

---

# Efficient Record Linkage Using Rapid Miner Tool

Dr.D. Vasumathi<sup>1</sup>, M. Ravikanth<sup>2</sup>, B. Mallikarjuna Reddy<sup>3</sup>

<sup>1</sup>Assoc.Prof of CSE & Addl.Controller. Of Exams JNTUH, AP, INDIA

<sup>2</sup>Assit Prof of CSE Dept, Malla Reddy Engineering college, AP, INDIA

<sup>3</sup>Software Engineer,Koushik web solutions, Bangalore, INDIA

---

## Abstract

Rapid miner tool discovers a knowledge flow modeled as an operator sub routines or tree generation. Tool provides the workflow information in the form of XML format representation process. Different kinds of XML files data collects and arranged as a sequence alignment environment process. This kind of arrangement is defined as a pipelining process. Through pipeline we are extract the large amount of records with high scalability environment process. Large amount of records we are starts the extraction process from multiple numbers of databases. Multiple numbers of databases communicates with remote environment process. All databases works based on aggregation operation implementation process. Linkage of records environment creation process starts with multi layer perception. Multi layer environment creates the good resolution record system generation. Resolution system creates with sequential decision making environment process. Sequential decision making provides good linkage strategy results. Those results are good interesting and performance results evidence mechanism identification process.

**Keywords:** Rapid Miner, XML files, Aggregation, Multi layer environment, and Linkage strategies.

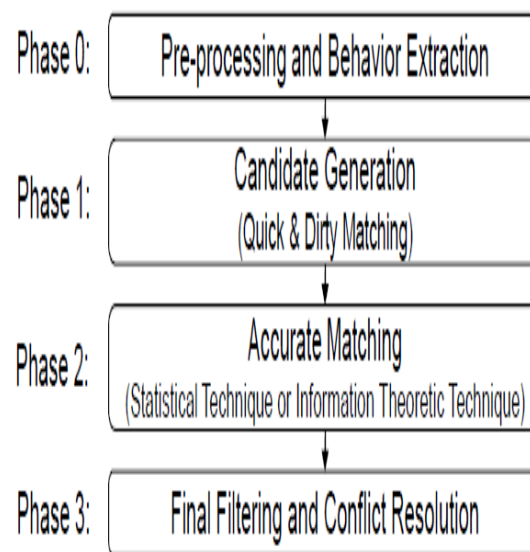
---

## Introduction:

We are introduces the two step record linkage environment for extracting efficient results identification. It can works based on clustering environment process. First step linkage process data we are divide into different number of clusters. Next to start the merging process of different kinds of clustering environment and provides the each and every cluster performance specification process. It can provide large cardinality results identification and provides exact matching results. Two step linkage environment shows identifies the good error identification also possible. It can increases the results under the implementation. Two step strategies focuses on the identification results based on the categorization wise specification process identification. It can works as a good reduction technique environment process and reduces the overhead problems. Two step linkage strategies provide faster results for the query extraction process inside the web databases. These kind of process techniques implementation possible like as a good web miner tool creation process.

We are introduces the global decision making analysis under implementation part specification process. Using global decision making environment provides good schema results specification environment process. In different kinds of databases we are apply the preprocessing. It can extract the results of information as a case based environment. After getting that results one more time we are apply the classification environment process. Next step implementation applies the some kind of rules or constraints specification process. It can reduce the similarity of record pair's environment process specification process. At last produces the efficient results identification process.

Many numbers of approaches are starts the extraction of record pair's of information identification.



**Fig 1: Multi Layer approach environment**

Using query matching implementation generates the candidate sets of results information. It can contain different number of states specification process.

**Phase0** starts with preprocessing environment and good encoding environment process identifies the accurate results identification process.

**Phase1** using the functions like quick and dirty and 2-dimensional environment we are apply here. It can shows inside the output results like compact environment process.

**Phase2** merge the different kinds of clusters environment and provides satisfiable generate results we are provides here. It can show index based results specification process. We are applying the mapping and distribute the results as a effective process results.

**Phase3** It can avoid all kinds of conflict results information identification process. It can show the final results as a low scoring environment process only. Whenever to follow all kinds of phases under the implementation we are increases 5 to 10% are increases under extraction process results specification

process. It can show all results as a effectiveness results. These are steps works as a pipeline representation process.

### **Related Work:**

We are shows the comparison in between of present record linkage techniques to previous linkage techniques specification process. Previously we are works on two files remove the duplication content and identify the content of information like un-duplication content and merge the files of content specification process.

We are implementing the record linkage process from one to one record mapping environment process. This is same procedure we are follow linear assignment of content identification.

After some number of days we are people starts the research in Bayesian learning principles. It can works based on the probability based conditions under the web databases implementation process. We are generating the blocks of content information and in between of one record to another record maintain the logical relationship environment process. Logical relationship content extract the information in the form hierarchical cluster results information.

We are starts the implementation procedure specification process with micro data analysis process environment creation. It can show the results of information as a duplicate data environment process. We are starts the capturing process and recapturing process results. It can works as a probabilistic record linkage environment process generation. We are assigning the weight for those particular records specification environment process. Which records are contains highly sufficient weighted those things are select under implementation procedure.

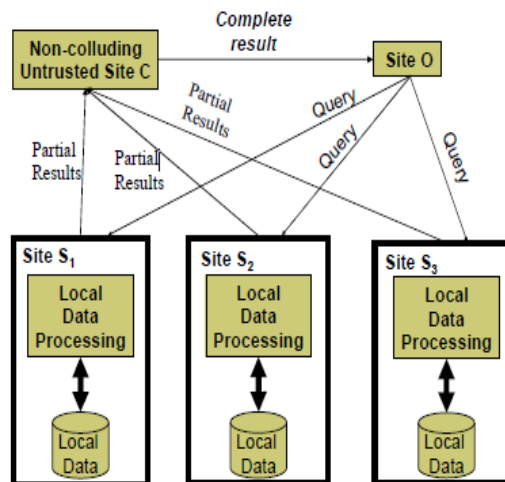
Some of the techniques also we are implemented here. Non statistical matching environment we are implemented here. It can start the extraction procedure with slow movement extraction pattern.

At last we are implementing automated record linkage techniques inside the implementation. It can derive the results as faster results.

### **Implementation work:**

We are starts the implementation with different kinds of conditional probabilities and agreement pattern implementation process identification. Generally we are starts the specification of extraction results as a final results.

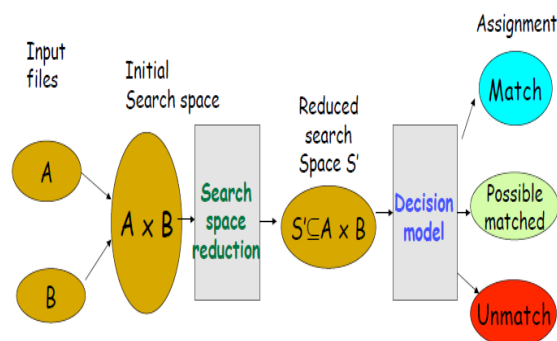
In two step record linkage systems we are implement two linkage rules. It can start the extraction results with three different categories procedure that is called as a link, non link and undecided environment process environment.



**Fig 2: Linkage rules implementation architecture**

### SVM based iterative clustering with de-duplication results identification

We are takes the input files inside the search engine specification environment process. Starts the searching process for identification of reduction results specification process. One file to another file we are apply the multiplication environment process implementation identification process generation. Many number of files features we are extract as a final results specification. All records are matched under the decision model implementation process. Decision model decides whether the results are correct or incorrect results like matched, unmatched and possible matched.



**Fig 3: Proposed Architecture**

### Automated approaches with Semantic Annotation

Annotations are provides the good referential results identification. It can show the different kinds of exact match results. Semantic annotations are extracting the results based on threshold, schema matched results specification identification process. We are trying to place with Expectation and

maximization environment. It can collect all kinds of features from different locations environment process.

1. We are collect the features like overall rank score and displayed inside the framework representation process.
2. It can select two or more number of references, it can work iteratively.
3. Eliminate redundant matching

Fig 4: Data sets selection process.

4. All kinds of results are placed inside the vector.
5. Apply the dice similarity.
6. Apply the dice similarity.

### Experimental Datasets:

### Conclusion:

Two step record linkage extract good semantic annotation results with reference sets, aggregation operations results environment process. It can shows the results as a label based environment process. One reference set to another reference set we are maintain the information mediator selection process. It can works as a good information extraction environment process identification process. Semantic annotation provides the good graphical model results in Output state.

Name	Source	Reference Set Match	Records
BFT	Bidding For Travel	Hotels	1,125
EBay	EBay Comics	Comics	776
Craigs List	Craigs List Cars	Cars, KBBCars (in order)	2,568
Boats	Craigs List Boats	None	1,099

It can allow the more number of large scale results environment process. Semantic annotations extracts as effective results without any burden as a label based results.

Name	Source	Attributes	Records
Fodors	Fodors Travel Guide	name, address, city, cuisine	534
Zagat	Zagat Restaurant Guide	name, address, city, cuisine	330
Comics	Comics Price Guide	title, issue, publisher	918
Hotels	Bidding For Travel	star rating, name, local area	132
Cars	Edmunds & Super Lamb Auto	make, model, trim, year	27,006
KBBCars	Kelly Blue Book Car Prices	make, model, trim, year	2,777

## References:

1. An efficient record linkage scheme using graphical analysis for identifier error detection, 2011
2. A Two-Step Classification Approach to Unsupervised Record Linkage, 2007
3. Iterative Record Linkage for Cleaning and Integration, 2011
4. Public Record Aggregation Using Semi-supervised Entity Resolution, 2011
5. Record Linkage Modeling in Federal Statistical Databases, 2011
6. Integrated Data System Person Identification: Accuracy Requirements and Methods
7. Behavior Based Record Linkage
8. Record Linkage: Theory and Practice
9. Linkage in Medical Records and Bioinformatics Data
10. Weka tutorials

## About The Authors



(1) Dr. D. VASUMATHI, currently working as Associate Prof. in department of computer science and engineering from JNTU College of engineering, Hyderabad. Her Research completed on Data Mining from JNTUH. She has 10 years of teaching experience. She has published several papers in international and national journals and conferences. She's the member of CSI, IEEE. Areas of Specialization are Data Mining, Networks, Cloud Computing, Image Processing, and Distributed Databases



(2) M. Ravikanth. Completed M.Tech from JNTUK in a Specialization Computer Science and Engineering..and PhD pursuing in JNTU Hyd. Presently Working as Asst. Prof in MREC, Hyderabad. Area of Specialization are Data Mining, Data Warehouses, DBMS, software engineering..



(3) B. Mallikarjun Reddy pursuing M.Tech from JNTU Hyd in a Specialization Computer Science and Engineering.. Presently Working as Asst. Prof in CRVCET, JNTU, and Hyderabad. Areas of Specialization are Data Mining, Data Warehouse, DBMS, and Automata Theory.