

# Proficient Information Method for Inconsistency Detection in Multiple Data Sources

**B. Bazeer Ahamed<sup>1</sup> and T. Ramkumar<sup>2</sup>**

<sup>1</sup>Department of Computer Science & Engineering, Sathyabama University,

Chennai, India

bazeerahamed@gmail.com

<sup>2</sup>School of Information Technology and Engineering, VIT University,

Vellore - 632 014, India

ramooad@yahoo.com

## ABSTRACT

*With the form of new data sources on the Internet, integrating data from heterogeneous data repositories has become critical. However, multiple sources of data introduce problems such as redundancy, conflicts, or missing data reports. The two major categories of challenges for large scale data integration systems are heterogeneous data and conflicting data. Inaccurate results and poor decision making may occur during the integration process; during integration process the data is redundant and inconsistent. The solutions for heterogeneous data have been researched for many years, but the challenges of conflicting data are not well explored yet. We aim at improving the quality of information integration via data inconsistency detection method and information design process through experimental results.*

**Keywords:** heterogeneous data sources, information integration, inconsistency detection, semantic integration.

**Mathematics Subject Classification:** 68P15, 68P20

**Computing Classification System:** H.3.3

## 1. INTRODUCTION

Depending on the specific data mining tasks, multiple information sources can refer to data collected from different places, data of different formats, or data that capture different aspects of the objects. To provide better web services on automobile industry like car manufacturer, dealer in decision making, we need to conduct integrative analysis of all the information sources, we need to conduct integrative analysis of all the information sources. As rich heterogeneous data can be collected in nearly every industry, people began to recognize.

The importance of integrating multiple data sources in the field of data integration and data fusion, many algorithms have been developed to effectively integrate or combine raw data. Many studies focused on how to match the schemas of different data sources, detect entities that refer to the same real-world objects, answer queries by searching from multiple data sources, or combine multiple redundant information sources into one reliable source. In numerous applications that own multiple data sources, it is crucial not only to integrate or combine multiple data sources, but also consolidate different concepts for intelligent decision making. Therefore, in the field of knowledge integration, many algorithms have been developed to merge and synthesize models, rules, patterns obtained from multiple sources by reconciling their differences. These methods conduct classification or clustering from multiple data sources to identify more reliable and meaningful label predictions.

We propose a general framework to detect inconsistencies across multiple heterogeneous information sources, as well as approaches to find objects performing inconsistently in information networks and distributed systems. In today's large-scale distributed systems, the same type of information may be collected from each machine in the system. Although some knowledge can be extracted from each individual information source, a much richer body of knowledge can only be obtained by exploring the correlations or interactions across different sources. The effectiveness of this framework is demonstrated in its ability of detecting anomalous events and locating problematic sources from real monitoring data of three companies' infrastructures.

The rest of the paper is organized as follows. Section 2 reviews various processes of inconsistent data in heterogeneous data. Consistency examination process during integrated data is evaluated in section 3. Experimental measures of inconsistency detection is worked out in section 4. Section 5 concludes the paper and shows our future directions on this topic.

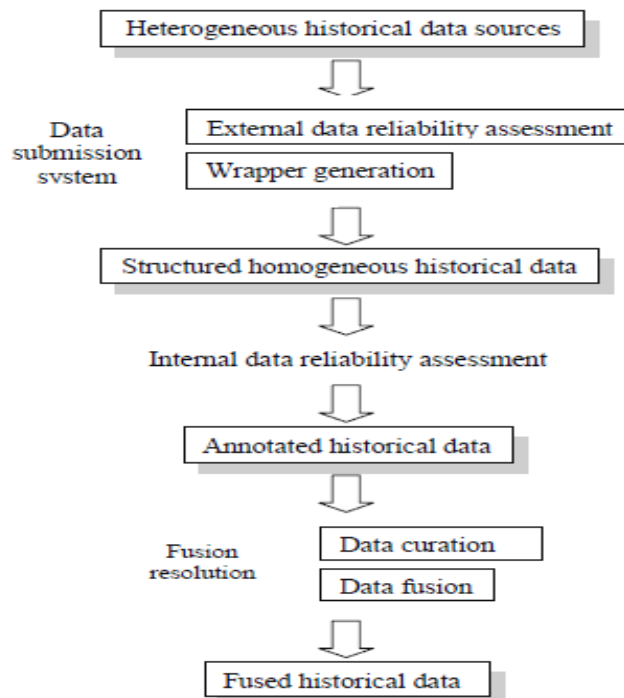
## 2. RELATED WORK

### 2.1. Efficient data fusion for heterogeneous data sources

Data integration from heterogeneous data sources requires a tremendous amount of work. The possible problems that users may encounter during data integration processes are inaccurate data, inconsistent data and redundant data. These problems are either caused by heterogeneous data sets or conflicting data sets (Gao et al., 2008). Heterogeneous data is defined as data stored in different

schemas or in different representations. Redundant data is defined as data stored in multiple databases with overlapping time, location, or name. These redundant data may result in inconsistency if the overlapping parts are inconsistent (i.e. temporal/spatial/naming inconsistency) (Zubko et al., 2009).

From the database point of view, data integration may be performed when there is heterogeneity at the schema level, tuple level, or value level. Information resulted from data integration process at different levels may have different representations, information types, and functionality, etc. Thus, when a designer starts to create a data integration system, the factors that needs to be considered includes the type of data, the algorithm of data merge and the level where the data integration process happens. A common approach to perform data integration involves the following steps: (1) identify the corresponding attributes in the sources, (2) differentiate objects that originate in different sources and if these data describe the same attributes, and (3) merge these sources into a single representation. Figure 1 illustrates the data integration architecture.



**Figure 1.** Data integration architecture in intelligent data centers.

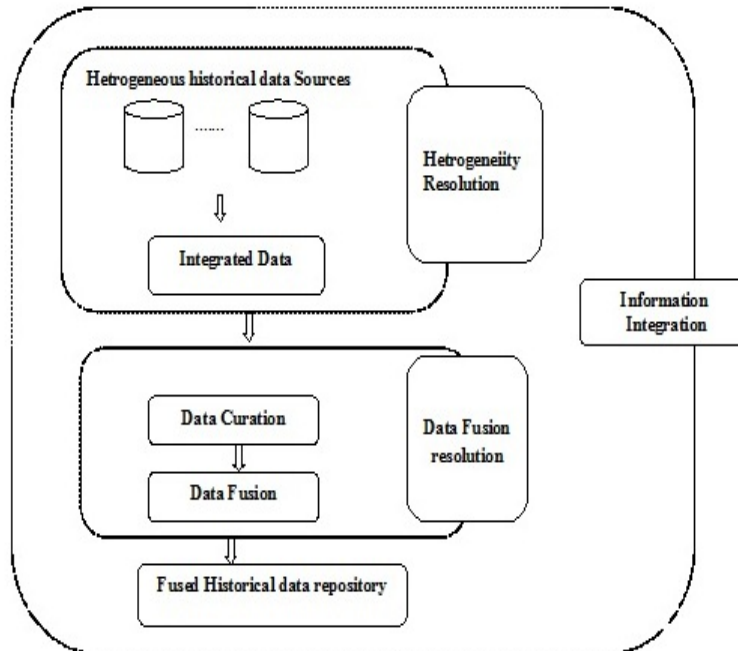
## 2.2. Combination of multiple heterogeneous models

The success of ensemble techniques has been proven theoretically and observed in real practice. However, the major challenge is that it is hard to obtain sufficient and reliable labelled data for effective training because they require the efforts of experienced human annotators. To tackle this challenge, the problem of consolidating multiple supervised and unsupervised information sources by negotiating their predictions to form a final superior classification solution (Ramachandran et al., 2009), (Jeh and Widom, 2002). The method can utilize all the available sources in the consensus

combination framework, no matter they contain labeled or unlabeled information. A global optimal label assignment for objects is derived by maximizing consensus among the given models. This consensus maximization approach crosses the boundary between supervised and unsupervised learning, and its effectiveness has been shown in many real-world problems, where the classification accuracy is significantly improved. In particular, the proposed method has been used to solve the following problems (Su et al., 2011):

- 1) user-generated video classification based on video, audio, and text features (Gao et al., 2011),
- 2) decision fusions of heterogeneous sensor nodes in sensor networks (Su et al., 2011), and
- 3) combination of heterogeneous anomaly detectors to improve performance over botnet or network traffic anomaly detection in cyber-security areas (Ramachandran et al., 2009).

Figure 2 describes the information integration architecture.



**Figure 2.** Information integration architecture.

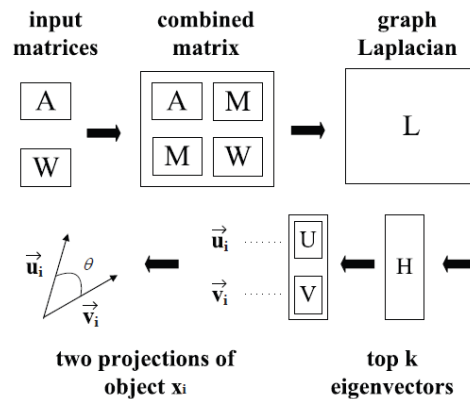
### 2.3. Inconsistency detection across multiple heterogeneous sources

Reaching consensus among heterogeneous information sources gives us the gains in classification performance. On the other hand, by exploring the differences among sources, identify something unusual and interesting, detect anomalies or inconsistencies across multiple information sources. Although the problem of anomaly detection has been widely studied (Zadorozhny et al., 2013), most of the existing approaches identify anomalies from one single data source. Different from existing work, we propose to identify inconsistencies across multiple information sources as a new type of meaningful anomalies in this part. In this sub-section, we define the general problem of inconsistency

detection across multiple heterogeneous sources and propose an effective spectral framework to identify such anomalies.

To detect an object that has inconsistent behavior among multiple heterogeneous sources (Zubko et al., 2010). A set of objects can be described from various perspectives (multiple information sources). The underlying clustering structure of normal objects is usually shared by multiple sources. However, anomalous objects belong to different clusters when considering different aspects. To identify such objects, computing the distance between different eigen decomposition results of the same object with respect to different sources as its anomalous score, and give interpretations from the perspectives of constrained spectral clustering and random walks over graph. Experimental results on several UCI as well as DBLP and MovieLens datasets demonstrate the effectiveness of the given approach (Bleiholder and Naumann, 2008).

The available information sources on the same set of objects have similar clustering structures, and thus if an object is assigned to different clusters when using various information sources, it can be regarded as a horizontal anomaly (Fan et al., 2005). This suggests that the first cluster the objects separately in each source and compare the clustering results. However, because clustering is unsupervised learning, we do not know the correspondence between clusters in different clustering solutions. To solve this problem by adding the constraint that the same object should be put into the same cluster by all the clustering solutions as often as possible. Figure 3 shows the horizontal anomaly detection.



**Figure 3.** Flow of horizontal anomaly detection.

The physical meaning of the proposed algorithm is explained from both constrained spectral clustering and random walk point of view. Experimental results show that the HOAD algorithm can consistently find horizontal anomalies from DBLP and MovieLens datasets (Zadorozhny and Lewis, 2013), where other anomaly detection methods fail to identify.

#### 2.4. Inconsistency detection for system debugging

This procedure is usually referred to as System Debugging. A distributed system consists of multiple connected machines, and monitoring data collected from each machine can be regarded as an information source (Dong and Li, 1999). Although some knowledge about system problems can be

extracted from each individual information source, a much richer body of knowledge can only be obtained by exploring the correlations or interactions across different sources (Gonzalez et al., 2010). Specially, the correlations between measurements collected across the distributed system can be used to infer normal system behavior, and thus a reasonable model to describe such correlations is crucially important in detecting and problems locating system. The effectiveness of this framework is demonstrated in its ability of detecting anomalous events and locating problematic sources from real monitoring data of three companies' infrastructures (Kang et al., 2011) The effectiveness and efficiency of the proposed method through experiments on a large collection of real monitoring data from three companies' infrastructure.

The feature space of monitoring data into grid cells and compute the transition probabilities among the cells adaptively according to the monitoring data (Lawson and Hanson, 1974). Compared with previous system monitoring techniques, the advantages of our approach include:

- 1) It detects the system problems considering both spatial and temporal information;
- 2) The model can output the problematic measurement ranges, which are useful for human debugging; and
- 3) The method is fast and can describe both linear and non-linear correlations.

Experiments on monitoring data collected from three real distributed systems involving 100 measurements from around 50 machines show the effectiveness in this method (Gao et al., 2008).

## **2.5. Data integration in heterogeneous database systems**

Data integration is the progression of extracting and integration data from multiple heterogeneous foundations to be loaded into an integrated in sequence supply. Solving structural, syntactical and semantic heterogeneities connecting source and intention data has been a composite problem for data integration for a number of years.

One solution to this problem has been urbanized through the use of a single global database schema that represents the integrated information with mappings from global schema to local schemas, where each query to the global schema is translated to queries to the local databases using these mappings.

The use of domain ontology, metadata, transformation rules, user, and system constraints have resolved the majority of the problems of domain mismatch associated with schematic integration and global schematic approaches.

However, even when all the mappings, semantic and structure heterogeneity are solved in the global schema, consistency may not have been achieved, because the information provided by the sources may be mutually inconsistent. This problem has remained because it is impossible to capture all the information and avoid null values. At the same time, each autonomous component database deals with its own properties or domain constraints on information, such as accuracy, reliability, availability, timeliness and cost of information access.

Several approaches to solve inconsistency between databases have been implemented, I. and II.:

I. By reconciliation of data, also known as data fusion: different values become just one using a fusion function (i.e. average, highest, and majority), depending on the data semantic.

II. On the basis of individual data properties: associated with each information source (i.e. cost of retrieving information, how recent is the information, level of authority associated with this source, or accuracy and completeness of information). These properties can be specified at different levels: the global schema design level, the query itself or in the users profile.

### 3. EXPERIMENTAL RESULTS

#### 3.1. Consistency examination process during integrated data

To provide inconsistency detection my model generates a system called typical substitution method. The goals of this method are to detect inconsistency occurrences and to provide proper values for each reported interval to moderate reduce of inconsistent data. When data sources are integrated, reports can be grouped in several systems depending on their related conditions. The unknown variable vector  $X$  represents unknown event density for each time:

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} \quad (1)$$

The size of vector  $X$  depends on the extended beyond condition of these process, in other words,  $n$  is different for every substitution method. The coefficient matrix  $M$  denotes:

$$M = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \quad (2)$$

The guide value of reports as a constraint value vector  $P$  is

$$P = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \quad (3)$$

Let us consider four reports from heterogeneous data sources of proceedings with chronological cover as shown in Table 1. When we position this process successively on the timeline by their occurrence time the timeline will be divided into smaller units of time intervals. Each interval must have a non-negative value shared by all the process. And the sum of matching intervals will be equal to the sum of the description values.

**Table 1.** Vehicle defect identification before inconsistency.

ID (Ri)	Vehicles	Manufacturer	From	To	Duration	Defective Value(v) out of 1000
R1	omni	Maruthi	1910	1980	70	700
R2	omni	Maruthi	1930	1980	50	500
R3	omni	Maruthi	1945	1985	40	600
R4	omni	Maruthi	1960	2000	40	700

These four ID's divide the timeline into common intervals. The number of intervals depends on the number of values and how they partly cover.

The above TSM in matrix form is

$$AX = P,$$

$$A = \begin{bmatrix} 111100 \\ 111110 \\ 000001 \\ 000111 \end{bmatrix}, X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, P = \begin{bmatrix} 700 \\ 500 \\ 600 \\ 700 \end{bmatrix} \quad (4)$$

where 1 represents the continuation (reporting) of a given value for a specific time intervals and 0 indicates that equivalent values that does not cover that interval. The TSM is a method used to determine the best solutions to maximize the profit or minimize the cost for a model that includes equations representing a list of constraint or requirements. In the proposed method, the objective is to minimize the difference between resolution sets and the real values for each interval, which can be referred to as limitations for these processes. Therefore, the following optimization problem is defined:

$$\begin{aligned} & \text{Max (orMin)} \\ & S_{1x1} + S_{2x2} + \dots + S_{nxn} \\ & a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \leq b_1 \\ & a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \leq b_2 \\ & \vdots \\ & a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \leq b_m \\ & x_1, x_2 \dots x_n \geq 0, b_1, b_2 \dots b_m \geq 0 \end{aligned} \quad (5)$$

where the TSM form uses  $S = [S_1, S_2 \dots S_n]$ , X expressed in (1), M expressed in (2) and P expressed in (3).



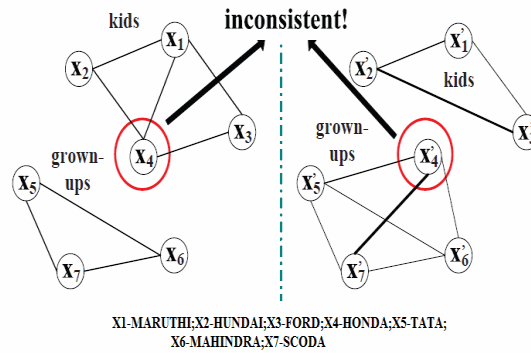
Here  $S$  is the unremitting cover process which belongs to  $R_m X_n$  is an  $m \times n$  matrix, and  $X$  that belongs to  $R_n$  is the vector of unknown variables. If  $P$  belongs to  $R_n$  of a given vector, it is considered that all constraints must be stated as equalities, in each side each constraint must be nonnegative and finally the variable is also be nonnegative value.

The solution to the optimization problem (4) can be found, and it is referred to as the best solution. In order to find an original best result set and make it the current values, it is checked if the solution roundup values are detected. If at least one of the solution sets is improved, it is made as the current value and next the algorithm will go to first step. Other optimization algorithms can also be used to solve the problem given in (5) (Precup et al., 2015), (Valdez et al., 2017).

These steps help the process to determine and identify the consistency in integrated data. The process of inconsistency detection will be next examined in detail.

### 3.2. Inconsistency detection in heterogeneous sources

We define the general problem of inconsistency detection across multiple heterogeneous sources and propose an effective spectral framework to identify such anomalies. The underlying clustering structure of normal objects is usually shared by multiple sources. Figure 4 clarifies the inconsistent detection methods with interrelated examples of car manufacturer company.



**Figure 4.** Inconsistent detection example based on car manufacturer company.

The TSM is the process for the features of system generation when there is a high degree of have common characteristics in previous method. In this method to identify the inconsistency in heterogeneous method we illustrate the upend switch method (USM) in terms of solving the square problem by a least squares values approach. A minimal solution can be found in the set  $X$  of the system  $AX = b$  implemented with an equation to minimize the linear programming method. The value of the inverse  $X$  is a vector with complex components  $x_1, x_2 \dots x_n$ . Differentiating with respect to  $X$  and setting the result to zero we get  $2X - A^T \lambda = 0$ . We combine the input matrices into one matrix that captures the information from each source, but also ensures that individual object

relationships are preserved. We then adopt USM technique to identify the key vectors of the graph of the combined matrix, and identify anomalies by computing cosine distance between the components.

To prove the optimal solution we can make  $A(X - X') = 0$ . This paper makes use of the unknown variables  $X$  can be computed by  $X = A^T (A^T)^{-1} b$ . The vector  $X'$  is used as follows to indicate the solution for non-negative values.

Considering:

INPUT: equivalent metrics  $A$  and  $X$ , number of non-negative values  $P$  and conflict values  $b'$ .

OUTPUT: Detection value of  $\mu$

the algorithm for USM process consists of the following steps:

1. Compare matrix  $X'$  value with  $b'$  value.
2. Compute  $AX' = b'$  values.
3. Conduct  $b'$  values as non-negative values.
4. Compute each object values  $b' \neq b$  based on  $X'$ .
5. Return value=0 as consistent and value $\neq$ 0 is inconsistent.
6. Having a set of objects  $X = \{x_1; x_2; \dots; x_n\}$  there are  $P$  values which provide information about the classification of  $X$  and  $A$  is the number of non metrics values which hold the conflict values of  $b'$ . While comparing  $X$  metrics with  $A$  metrics undetermined values is identical which is not equivalent to  $b'$  values when these values are computed on  $X'$ :

$$X' = \begin{cases} 1 & \text{is assigned to group } b \text{ by a model} \\ 0 & \text{otherwise} \end{cases}$$

The first term ensures that if an object  $x_i$  is assigned to group  $b$ , when  $b=1, 2, \dots, n$  from the metrics  $A$ . the second term imposes the constraints on the other group  $b'$ . When the two metrics are compared the values of prediction is identified by the TSM values which are equivalent to the  $X$  metrics if  $b \neq b'$  values. Thus the constraint value is 0, if the value is  $\neq 0$  the value is determined as inconsistent. Finally the value is said to be  $\mu$ .

An example is given as follows. Using

$$X' = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 200 \\ 0 \\ 0 \\ 500 \end{bmatrix} \quad (6)$$

the same  $X'$  is generated in  $b'$  i.e.

$$AX' = b' = \begin{bmatrix} b1 \\ b2 \\ b3 \end{bmatrix} = \begin{bmatrix} 200 \\ 0 \\ 0 \\ 500 \end{bmatrix} \quad (7)$$

The reliable value of zero can be generated by an expression. If I manipulate the report values of R2 from 500 to 800 to introduce some inconsistency as shown in Table 2. Report R2 has shorter time duration but higher value of report compare with R1.

**Table 2.** Vehicle defect identification after inconsistency.

ID (Ri)	Vehicles	Manufacturer	From	To	Duration	Value(v) out of 1000 product
R1	Omni	Maruthi	1910	1980	70	700
R2	Omni	Maruthi	1930	1980	50	800
R3	Omni	Maruthi	1945	1985	40	600
R4	Omni	Maruthi	1960	2000	40	700

The new TSM value will be

$$AX = P \begin{bmatrix} 111100 \\ 111110 \\ 000001 \\ 000111 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \\ x3 \\ x4 \end{bmatrix} = \begin{bmatrix} 700 \\ 900 \\ 600 \\ 700 \end{bmatrix} \quad (8)$$

and the value for each interval is generated by non-negative values as

$$X' = \begin{bmatrix} x1 \\ x2 \\ x3 \\ x4 \end{bmatrix} = \begin{bmatrix} 200 \\ 0 \\ 0 \\ 500 \end{bmatrix} \quad (9)$$

the equivalent b' matrix as

$$AX' = b' = \begin{bmatrix} b1 \\ b2 \\ b3 \end{bmatrix} = \begin{bmatrix} 200 \\ 0 \\ 0 \\ 500 \end{bmatrix} \quad (10)$$

the value of  $b' \neq b$ , therefore the difference value will be treated as

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} 700-200 \\ 900-800 \\ 600-600 \\ 700-700 \end{bmatrix} = \begin{bmatrix} 500 \\ 100 \\ 0 \\ 0 \end{bmatrix} \quad (11)$$

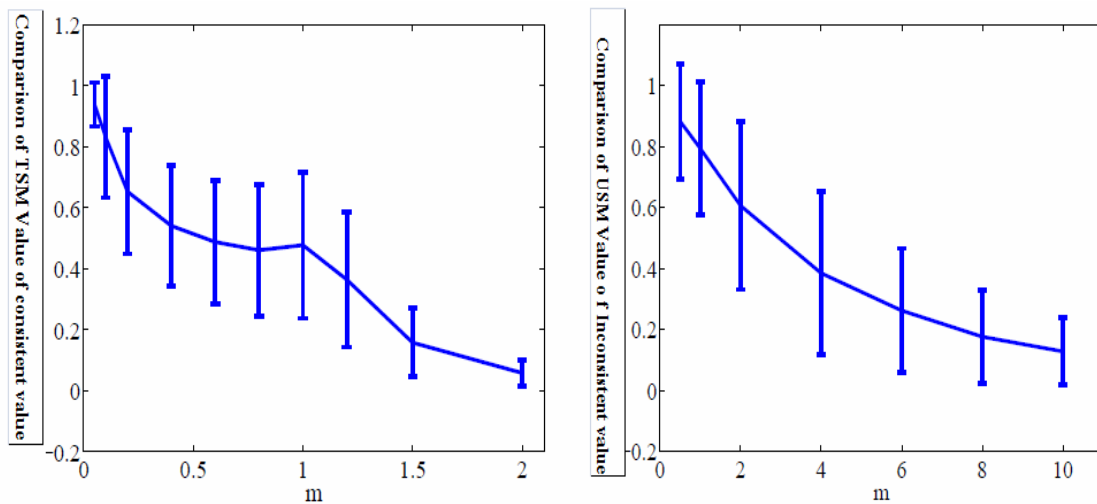
If we only consider the largest subset which contains all reports, some consistent reports may be punished by a nonzero  $\mu$ . But cannot find a perfect solution set for all inconsistent equalities. By comparing the list of conflict report to non-conflict list of reports we can find the precise list of reports that are in conflict with one specific report.

There is no feasible solution; the merged data contain report values that conflict with others. The matrix  $b'$  generated by solution set  $X'$  cannot satisfy all linear equations since  $\mu$  is  $b' \neq b$ , so we examined the combination of report values as illustrated in Table 3.

**Table 3.** Relation between difference and non-difference variation identification

ID	Differences values	Non-differences values	Variation
R1	R2,R3,R4	R3,R4	R2
R2	R1,R3,R4	R3,R4	R1
R3	R1,R2,R4	R1,R2,R4	0
R4	R1,R2,R3	R1,R2,R3	0

From the TSM value we can compute the inconsistent values, and the values are compared with switch method (USM). The variation is compared with the chart presented in Figure 5.



**Figure 5.** Parameter analysis of TSM versus USM.

#### 4. CONCLUSIONS

In the inconsistency detection study the data can hold inconsistency event detection and inconsistent values identification. An efficient approach of using the typical substitution method (TSM) is carried out. When data sources are integrated, the values can be grouped in several systems depending on

their related conditions based on the TPM method we can organize the consistency in large database by S.

The goal of this paper has been to find an approach that can estimate interval values accurately, satisfy most constraints of inconsistency values in heterogeneous data sets, I used the (USM) upend switch method i.e. inverse the TPM value to identify the inconsistency value  $X'$ . Sometimes  $b \neq b'$ , so USM will produce a consistent report on the relation  $R_1, R_2, \dots, R_n$  by comparing the variant and non-variant result we can justify that the values are inconsistent. Concluding,  $\mu$  - non-zero values are inconsistent values.

Experimental results show that the proposed USM algorithm can detect the inconsistent datasets in multiple data sources. Meanwhile these USM techniques can be applied in parallel computing for better results and speed up computation. By synthesizing or comparing multiple information channels, we can identify insightful knowledge and provide users a robust and accurate solution.

## REFERENCES

- Bleiholder J., Naumann F., 2008, Data fusion. *ACM Computing Surveys* **41**(1), 1–41.
- Dong G., Li J., 1999, Efficient mining of emerging patterns: discovering trends and differences. *Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'09)*, 43–52.
- Fan W., Greengrass E., McCloskey J., Yu P.S., Drummey K., 2005, Effective estimation of posterior probabilities: Explaining the accuracy of randomized decision tree approaches. *Proceedings of IEEE International Conference on Data Mining (ICDM'05)*, 154–161.
- Gonzalez H., Halevy A., Jensen C.S., Langen A., Madhavan J., Shapley R., Shen W., 2010, Google fusion tables: data management, integration and collaboration in the cloud. *Proceedings of 1<sup>st</sup> ACM Symposium on Cloud Computing (SoCC)*, 175–180.
- Gao J., Ding B., Fan W., Han J., Yu P.S., 2008, Classifying data streams with skewed class distributions and concept drifts. *IEEE Internet Computing* **12**(6), 37–49.
- Gao J., Fan W., Turaga D., Verscheure O., Meng X., Su L., Han J., 2011, Consensus extraction from heterogeneous detectors to improve performance over network traffic anomaly detection. *Proceedings of IEEE International Conference on Computer Communications Mini-Conference*, 181–185.
- Jeh G., Widom J., 2002, Simrank: a measure of structural-context similarity. *Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'02)*, 538–543.
- Kang U., Meeder B., Faloutsos C., 2011, Spectral analysis for billion-scale graphs: Discoveries and implementation. *Proceedings of Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD'11)*, 13–25.
- Lawson C.L., Hanson R.J., 1974, *Solving Least Squares Problems*. Philadelphia, PA: Prentice-Hall.

Precup, R.-E., Sabau, M.-C., Petriu, E.M., 2015, Nature-inspired optimal tuning of input membership functions of Takagi-Sugeno-Kang fuzzy models for anti-lock braking systems. *Applied Soft Computing* **27**, 575–589.

Ramachandran C., Malik R., Jin X., Gao J., Nahrstedt K., Han J., 2009, Videomule: A consensus learning approach to multi-label classification from noisy user-generated videos. *Proceedings of ACM International Conference on Multimedia (ACM MM'09)*, 721–724.

Su L., Yang Y., Ding B., Gao J., Abdelzaher T.F., Han J., 2011, Hierarchical aggregate classification with limited supervision for data reduction in wireless sensor networks. *Proceedings of ACM International Conference on Embedded Networked Sensor Systems (Sensys'11)*, 40–53.

Valdez F., Vázquez J.C., Melin P., Castillo O., 2017, Comparative study of the use of fuzzy logic in improving particle swarm optimization variants for mathematical functions using co-evolution. *Applied Soft Computing* **52**, 1070–1083.

Zadorozhny V., Lewis M., 2013, Information fusion for USAR operations based on crowd sourcing. *Proceedings of 16<sup>th</sup> IEEE International Conference on Information Fusion*, 1450–1457.

Zadorozhny V., Manning P., Bain D.J., Mostern R., 2013, Collaborative for historical information and analysis: vision and work plan. *Journal of World-Historical Information* **1**(1), 1–14.

Zubko V., Leptoukh G.G., Gopalan A., 2010, Study of data-merging and interpolation methods for use in an interactive online analysis system: MODIS terra and Aqua Daily Aerosol case. *IEEE Transactions on Geoscience and Remote Sensing* **48**(12), 4219–4235.