# Fuzzy motivated data linkage for Self-Adaptive Systems

## G. Roja Ramani, M. SrujanKumar Reddy

**[1]PG Student, Department of Computer Science and Engineering, Madanapalle Institute of Technology & Science, India.**

**[2]Assistant Professor, Department of Computer Science and Engineering, Madanapalle Institute of Technology & Science, India.**

***Abstract:*** *Data-linkage is an extremely important task in several domains. Data-linkage means that, identifying the different data entries that are of same entity from different data sources. Joins (i.e., matches) related data entries or data sets which do not share a common identifier as it's goal. Data linkage commonly divided into two types. Namely One-to-One and One-to-many. In first, we relate one row of tableA with a single matching row of tableB. In case of second, we relate one row of A with many rows of B. In these scenarios data linkage is done across the same type of entities. It is necessary and important to match entities of different types too. For that, we are proposing a new Many-to-Many data-linkage scheme for entities of different types in self-adaptive systems using Fuzzy motivated data linkage. Each and every entity of A is weighted with each and every entity of B. New method uses, Clustering Tree technique and the tree characterizes that entities should be linked together. The tree can be easily transformed into association rules and understandable.*

**Keywords***: Data-linkage, Many-to-Many linkage, fuzzy motivation, self adaptive systems*

## I. INTRODUCTION

Data linkage is the process of identifying dissimilar data entries, which refers to the same entity among different data sources. The aim of this task is, joining the data sets which do not share a common identifier (i.e., a foreign key). Common data linkage scenarios include: linking data when combining the two different databases; data de-duplication, linking related DNA sequences and matching astronomical objects from different catalogues. Data linkage classified into two types: one-to-one and one-to-many. In an earlier, the goal is to link an entity from one data set with a single matching entity in another data set. In the later scenario, the goal is to link an entity from the first data set with a group of matching entities from the other data set.

In the proposed Fuzzy motivated linkage, the linkage takes place between the entities of different types(heterogeneous data sets) of entities. The relation between each and every entity of table A is weighted with each and every entity of table B using one-class clustering tree method developed by fuzzy logic. This method is different from the traditional decision tree method because, in clustering tree every leaf contains a cluster rather than a single classification and the inner nodes represents the features of the entities of first dataset. Every leaf of the tree represents features of the second set that are matching. Whereas decision trees used for regression tasks and classification. Training set which we used for inducing the tree must be labeled. But acquiring a labeled data set is an expensive task. Tree can be easily transformed into association rules and understandable. For inducing tree, we are using Prepruning method.

By an appropriate reorganizing themselves internally, the fuzzy control systems are able to function within a changing outside world. Fuzzy logic used by designers for developing sophisticated control systems are discovering support form associated technologies. By including adaptive control in their planning, fuzzy control systems are proficient to design the systems which can adjust to environmental changes - an important factor in wide array of applications.

Adaptive systems have capability to explain their reasoning and to learn, in additional they get capacity to be modified and extended. Due to the balance among explicit human knowledge and learned responsiveness capabilities, these make the systems extensible, very robust and suitable for solving a variety of troubles.

## II. EXISTING SYSTEM

One-to-Many data linkage, an important task in several domains. Here linkage is performed between entities of same or different types using one-class clustering tree method, where in the linkage problem only two tables are involved. We relate one row of table A with many rows of table B. Here OCCT(one class clustering tree) is used for the induction of tree. But it is also very important to develop linkage techniques that link among matching entities of different types as well.

*Drawbacks:*
- ✓ Data linkage is performed between the entities of the identical or different types.
- ✓ Performs one-to-many data linkage only.

## III. PROPOSED SYSTEM

We are propose a new Many-to-Many data linkage scheme using Fuzzy motivated data linkage, which links among the entities of different types. The proposed scheme is based on a clustering tree that characterizes the entities that should be linked together. Tree is constructed to understand and transform entities into association rules easily, i.e., the interior nodes consists of features describing the first dataset, and the tree leaves describes features of their related entities from the second data set. For inducing tree we are using four splitting methods and one pruning criteria i.e., Prepruning.

Fuzziness is an alternative to randomness. Randomness describes the uncertainty of event occurrences. It describes the event ambiguity and measures the degree to which an event can occur, but not about the event i.e., event occurs or not. If an event occurs, that is "random" and to which degree it occurs is "fuzzy".

For this purpose we are using the real Weather Information data received from the revenue department as the data for the tree induction as training data and the corresponding disease cases(i.e., Malaria occurrence for corresponding weather conditions). As the weather conditions changes the occurrence of disease cases also changes i.e., adaptive nature. For the given data weights of each data item will be calculated so that the items which are having less threshold than expected threshold are pruned using Miner Algorithms and Fuzzy Membership Function given in modules.

*Advantages:*

- Links between entities of different types
- More efficient and Scalable
- Convenience and Portability
- Flexibility
- Global Opportunities

**Algorithm:**

BEGIN
  1.Initializae Ftree
  2.Read Itemset
  3.Calculate WI-Support
  4.Insert all items into Ftree
  5.Join prefix items
  6.Identify prunable items
   **7. If** FIWI-Support(i)>MST **then**
     8.insert I to Ftree
       9.Else
          10.prunable items

    11.Return tree
  12.End **if**
END

## IV.    ARCHITECTURE DIAGRAM



## V.    IMPLEMENTATION

*Modules*
  1. Item Set Mining
  2. Weighted Transaction Equivalence
  3. The Infrequent Weighted Item set Miner Algorithm
  4. The Minimal Infrequent Weighted Item set Miner Algorithm

*Modules Description*
  ### 1.    *Item Set Mining:*
Item set mining is an exploratory data mining technique widely used for discovering valuable correlations among data. The first attempt to perform item set mining was focused on discovering frequent item sets, i.e., patterns whose observed frequency of occurrence in the source data (the support) is above a given threshold. Frequent Item sets find application in a number of real-life contexts (e.g., market basket analysis , medical image processing , biological data analysis). However, many traditional approaches ignore the influence/interest of each item/transaction within the analyzed data. To allow treating items/transactions differently based on their relevance in the frequent item set mining process, the notion of weighted item set has also been introduced . A weight is associated with each data item and characterizes its local significance within each transaction.

### 2.    *Weighted Transaction Equivalence*
The weighted transaction equivalence establishes an association between a weighted transaction data set T, composed of transactions with arbitrarily weighted items within each transaction, and an equivalent data set TE in which each transaction is exclusively composed of equally weighted items. To this aim, each weighted transaction tq corresponds to an equivalent weighted transaction set, which is a subset of TE's transactions. Item weights in tq are spread, based on the irrelative significance, among their equivalent transactions in TEq. The proposed transformation is particularly suitable for compactly representing the original data.

### 3.    *The Infrequent Weighted Item set Miner Algorithm*
        A weighted transactional data set and a maximum IWI-support (IWI-support-min or IWI-support-max)threshold, the Infrequent Weighted Item set Miner algorithm extracts all IWIs whose IWI-support satisfies. Since the IWI Miner mining steps are the same by enforcing either IWI-support-min or IWI-support-max thresholds, we will not distinguish between the two IWI support measure types in the rest of this section. IWI Miner is a FP-growth-like mining algorithm that performs projection-based item set mining. Hence, it performs the main FP-growth mining steps:
        (a) FP-tree creation and
        (b) Recursive item set mining from the FP tree index. Unlike FP-Growth, IWI Miner discovers infrequent weighted item sets instead of frequent (unweighted) ones. To accomplish this task, the following main modifications with respect to FP-growth have been introduced:
        (i) A novel pruning strategy for pruning part of the search space early and

3680

(ii) A slightly modified FP tree structure, which allows storing the IWI-support value associated with each node.

### 4. The Minimal Infrequent Weighted Item set Miner Algorithm

A weighted transactional data set and a maximum IWI-support (IWI-support-min or IWI-support-max) threshold , the Minimal Infrequent Weighted Item set Miner algorithm extracts all the MIWIs that satisfy . The pseudo code of the MIWI Miner algorithm is similar to the one of IWI Miner. The MIWI Mining procedure is invoked instead of IWI Mining. The MIWI Mining procedure is similar to IWI Mining. However, since MIWI miner focuses on generating only minimal infrequent patterns, the recursive extraction in the MIWI Mining procedure is stopped as soon as an infrequent item set occurs. In fact, whenever an infrequent item set I is discovered, all its extensions are not minimal.

## VI.    RESULTS



Fig 1 Precision and Recall graph for Fuzzy linkage in self-adaptive systems

***Result analysis:*** As we know that precision and recall shows the best performance of a system. Hence, the above graph shows the best performance for self-adaptive systems by using fuzzy linkage when compared to traditional OCCT: One-to-Many data linkage method.



Fig2. Reading data items



Fig 3. Setting attributes for data items



Fig 4. After setting attributes weighted will be allocated to them



Fig 5. Fuzzy function

When the weights were set, the actual Fuzzy membership function works means that the weights randomly shared for many-to-many linkage among data items.
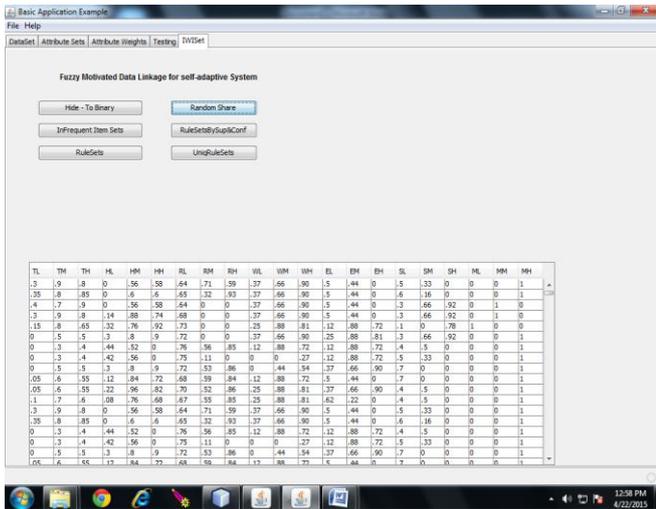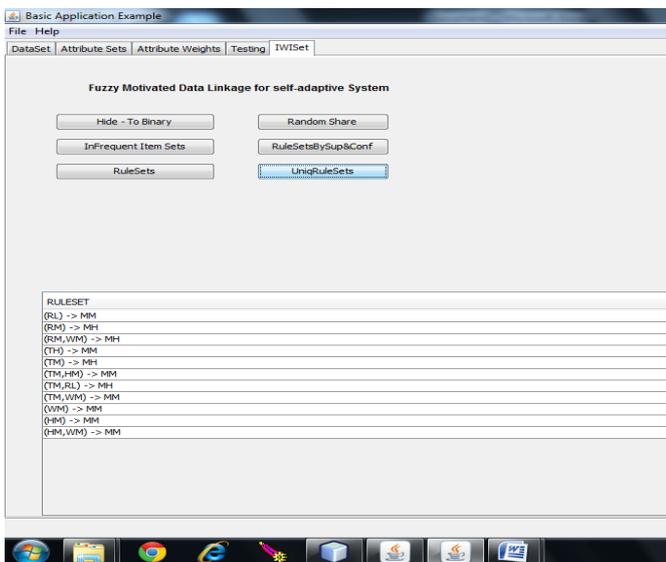
Fig 6. Random sharing of weights



Fig 7. Unique rule sets

When the random share is done, then rules will be generated as per support & confidence. There after unique rules will be generated.
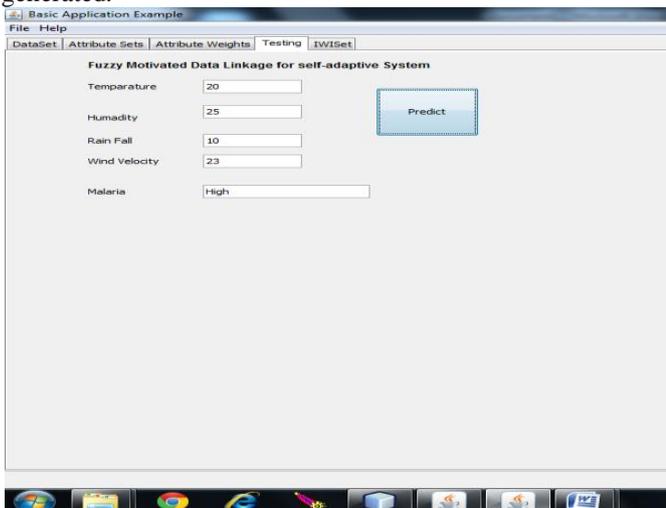


Fig 8. Figure showing final result

After generating unique rules, the testing takes place i.e., gives prediction results for disease occurrence may be like high/medium/low for the given weather information.

## VII.  CONCLUSION AND FUTURE WORK

We are providing Fuzzy mechanism for developing many to many data linkage in heterogeneous datasets using clustering tree for the induction of tree. The linkage takes place between the entities of different types(heterogeneous data sets) of datasets. The relation between each and every entity of A is weighted with each and every entity of B and linkage is done among entities by fuzzy logic. This method is different from the traditional decision tree method. On including adaptive control systems to this, we design the system which can adjust to environmental changes - an important factor in wide array of applications. And it improves efficiency of linkage capacity over different data sets.

In future, we are also trying to expand the scope of this system involving different environments.

## VIII.  REFERENCES

[1]. Ma'ayan Dror, Asaf Shabtai, Lior Rokach, and Yuval Elovici, "OCCT: A One-Class Clustering Tree for Implementing One-to-Many Data Linkage" IEEE tractions on knowledge and Data Engineering, vol. 26, no. 3, March 2014.

[2]. Hendrik Blockeel, Luc De Raedt, Jan Ramon, ""Top-Down Induction of Clustering Trees," ArXiv Computer Science e-prints, pp. 55-63, 1998.

[3]. M. Yakout, A.K. Elmagarmid, H. Elmeleegy, M. Quzzani, A. Qi, "Behavior Based Record Linkage," Proc. VLDB Endowment, vol. 3, nos. 1/2, pp. 439-448, 2010.

[4]. I.P. Fellegi and A.B. Sunter, "A Theory for Record Linkage," J. Am.Statistical Soc., vol. 64, no. 328, pp. 1183-1210, Dec. 1969.

[5]. A.J. Storkey, C.K.I. Williams, E. Taylor, and R.G. Mann, "An Expectation Maximisation Algorithm for One-to-Many Record Linkage," Univ. of Edinburgh Informatics Research Report, 2005.

[6]. Dick, S.; Dept. of Electr. & Comput. Eng., Univ. of Alberta, Edmonton, Canada "Towards complex fuzzy logic" IEEE tractions on Fuzzy Systems,Volume:13, Issue:3, June 2005.

[7]. Giannini, J.A. ; Appl. Res. Lab., Pennsylvania State Univ., University Park, PA, USA ; Kilgus, C. "A fuzzy logic technique for correcting climatological ionospheric models" IEEE transactions on Fuzzy Systems, (Volume:35, Issue:2 ), March 1997.

3682