of  $RM_{n \times n}^{K_{B_t}}$ 
10    output.collect( $key'$ ,  $value'$ );
11  end
12 end

```

Algorithm 4: S1-Reduce(*key*, *value*)

Input:

//*key*: Row number in sub-relation matrix

//*value*: A list of byte array V

Output:

key': Row number in sub-relation matrix

value': A byte array

```

1 begin
2    $key' = key$ ;
3   Let byte  $value' = \{1\}$ ; //  $1 \times n$  vector
4   for each  $v \in V$  do
5      $value' = value' \& v$ ; //  $\&$  is And-operation
6     of vectors
7   end
8   output.collect( $key'$ ,  $value'$ );
9 end

```

the column/row decomposition strategy¹. Be different to the original matrix multiplication in Twister, we combine two matrix multiplications as one that can accelerate the process of the computation. According to Corollary 5, given the relation matrix $RM_{n \times n}^{K_B} = (m_{ik})_{n \times n}$ and the decision matrix $DM = (d_{kj})_{n \times r}$, one can write the basic matrix with the alternative form as

$$\begin{aligned}
 HD &= \Lambda_{n \times n}^{K_B} \bullet (RM_{n \times n}^{K_B} \bullet DM_{n \times r}) \\
 &= \text{diag} \left(\frac{1}{\sum_{k=1}^n m_{1k}}, \frac{1}{\sum_{k=1}^n m_{2k}}, \dots, \frac{1}{\sum_{k=1}^n m_{nk}} \right) \bullet
 \end{aligned}$$

1. http://www.iterativemapreduce.org/samples.html#Matrix_Multiplication

$$\begin{aligned}
 &\begin{pmatrix} \sum_{k=1}^n m_{1k}d_{k1} & \dots & \sum_{k=1}^n m_{1k}d_{kr} \\ \sum_{k=1}^n m_{2k}d_{k1} & \dots & \sum_{k=1}^n m_{2k}d_{kr} \\ \vdots & \ddots & \vdots \\ \sum_{k=1}^n m_{nk}d_{k1} & \dots & \sum_{k=1}^n m_{nk}d_{kr} \end{pmatrix} \\
 &= \begin{pmatrix} \frac{\sum_{k=1}^n m_{1k}d_{k1}}{\sum_{k=1}^n m_{1k}} & \dots & \frac{\sum_{k=1}^n m_{1k}d_{kr}}{\sum_{k=1}^n m_{1k}} \\ \frac{\sum_{k=1}^n m_{2k}d_{k1}}{\sum_{k=1}^n m_{2k}} & \dots & \frac{\sum_{k=1}^n m_{2k}d_{kr}}{\sum_{k=1}^n m_{2k}} \\ \vdots & \ddots & \vdots \\ \frac{\sum_{k=1}^n m_{nk}d_{k1}}{\sum_{k=1}^n m_{nk}} & \dots & \frac{\sum_{k=1}^n m_{nk}d_{kr}}{\sum_{k=1}^n m_{nk}} \end{pmatrix}
 \end{aligned}$$

Therefore, after finishing the computations of the relation matrix and the decision matrix, the basic matrix HD can be computed directly. The alternative form also shows that the one-step matrix operations for computing HD has the similar form with the matrix multiplication. The column/row decomposition strategy can also be used here.

4.2 S2: An incremental parallel strategy based on MapReduce

To accelerate the process of computation of rough set approximations, Zhang et al. proposed an incremental algorithm to update the relation matrix when adding or deleting attributes in set-valued information systems [31]. In S1, we know that each condition attribute generates a sub-relation matrix each time and they are merged in the Sub-Merge of Map phase and the Reduce phase. The process of merging can be viewed as a process of adding attributes one by one. It is a typical incremental process. Based on these analysis, we present the incremental parallel strategy based on MapReduce, called S2.

Corollary 9: [31] Let $P, Q \subseteq C$ and $Q \cap P = \emptyset$. Suppose $M_{n \times n}^{K_{P \cup Q}} = (m_{ij}^\uparrow)_{n \times n}$ is a relation matrix representing $K_{P \cup Q}$, and $M_{n \times n}^{K_P} = (m_{ij})_{n \times n}$ is a relation matrix representing K_P . The relation matrix by adding Q to P can be updated as:

$$m_{ij}^\uparrow = \begin{cases} 0, & m_{ij} = 0 \vee (x_i, x_j) \notin K_Q \\ 1, & m_{ij} = 1 \wedge (x_i, x_j) \in K_Q \end{cases} \quad (37)$$

Remark 3: When adding Q to P , m_{ij} is constant while $m_{ij}^\uparrow = 0$.

Here, the flow chart of S2 is same to S1. The difference between S2 and S1 is that we use an incremental technique to accelerate the construction of the relation matrix in the Map phase by Corollary 9 as outlined in Algorithm 5. The other steps are same to those in S1.

4.3 S3: An incremental parallel strategy using the sparse matrix based on MapReduce

A sparse matrix is a matrix populated primarily with zeros as elements of the table [34]. The native data structure for a matrix is a two-dimensional array, which is used in S1 and S2. Each entry in the array represents an element m_{ij} of the matrix and can be accessed by the two indices i and j . In tradition, i represents the row number (top-to-bottom), while j represents the column number (left-to-right) of each element in the table. For an $m \times n$ matrix,

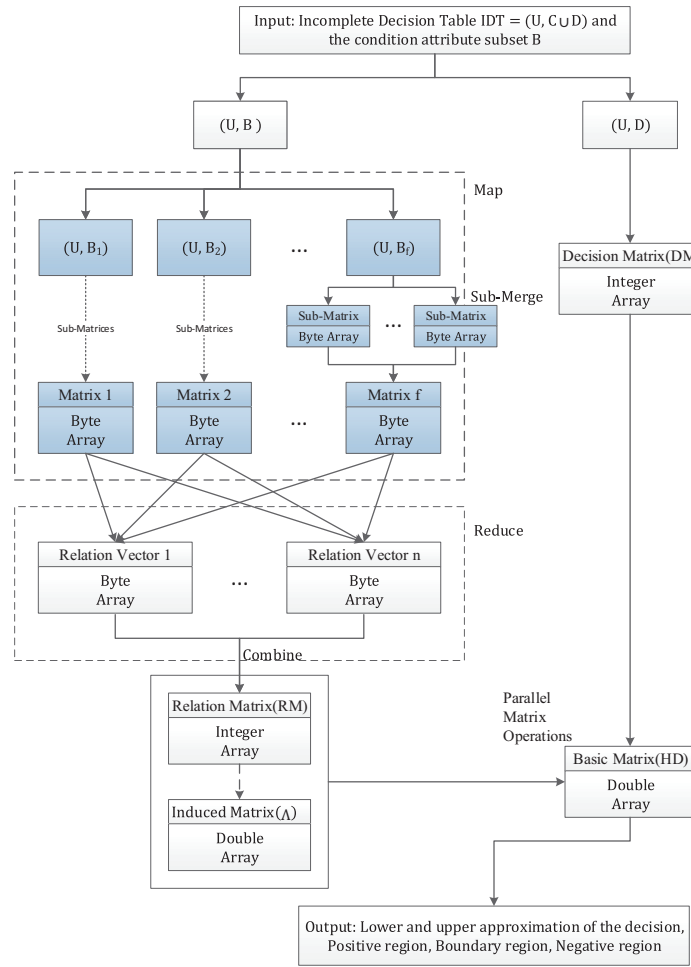


Fig. 3. A parallel strategy based on MapReduce. To reduce space complexity, we use byte arrays to storage the sub-relation matrices.

enough memory to store up to $(m \times n)$ entries to represent the matrix is required. As the number of condition attributes increases, there are more and more zero entries in the relation matrix. Reduction of substantial memory requirement can be realized by storing only the non-zero entries. Depending on the number and distribution of the non-zero entries, different data structures can be used and yield huge savings in memory when compared to the native method.

Storing formats of the sparse matrix can be divided into two groups: (1) those that support efficient matrix operations and (2) those that support efficient modification. The incremental construction of the relation matrix in S2 shows the most operations are matrix update (modification). Thus, we select the latter format. The efficient modification group includes List of lists (LIL), Dictionary of keys (DOK) and Coordinate list (COO), which is typically used to construct the matrix. In our implementation, we use LIL to construct the relation matrix in parallel. LIL stores one list per row, where each entry stores a column index as the value. Typically, these entries are sorted by the column index for faster lookup. This is another format which is good for incremental

matrix construction. Since our relation matrix is $(0, 1)$ matrix, we use a simplified LIL that also stores one list per row but each entry only stores a column index. We here define the sparse degree of the relation matrix as follows.

Definition 19: Given a relation matrix $RM_{n \times n}^{K_B} = (m_{ij})_{n \times n}$. K_B is a characteristic relation on U . The sparse degree of the relation matrix, denoted by SD_B , is defined as

$$SD_B = 1 - \frac{\sum_{i=1}^n \sum_{j=1}^n m_{ij}}{n \times n} \quad (38)$$

Lemma 3: Let $P \subseteq Q \subseteq C$, $M_{n \times n}^{K_P} = (m_{ij})_{n \times n}$ and $M_{n \times n}^{K_Q} = (m'_{ij})_{n \times n}$ be relation matrices representing K_P and K_Q , respectively. Their sparse degrees are SD_P and SD_Q , respectively. Then

$$SD_P \leq SD_Q \quad (39)$$

Proof. According to Corollary 9, when adding attributes, m_{ij} is constant while $m'_{ij} = 0$. Since $P \subseteq Q$, we have $m_{ij} \geq m'_{ij}$, $1 \leq i, j \leq n$. Hence, $\sum_{i=1}^n \sum_{j=1}^n m_{ij} \geq \sum_{i=1}^n \sum_{j=1}^n m'_{ij}$.

Algorithm 5: S2-Map(*key*, *value*)

Input:
//*key*: The condition attribute B_t
//*value*: The sub-decision table $S_t = (U, B_t)$
Output:
key': Row number in sub-relation matrix
value': A byte array

```

1 begin
2   Let byte [][]  $RM_{n \times n}^{K_{B_t}} = \{m_{ij}\} = 1$ ; //  $n \times n$  matrix
3   for each  $b \in B_t$  do
4     for  $i = 1$  to  $n$  do
5       for  $j = 1$  to  $n$  do
6         if  $m_{ij} == 1$  and  $(x_i, x_j) \notin K_b$  then
7            $m_{ij} = 0$ ;
8         end
9       end
10    end
11    for  $i = 1$  to  $n$  do
12       $key' = i$ ;
13       $value = RM_{n \times n}^{K_{B_t}}(i)$ ; // the  $i$ th row of  $RM_{n \times n}^{K_{B_t}}$ 
14      output.collect( $key'$ ,  $value'$ );
15 end
```

Therefore, $1 - \frac{\sum_{i=1}^n \sum_{j=1}^n m_{ij}}{n \times n} \leq 1 - \frac{\sum_{i=1}^n \sum_{j=1}^n m'_{ij}}{n \times n}$, in other words, $SD_P \leq SD_Q$. \square

Remark 4: When there are more attributes in IIS, the sparse degree of the relation matrix would be higher, in other words, the relation matrix would be sparser.

According to these observations, we propose the incremental parallel strategy using the sparse matrix based on MapReduce, called S3. We here use LIL to redesign our Map and Reduce function, as outlined in Algorithms 6 and 7, respectively. The other steps are similar to those in S1. The most difference is that we use the sparse matrix to replace the matrix and some matrix operations, *e.g.*, matrix multiplication, are replaced as the sparse matrix operations. The column/row decomposition strategy is also used here.

5 EXPERIMENTAL ANALYSIS

Our experiments run on the representative MapReduce runtime system Twister [4]. Its version is 0.9. And we code our algorithms (a pure sequential version and a parallel Twister version) in JAVA and compile the programs with JDK1.7.0_05. The experiments on Twister are carried out on Cheetah cluster [35] using 9 nodes, which contains 1 head node and 8 compute nodes. 4 compute nodes have 16 GB main memory and use AMD Opteron(TM) Processor 2376 CPUs (8 cores in all, each has a clock frequency of 2.3 GHz). Another 4 nodes have 8 GB main memory and use Intel Xeon CPU E5410, comprising two Quad-Core CPUs (8 cores in all, each

Algorithm 6: S3-Map(*key*, *value*)

Input:
//*key*: The condition attribute B_t
//*value*: The sub-decision table $S_t = (U, B_t)$
Output:
key': Row number in sub-relation matrix
value': A list

```

1 begin
2   for  $i = 1$  to  $n$  do
3     Let  $LIL(i) = Null$ ; // LIL is list of list,
4     // LIL(i) is the  $i$ th list in LIL
5   end
6   Select an attribute  $b'$  from  $B$ ;
7   for  $i = 1$  to  $n$  do
8     for  $j = 1$  to  $n$  do
9       if  $(x_i, x_j) \in K_{b'}$  then  $LIL(i).add(j)$ ;
10    end
11    for each  $b \in \{B_t - b'\}$  do
12      for  $i = 1$  to  $n$  do
13        for  $j = 1$  to  $n$  do
14          if  $(x_i, x_j) \notin K_b$  then
15             $LIL(i).remove(j)$ ;
16          end
17        end
18      end
19      for  $i = 1$  to  $n$  do
20         $key' = i$ ;
21         $value = LIL(i)$ ;
22        output.collect( $key'$ ,  $value'$ );
23 end
```

Algorithm 7: S3-Reduce(*key*, *value*)

Input:
//*key*: Row number in sub-relation matrix
//*value*: A list of list V
Output:
key': Row number in sub-relation matrix
value': A list

```

1 begin
2    $key' = key$ ;
3   Let  $value' = ALL = \{1, 2, \dots, n\}$ ;
4   for each  $v \in V$  do
5      $value' = value' \cap v$ ; //  $\cap$  is the intersection
6     // operation of two lists
7   end
8   output.collect( $key'$ ,  $value'$ );
9 end
```

has a clock frequency of 2.33 GHz). We here compare the performance using 1, 2, 4, 8, 16 and 32 cores².

We download two data sets from machine learning data repository, University of California at Irvine [36]. They are “Balance Scale” and “Mushroom”. In data set “Balance Scale”, there are 625 samples, 1 decision attribute and 4 condition attributes. Another data set “Mushroom” consists of 8124 samples. Each sample consists of 1 decision attribute and 22 condition attributes, which includes missing values. We only test the performance in this paper. Hence, we set half of missing values as “*”, and the other half as “?”. To test the performance of parallel methods, the data sets are used as input many times to stimulate the large input condition. Take the data set “Mushroom” as an example, we keep the number of sample constant and input the condition attributes repeatedly. If we repeat 32 times, it would have 704 (22×32) attributes. We call this data set S8k-A0.7k. All simulated data sets are outlined in Table 2.

TABLE 2
A description of data sets

Data Sets	Rename	Samples	Attributes	Classes	Size
Balance Scale	S625-A8k	625	8192	3	16MB
	S625-A82k	625	81920	3	158MB
	S625-A819k	625	819200	3	1.37GB
	S625-A4096k	625	4096000	3	5.43GB
Mushroom	S8k-A0.7k	8124	704	2	11MB
	S8k-A2.8k	8124	2816	2	45MB
	S8k-A23k	8124	22528	2	357MB
	S8k-A45k	8124	45056	2	715MB

Note: S625-A8k means 625 samples and about 8k attributes.

5.1 Performance

The aforementioned three parallel strategies (see Section 4) are executed on the MapReduce runtime system Twister.

In what follows, we first test the performance of Strategy 1, which is a parallel method without the incremental technique. We examine its speedup characteristics. It is defined by the following formula [37]:

$$Speedup(p) = \frac{T_1}{T_p} \quad (40)$$

where p is the number of cores, T_1 is the execution time on one core, and T_p is the execution time on p cores.

The linear speedup is a system with p times the number of computers yields a speedup of p . It is ideal but difficult to achieve because the communication cost increases when the number of clusters becomes large.

We perform the execution time and the speedup evaluation on data sets with various sizes and structures. The number of cores varies from 1 to 32. According to the parallel method of Section 4.1, the only sequential

portion in Strategy 1 is the construction of the decision matrix. Therefore, we recorded the execution time of the construction of the decision matrix and other steps separately when testing the pure sequential version. Then, the expected speedup number can be obtained by using Amdahl’s Law with the following formula [38]:

$$Expected\ Speedup(p) = \frac{1}{R_S + \frac{R_P}{p}} \quad (41)$$

where $R_S + R_P = 1$, R_S represents the ratio of the sequential portion in the program.

Figures 4 and 5 show the execution time and the speedup for all these data sets. The expected speedup of S625-A8k (Balance Scale) and S8k-A0.7k (Mushroom) are shown in Figure 4(b) and Figure 5(b), respectively. Both are very close to the linear speedup because the only sequential portion (the construction of the decision matrix) can be computed efficiently. The expected speedup of the other three simulated data sets of Balance Scale and Mushroom are not shown here because it is obvious that they are better than S625-A8k and S8k-A0.7k, respectively. As the results show, the proposed parallel algorithms have a great speedup performance. As the size of the data set increases, the speedup becomes better. For example, the size of the data set S625-A4096k is 5.43 GB. Its speedup on 32 core is greater than 18. Therefore, the parallel Strategy 1 can process large-scale data efficiently.

5.2 A comparison of different parallel strategies

In this subsection, we show a comparison of different parallel strategies. We here only compare the execution time of different parallel strategies rather than speedup because Strategies 2 and 3 (S2, in Section 4.2 and S3, in Section 4.3) use an incremental technique to accelerate the update process, whose time complexity depends on the characteristics of the testing data. Figure 6 shows the execution time of Strategies 1, 2 and 3 versus the number of cores. It is obvious that S2 and S3 always have the better performance than S1. That is because we utilize the incremental method in S2 and S3 to facilitate the Sub-Merge operation. With the incremental method, further processes of merge (update) operation are only executed while the values in temporary result matrix are 1. Instead, it always executes $n \times n$ (n is the numbers of objects) operations in S1. Therefore, many redundant processes can be skipped in S2 and S3.

In S3, we use the sparse matrix to optimize the storing space, which is different from the data structure used in S1 and S2, since the relation matrix constructed from large data is usually sparse. Thus, we implement the sparse relation matrix with the data structure LIL in S3, which allows us to change the size of a sparse matrix dynamically, and the space requirement can be decreased as well.

From Figure 6, even though the execution time using S3 is faster than that using S2 in most conditions, S2

2. The single core is used to test the performance of the pure sequential version.

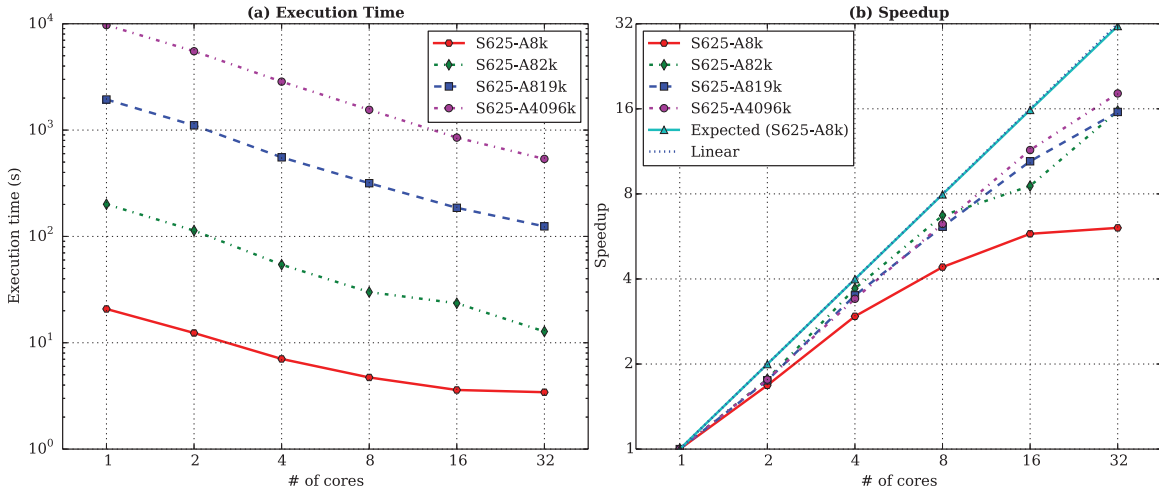


Fig. 4. Execution Time and Speedup plots using parallel Strategy 1 on simulated data sets of Balance Scale.

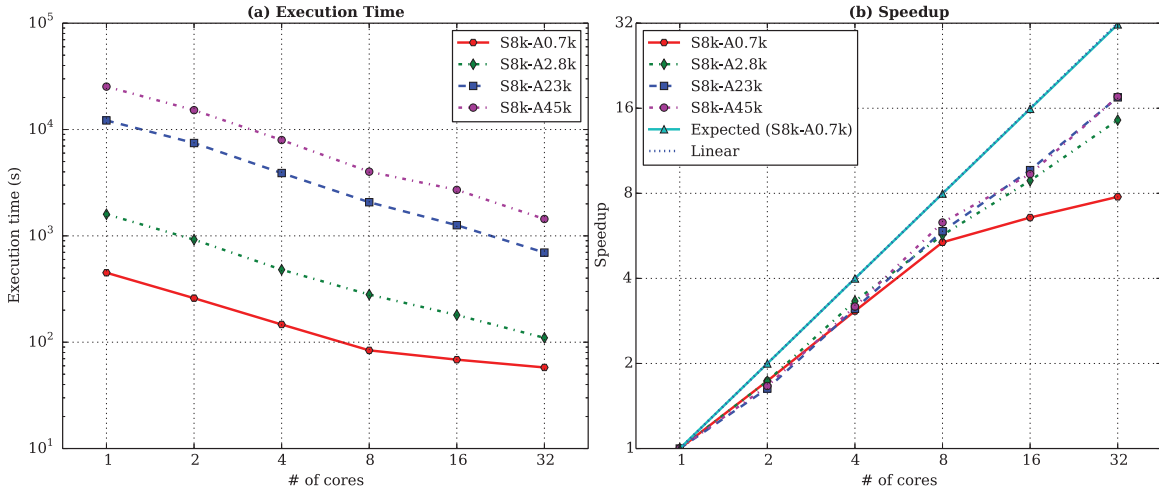


Fig. 5. Execution Time and Speedup plots using parallel Strategy 1 on simulated data sets of Mushroom.

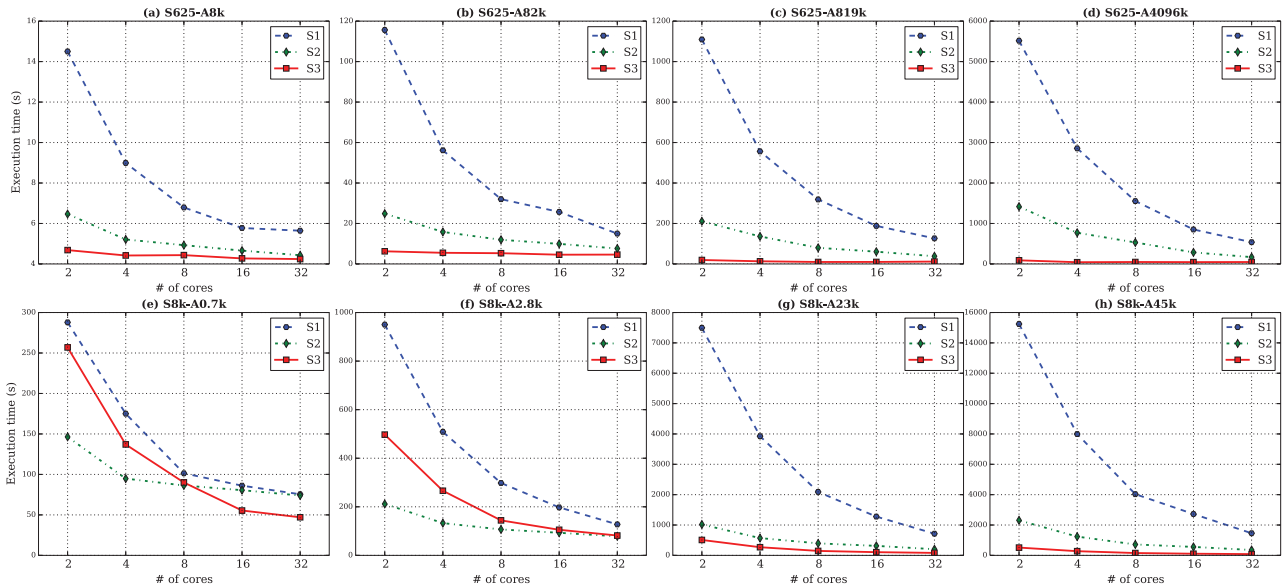


Fig. 6. Execution time of different parallel strategies versus the number of cores. S1: Strategy 1; S2: Strategy 2; S3: Strategy 3.

still has the better performance than S3 while the testing data contains huge numbers of objects and only a few attributes. Since the matrix structure is used in S2, a traverse process in an $n \times n$ relation matrix requires $n \times n$ elements to be examined. However, LIL structure traversal in S3 only needs to look over the elements remain in the structure. Owing to the feature of LIL structure, the sparse degree of the relation matrix has a significant effect on the performance of S3 which means the higher sparse degree of the LIL structure has, the fewer elements need to be traversed. In this condition, S3 performs better than S2. But, in some cases, S2 is superior to S3 such as Figure 6(f). Two possible reasons are given as follows: 1) S3 takes much more time to initialize the matrix while the size of objects increases significantly; 2) The size of generated intermediate matrix will not be shrunk significantly because of only few attributes using S3 which causes the communication overhead.

Based on these, we have the following qualitative results: 1) S3 performs better than S2 when the sparse degree of the relation matrix is higher that may cause by a large number of attributes; 2) S2 performs better than S3 when the sparse degree of the relation matrix is lower. In real applications, we can choose the incremental strategies S2 or S3 to accelerate processes of execution according to different sparse degrees.

6 CONCLUSIONS

In this paper, to process large-scale incomplete data with rough set theory, we first introduced the matrix representation of lower and upper approximations in IIS. According to the characteristics of the matrix, we proposed three parallel strategies based on MapReduce to compute approximations. All of them were implemented on the MapReduce runtime system Twister. The results of testing speedup of S1 show that the parallel strategy S1 has a good performance on processing large-scale data. In examining an incremental method to accelerate the processing procedures in S2 and S3, we found that S2 and S3 always have better performance than S1, and, in most cases, S3 outperforms S2. But, if the input data contains huge quantities of objects and only a few attributes, S2 would do better than S3. Our extensive experimental evaluation demonstrates that our proposed parallel methods are more efficient in analyzing data with large amount of attributes. As we can see, computing lower and upper approximations is the key to rule induction and feature selection when utilizing rough set-based methods. We plan in our future work to further investigate the knowledge discovered and features selection from large-scale, incomplete data.

ACKNOWLEDGMENTS

We would like to thank Alexander Ororbia and all anonymous reviewers for their constructive comments, which helped improve the presentation of this paper.

This work is supported by the National Science Foundation of China (Nos. 60873108, 61175047, 61100117) and NSAF (No. U1230117), the Fundamental Research Funds for the Central Universities (No. SWJTU11ZT08), the Research Fund of Traction Power State Key Laboratory, Southwest Jiaotong University (No. 2012TPL_T15), the Science and Technology Planning Project of Sichuan Province (No. 2012RZ0009), China and the Fostering Foundation for the Excellent Ph.D. Dissertation of Southwest Jiaotong University (No. 2012ZJB), China.

REFERENCES

- [1] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," in *Proceedings of the 6th conference on Symposium on Operating Systems Design & Implementation - Volume 6*, ser. OSDI'04. Berkeley, CA, USA: USENIX Association, 2004, pp. 10–10.
- [2] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, Jan. 2008.
- [3] T. White, *Hadoop: The Definitive Guide*, 2nd ed. O'Reilly Media / Yahoo Press, 2010.
- [4] J. Ekanayake, H. Li, B. Zhang, T. Gunarathne, S.-H. Bae, J. Qiu, and G. Fox, "Twister: a runtime for iterative mapreduce," in *Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing*, ser. HPDC '10. New York, NY, USA: ACM, 2010, pp. 810–818.
- [5] Amazon Elastic MapReduce, <http://aws.amazon.com/elasticmapreduce/>.
- [6] Twister4Azure, <http://salsahpc.indiana.edu/twister4azure/index.html>.
- [7] Z. Pawlak, *Rough Sets: Theoretical Aspects of Reasoning about Data, System Theory, Knowledge Engineering and Problem Solving*. Dordrecht: Kluwer Academic Publishers, 1991, vol. 9.
- [8] Z. Pawlak and A. Skowron, "Rudiments of rough sets," *Information Sciences*, vol. 177, no. 1, pp. 3–27, 2007.
- [9] Z. Pawlak and A. Skowron, "Rough sets: Some extensions," *Information Sciences*, vol. 177, no. 1, pp. 28–40, 2007.
- [10] Z. Pawlak and A. Skowron, "Rough sets and boolean reasoning," *Information Sciences*, vol. 177, no. 1, pp. 41–73, 2007.
- [11] J.Y. Liang, F. Wang, C.Y. Dang, and Y.H. Qian, "A group incremental approach to feature selection applying rough set technique," *Knowledge and Data Engineering, IEEE Transactions on*, 2012 (Accepted).
- [12] J. W. Grzymala-Busse and W. Ziarko, "Data mining and rough set theory," *Communications of the ACM*, vol. 43, no. 4, pp. 108–109, Apr. 2000.
- [13] W. Ziarko, "Discovery through rough set theory," *Communications of the ACM*, vol. 42, no. 11, pp. 54–57, Nov. 1999.
- [14] Q.H. Hu, W. Pedrycz, D.R. Yu, and J. Lang, "Selecting discrete and continuous features based on neighborhood decision error minimization," *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, vol. 40, no. 1, pp. 137–150, feb. 2010.
- [15] Q.H. Hu, Z.X. Xie, and D.R. Yu, "Hybrid attribute reduction based on a novel fuzzy-rough model and information granulation," *Pattern Recognition*, vol. 40, no. 12, pp. 3509–3521, Dec. 2007.
- [16] Q.H. Hu, D.R. Yu, J.F. Liu, and C.X. Wu, "Neighborhood rough set based heterogeneous feature subset selection," *Information Sciences*, vol. 178, no. 18, pp. 3577–3594, Sep. 2008.
- [17] Y.H. Qian, C.Y. Dang, J.Y. Liang, and D. Tang, "Set-valued ordered information systems," *Information Sciences*, vol. 179, no. 16, pp. 2809–2832, Jul. 2009.
- [18] Y.H. Qian, J.Y. Liang, W. Pedrycz, and C.Y. Dang, "An efficient accelerator for attribute reduction from incomplete data in rough set framework," *Pattern Recogn.*, vol. 44, no. 8, pp. 1658–1670, Aug. 2011. Available:
- [19] Y.H. Qian, J.Y. Liang, W. Pedrycz, and C.Y. Dang, "Positive approximation: An accelerator for attribute reduction in rough set theory," *Artificial Intelligence*, vol. 174, no. 9–10, pp. 597–618, Jun. 2010.

- [20] J. Zhang, T. Li, D. Ruan, and D. Liu, "Neighborhood rough sets for dynamic data mining," *International Journal of Intelligent Systems*, vol. 27, no. 4, pp. 317–342, 2012.
- [21] J. Zhang, T. Li, D. Ruan, Z. Gao, and C. Zhao, "A parallel method for computing rough set approximations," *Information Sciences*, vol. 194, no. 0, pp. 209–223, Jul. 2012.
- [22] J. Zhang, T. Li, and Y. Pan, "Parallel rough set based knowledge acquisition using mapreduce from big data," in *Proceedings of the 1st International Workshop on Big Data, Streams and Heterogeneous Source Mining: Algorithms, Systems, Programming Models and Applications*, ser. BigMine '12. New York, NY, USA: ACM, 2012, pp. 20–27.
- [23] T. Li, D. Ruan, W. Geert, J. Song, and Y. Xu, "A rough sets based characteristic relation approach for dynamic attribute generalization in data mining," *Knowledge-Based Systems*, vol. 20, no. 5, pp. 485–494, Jun. 2007.
- [24] M. Kryszkiewicz, "Rough set approach to incomplete information systems," *Information Sciences*, vol. 112, no. 1–4, pp. 39–49, Dec. 1998.
- [25] M. Kryszkiewicz, "Rules in incomplete information systems," *Information Sciences*, vol. 113, no. 3–4, pp. 271–292, Feb. 1999.
- [26] J. Stefanowski and A. Tsoukias, "On the extension of rough sets under incomplete information," in *Proceedings of the 7th International Workshop on New Directions in Rough Sets, Data Mining, and Granular-Soft Computing*, ser. RSFDGrC '99. London, UK, UK: Springer-Verlag, 1999, pp. 73–81.
- [27] J. Grzymala-Busse, "Characteristic relations for incomplete data: A generalization of the indiscernibility relation," in *Transactions on Rough Sets IV*, ser. Lecture Notes in Computer Science, J. Peters and A. Skowron, Eds. Springer Berlin / Heidelberg, 2005, vol. 3700, pp. 58–68.
- [28] C. Ranger, R. Raghuraman, A. Penmetsa, G. Bradski, and C. Kozyrakis, "Evaluating mapreduce for multi-core and multi-processor systems," in *Proceedings of the 2007 IEEE 13th International Symposium on High Performance Computer Architecture*, ser. HPCA '07. Washington, DC, USA: IEEE Computer Society, 2007, pp. 13–24.
- [29] B. He, W. Fang, Q. Luo, N. K. Govindaraju, and T. Wang, "Mars: a mapreduce framework on graphics processors," in *Proceedings of the 17th international conference on Parallel architectures and compilation techniques*, ser. PACT '08. New York, NY, USA: ACM, 2008, pp. 260–269.
- [30] G.L. Liu, "The axiomatization of the rough set upper approximation operations," *Fundamenta Informaticae*, vol. 69, no. 3, pp. 331–342, 2006.
- [31] J. Zhang, T. Li, D. Ruan, and D. Liu, "Rough sets based matrix approaches with dynamic attribute variation in set-valued information systems," *International Journal of Approximate Reasoning*, 2012.
- [32] G.L. Liu, "Axiomatic systems for rough sets and fuzzy rough sets," *International Journal of Approximate Reasoning*, vol. 48, no. 3, pp. 857–867, Aug. 2008.
- [33] M.-L. Hou, S.-L. Wang, X.-L. Li, and Y.-K. Lei, "Neighborhood rough set reduction-based gene selection and prioritization for gene expression profile analysis and molecular cancer classification," *Journal of Biomedicine and Biotechnology*, vol. 2010, 2010.
- [34] J. Stoer and R. Bulirsch, *Introduction to Numerical Analysis*, 3rd ed. Berlin, New York: Springer-Verlag, 2002.
- [35] *Cheetah cluster in Georgia State University*, <http://cs.gsu.edu/?q=cheetah>.
- [36] D. Newman, S. Hettich, C. Blake, and C. Merz, *UCI Repository of Machine Learning Databases*, University of California, Department of Information and Computer Science, Irvine, CA, 1998. < <http://archive.ics.uci.edu/ml/> >.
- [37] X. Xu, J. Jäger, and H.-P. Kriegel, "A fast parallel clustering algorithm for large spatial databases," *Data Mining and Knowledge Discovery*, vol. 3, pp. 263–290, 1999, 10.1023/A:1009884809343.
- [38] G. M. Amdahl, "Validity of the single processor approach to achieving large scale computing capabilities," in *Proceedings of the April 18–20, 1967, spring joint computer conference*, AFIPS '67. New York, NY, USA: ACM, 1967, pp. 483–485.



Junbo Zhang received his B.Eng. degree from Southwest Jiaotong University, China, in 2009. He is currently a Ph.D. candidate at the School of Information Science and Technology, Southwest Jiaotong University, China, supervised by Prof. Tianrui Li. He is also a visiting Ph.D. student at the Department of Computer Science, Georgia State University, USA, under advisor Prof. Yi Pan. He is interested in data mining, cloud computing, granular computing and rough sets. He is now a student member of ACM and CCF. Since 2009, he has published over 30 research papers in referred journals (e.g., *Information Sciences*, *International Journal of Approximate Reasoning*, *International Journal of Intelligent Systems*), and conferences.



Jian-Syuan Wong received his B.S. degree from National Sun Yat-Sen University, Taiwan, in 2008 and M.S degree from Georgia State University, USA in 2012. He is currently a Ph.D. student at the College of Information Sciences and Technology, Pennsylvania State University. He was a research assistant at the Morehouse School of Medicine, USA. His research interests include cloud computing, parallel and distributed systems, and data mining.



Yi Pan was born in Jiangsu, China. He entered Tsinghua University in March 1978 with the highest college entrance examination score among all 1977 high school graduates in Jiangsu Province. Currently, he is the chair and a professor in the Department of Computer Science and a professor in the Department of Computer Information Systems at Georgia State University. Dr. Pan received his B.Eng. and M.Eng. degrees in computer engineering from Tsinghua University, China, in 1982 and 1984, respectively, and his Ph.D. degree in computer science from the University of Pittsburgh, USA, in 1991.

Dr. Pan's research interests include parallel and cloud computing, wireless networks, and bioinformatics. Dr. Pan has published more than 140 journal papers with about 50 papers published in various IEEE/ACM journals. In addition, he has published over 150 conference papers and edited/authored 37 books. He has received many awards from organizations such as IEEE, NSF, AFOSR, JSPS, IBM, ISIBM, IISF and Mellon Foundation. Dr. Pan has served as an editor-in-chief or editorial board member for 15 journals including 6 IEEE Transactions and a guest editor for 10 special issues for 9 journals including 2 IEEE Transactions. He has organized numerous international conferences and workshops and has delivered over 20 keynote speeches at international conferences.

Dr. Pan is a "Great Master Face-to-Face" Series Speaker (2012), an IEEE Distinguished Speaker (2000–2002), a Yamacraw Distinguished Speaker (2002), a Shell Oil Colloquium Speaker (2002), and a senior member of IEEE. He is listed in Men of Achievement, Who's Who in Midwest, Who's Who in America, Who's Who in American Education, Who's Who in Computational Science and Engineering, and Who's Who of Asian Americans.



Tianrui Li received his B.S. degree, M.S. degree and Ph.D. degree from the Southwest Jiaotong University, China in 1992, 1995 and 2002 respectively. He was a Post-Doctoral Researcher at Belgian Nuclear Research Centre (SCK • CEN), Belgium from 2005-2006, a visiting professor at Hasselt University, Belgium in 2008 and the University of Technology, Sydney, Australia in 2009. And, he is presently a Professor and the Director of the Key Lab of Cloud Computing and Intelligent Technique of Sichuan Province,

Southwest Jiaotong University, China.

Since 2000, he has co-edited 4 textbooks, 12 proceedings, 6 special issues of international journal and published over 120 research papers in refereed journals and conferences. He is a senior member of IEEE, the chair of IEEE CIS Chengdu Chapter and the area editor of International Journal of Computational Intelligence Systems (IJCIS, SCI indexed). He has served as ISKE2007, ISKE2008, ISKE2009, ISKE2010, ISKE2011, JRS2012, ISKE2012, ISKE2013 program chairs, IEEE GrC 2009 program vice chair and RSKT2008, FLINS2010 organizing chairs and has been a reviewer for several leading academic journals.