

System for Locating Record Duplication Using Genetic Programming

Mr. Satish Wagh, Prof. Dr. Varsha Patil,
Department of Computer Engineering,
Matoshri College of Engineering and Research Center, Eklahere, Nashik (MS), India, MCERC, Nashik,

Abstract—

With rapid advancement in technology enables high uses of data base which causes duplication of database management. The replicated data records generate multiple copies of dirty data and contain improper spelling and punctuation, incorrect data is associated with record, incomplete and also outdated or same data is duplicated in the database. For removing the dirty data various data cleaning software and applications are used. The proposed generic programming approach to record deduplication helps in combining several different pieces of evidence extracted from the data deduplication function helps in identifying whether two entries are in repository replica or not. In this approach for automatic adaptation of function are given fixed replica identification boundary for improving accuracy in terms of number of duplicate records found versus the actual number of duplicate records. This work has been tested on the blood bank database to duplicate the records. System used for locating record deduplication with the help of genetic programming helps in combining evidence extracted from the content and deduplication function enables in identification whether two entries present in repository are replica or not. Additionally genetic programming is capable for automatically adapting these function to a specific replica identification boundary this genetic programming approach is applied for various database management to find deduplicated records.

Index Terms— Database integration, Evolutionary computing and Genetic algorithms, Database integration

I. INTRODUCTION

SEVERAL systems such as digital libraries and other database systems like organization databases are affected by the duplicates. The system for locating record deduplication using genetic programming approach to find a deduplication function that is able to identify whether two entries in a repository are replica or not. Deduplication is a task of identifying the duplicate data in a repository that refers to the same real world entity or object and systematically substitute the reference pointers for the redundant blocks; also known as

storage capacity optimization. Dirty data is defined in various categories performance degradation - as additional useless data demand more processing, more time is required to answer user queries; quality loss the presence of replicas and other inconsistencies lead to distortions in reports and misleading conclusions based on the existing data; increasing operational costs because of the additional volume of useless data, extra costs are required on more storage media and extra computational processing power to keep the response time levels acceptable. To avoid these problems, it is necessary to study the causes of dirty data in repositories. A major root is the presence of duplicates, or near duplicates

in these repositories, those constructed by the aggregation or integration of distinct data sources. The issue of detecting and removing duplicate entries in a repository is generally known as record deduplication. In our project to remove the dirty data in the blood bank management system. As a part of genetic programming approach the gaining concepts and the entropy calculations are used to deduplicate the records.

II. LITERATURE SURVEY

A literature has been dedicated to the Record Deduplication and tremendous progress has been made ranging from efficient and scalable algorithm Record Deduplication in Cora and Restaurant dataset. Record deduplication is a growing research topic in database and related fields such as digital libraries. This problem arises mainly when data are collected from disparate sources using different information description styles and metadata standards. Common place for replicas are found in data repositories created from OCR documents. These situations can lead to inconsistencies that may affect many systems such as those that depend on searching and mining tasks. To solve these inconsistencies it is necessary to design a deduplication function that combines the information available in the data repositories in order to identify whether a pair of record entries refer to the real-

world identity. In realm of bibliographic citations, for instance, this problem was discussed by Lawrence et al. [19], [20]. Lawrence proposed a number of algorithms for matching citations from different sources based on word matching, phrase matching, and field extraction. As more strategies for extracting disparate pieces of evidence become available, Elmagarmid et al. [21] following two categories:

A. Ad-Hoc or Domain Knowledge Approaches

This category includes approaches that usually depend on specific domain knowledge or specific string distance metrics. Technique make use of declarative languages [21] can be also classified in this category.

Training-based Approaches

This category includes all approaches that depend on some sort of training supervised or semi-supervised in order to identify the replicas. Machine learning approaches fall into this category. Next, briefly comment on some work based on these two approaches (domain knowledge and training-based), those that exploit the domain knowledge and those that are based on probabilistic and machine learning techniques.

Pros- Closest matching approach.

Cons-

1. For record matching high weight tokens are required. 2. Identification of record replication is done on individual basis.

B. Probabilistic Approaches

W.W.Cohen (2002) was proposed, one to address the recorded duplication problem as a Bayesian inference problem (a probabilistic problem) and proposed the first approach to automatically handle replicas and recorded duplication problem as a Bayesian inference problem (a probabilistic problem) and proposed the first approach to automatically managed duplication.

To elaborate statistical method. Work, Fellegi and Sunter [5] proposed an elaborated statistical approach to deal with the problem of evidence. These methods rely on the definition of two boundary values that are used to classify a pair of records as being duplicate records or not.

Febre [2] to implement Tools for this method, such as, work with w boundaries as follows:

Pros- The positive identification boundary is the similar value lies above the boundary, The records are considered as replicas
Cons-

The negative identification boundary is the similar value lies below the boundary, The records are considered as not duplicate records

C. Machine Learning Approaches

M. Bilenko (2003) use machine learning technique to improve both the similarity functions that are applied to compare record fields and the way the pieces of evidence are combined. In their system, called Marlin, [21][22]. The training dataset is assumed to have similar characteristics to the test dataset, which makes solution using machine learning techniques and to generalize their solutions for duplicate records.

Pros- Using SVM Classifier for betterment of replica identification.

Cons-

Active Atlas [28] system required for record mapping in order to establish relationship.

III. IMPLEMENTATION DETAIL

- 1) Initialize and set the population (random or user provided individuals).
- 2) Evaluate all individuals in the current population, apply a numeric rating or fitness value to each one.
- 3) If the termination criterion is fulfilled, then execute the last step and continue.
- 4) Reproduce the best individuals in the next generation population.
- 5) Select individuals that will process the next generation with the best parents.
- 6) Apply the genetic operation to all individuals selected and compose the next population. Replace the existing generation of the generated population and go back to Step 2
- 7) Present best individuals in the population as the output of the evolutionary process.

A. Mathematical Model

System for locating recorded duplication using genetic programming approach illustrates the mathematical model that contains duplicate records. By using Genetic Programming, cosine similarity can be derived from this dataset and duplicate functions are formed and it is being constructed in the form of tree.

Dataset (ds): - Any Dataset containing different records from different domain
Extraction (fe): - Extract Feature from input dataset
Similarity Function (sf): - Find the Cosine Similarity between Records.

TreeConstruction(tc):-deduplicationFunction

$S = \{P, FE, SF, TC, GO, O\}$

Let $D: \{d_1, d_2, \dots, d_n\}$ where D consists of no of records Let $FE: \{f_1, f_2, f_3, \dots, f_n\}$ where FE is function which extract features

Let $SF: \{v_1, v_2, v_3, \dots, v_n\}$ where v consists of similarity function records

Let $TC: \{n_1, n_2, \dots, n_n\}$ where tc is a tree constructor or a function which constructs a tree

Let $GO: \{r_1, r_2, \dots, r_n\}$ where GO function is used to generate best fitness result

Let $O: \{o_1, o_2, \dots, o_n\}$ where O consists of output

Function $F1$ returns the features extracted from dataset record

$F1(d) \rightarrow FE$ for downloaded dataset

$F1(f) \rightarrow \{f_1, f_2, \dots, f_n\}$ ϵ to FE

Function $F2$ returns the similarity values of each record

$F2(FE) \rightarrow SF$

e.g. $F2(FE) \rightarrow \{v_1, v_2, \dots, v_n\}$ ϵ to SF . Function $F3$ returns the constructed tree or node.

$F3(SF) \rightarrow TC$

e.g. $F3(SF) \rightarrow \{n_1, n_2, \dots, n_n\}$ ϵ to TC . Function $F4$ returns the genetic operation result. $F4(SF) \rightarrow GO$

e.g. $F4(SF) \rightarrow \{r_1, r_2, \dots, r_n\}$ ϵ to GO . Function $F5$ returns the output

$F5(GO) \rightarrow O$

e.g. $F5(GO) \rightarrow \{o_1, o_2, \dots, o_n\}$ ϵ to O .

Functional Dependency of the above functions

	F1	F2	F3	F4	F5
F1	1	0	0	0	0
F2	1	1	0	0	0
F3	0	1	1	0	0
F4	0	0	1	1	0
F5	0	0	0	1	1

I. SYSTEM DESIGN

A. Data Preprocessing

Cosine similarity is a measure of similarity between two vectors of an inner product space that measures the cosine of the angle between them. The cosine similarity value of 0 is 1, and less than 1 for any angle; the minimum value of the cosine is -1. The cosine angle between two vectors is determined whether the two vectors are pointing in roughly the same direction. Cosine similarity between two strings is calculated using Levenshtein distance and Soft-TFIDF similarity method.

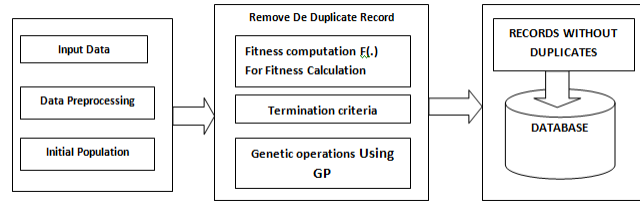


Fig.1. System Architecture

The cosine similarity function using Genetic programming approach automatically selects populations. Cosine similarity function is also used in the training phase to capture the characteristics of the dataset. Based on the characteristics of the dataset, only populations are selected.

B. Feature Vector Extraction

populations is the feature vectors selected based on cosine similarity functions. If the fitness function does not reach the fitness value, then populations are changed. e.g. (att1, att2, att3) Selected fitness function by machine learning approach has to reach the fitness value.

If the fitness function does not reach the fitness value, then have to change the fitness function using genetic operations. Genetic operations are mutation, crossover, reproduction. Selected fitness function has to be represented in tree format, for applying genetic operations easily.

C. Genetic Operations

If the function does not reach a fitness value, have to apply genetic operations again to change the fitness function.

D. Reproduction

Reproduction is the operation of Genetic programming that copies individuals without modifying them. The operator is used to implement elitist strategy that is adopted to keep the genetic code of the fittest individuals across the changes in the generations. If a good individual is found in previous generations, it will not be lost while executing the evolutionary process.

E. Crossover

The crossover operation allows genetic content which are the process of exchanging two parents, in this process that can generate two more children. In a Genetic programming evolutionary process, two parent trees are selected according to a matching (or pairing) policy and, then, a random subtree is selected in each parent.

Mutation

The mutation operation has been implemented for keeping a minimum diversity level of individuals in the population thus it avoids premature convergence.

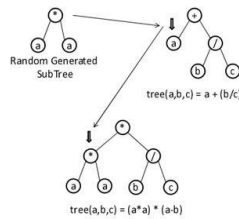


Fig.2.RandomSubtreeMutation

F. Removing Duplicate Records

After selecting a fitness function first calculate fitness value of all records using a fitness function. The Fitness Function selected using existing machine learning approach. These selected fitness function can remove the records. Before storing records into a database calculate fitness values for all records if two records match with same fitness value then remove one record.

II. GENETICALGORITHM

In production scheduling, population of solutions consists of many answers that may have different sometimes conflicting objectives. It operates on a population of solutions rather than a single solution. For example, in one solution it may be optimizing a production process to be completed in a minimum time. In another solution it may optimize for a minimum defects in our final product. As an increase in the number of objectives system are trying to achieve and also increase the number of constraints on the problem and similarly increase the complexity. Genetic programming are better for these types of problems where the search space is more large and the number of feasible solutions is very small. Modeling the Record Deduplication Problem using GP while using Genetic Programming to solve a problem, there are some ordinary requirements that must be fulfilled, order to successfully explore this technique which are based on the data structure used to represent the solution of record deduplication. It uses a tree-based GP representation for the evidence combination function, since it is a natural representation for this type of defined fu

ctions.

These requirements are the following:

- 1) The problem solution must be modeled as a tree structure.
- 2) The evolutionary operations applied over the modeled tree must beat the end of the result and converted into a valid tree finally.
- 3) The modeled tree must be automatically evaluated in order to make the use of this technique viable.

Evidence E_i is a pair $(i, \text{attribute})$, similarity function f_i that represents the use of a specific similarity function over the values of a specific attribute found in the data repositories being analyzed. e.g., it was used to deduplicate a database table with four attributes (e.g., name, surname, address, and postal code) using a specific similarity function (e.g., the Jaro function [Koudas et al. 2006]), would have the following pieces of evidence:

Algorithm 1 Algorithm 1: Genetic Programming

```

Let database (DB) is the set of deduplicated records;
Let sim (Similarity) is the similarity function of Duplicate records;
/* Data Preprocessing */
Generate a set of pairs  $P = p_1 :: p_n$ 
from DB; Compute  $sim(p)$  for each pair  $p \in P$ ;
 $R \leftarrow rank P$  according to  $sim(p)$  values; Let  $T$  be the top  $k$  pairs in  $R$ ;
Let  $B$  be the bottom  $k$  pairs in  $R$ ;
Use  $r$  labels each pair  $p \in T \cup B$  with  $L_p \in \{T, F\}$ ; Compute a weight  $W_p$  for each  $p \in P$ ;
/* Process Evaluation */
 $Gen_0 \leftarrow$  Generate  $m$  functions;
Evaluate ( $Gen_0, T \cup B$ );
Compute a weight  $W_f$  for each function  $f \in Gen_0$ ;
 $Com_0 \leftarrow \emptyset$ ;  $i = 0$  to a predefined number of generations do
   $Com_i \leftarrow$  top  $C$  functions in  $Gen_i$ ;
  /* Active Learning */ each pair  $p \in P$  do
     $Com_i$  labels  $p$  with  $L_p \in \{T, F, D\}$ ;
    if  $L_p = D$  then Use  $r$  labels  $p$  with  $L_p \in \{T, F\}$ ;
  /* Reinforcement */
  Update  $W_p$  for each  $p \in P$ ;
   $Gen_{i+1} \leftarrow$  Crossover Mutation ( $Gen_i, Com_i$ );
   $Gen_{i+1} \leftarrow Gen_{i+1} \cup Com_i$ ; Evaluate ( $Gen_{i+1}, P$ );
  Compute  $W_f$  for each function  $f \in Gen_{i+1}$ ;
   $Com_i \leftarrow$  top  $C$  functions in  $Gen_i$ ;
  foreach pair  $p \in P$  do do
     $Com_i$  labels  $p$  with  $L_p \in \{T, F\}$ 

```

E_1 (name, Jaro), E_2 (surname, Jaro), E_3 (address, Jaro), and E_4 (postal code, Jaro) For this example, a very simple function is a linear combination such as $F_s(E_1; E_2; E_3; E_4) = E_1 + E_2 + E_3 + E_4$.

I. EXPERIMENTS

A. Dataset

Apart from the tasks related to the logical and the physical side of the database system, Database Administrator may also take part in database System operations. main role is to give developers recommendations about the DBMS specificities, thus helping them avoid any data limiting database performance issues. Other important work of the DBAs are related

to data modeling to optimizing the system, as well as to the creation and analysis of new databases. computer science, is evolutionary computation subfield of artificial intelligence (more particularly computational intelligence) that involves combinatorial optimization problems.

To used Restaurant and Cora dataset to analyze the proposed algorithm and the performance of the proposed algorithm is compared against the genetic programming technique with the help of evaluation metrics. In this experiment select datasets from the cora dataset data repository and the datasets used is Restaurant dataset. The datasets, which are used in proposed approach, is detailed below. **Dataset 1 (CORA):** The Cora dataset consists of duplicate and non-duplicate data records which is cited to 122 Conference paper, it is divided into various attributes (name, year, title, and other information) **Dataset 2 (Restaurant):** This dataset contains 500 records (400 originals and 100 duplicates), duplicate record based on one original record (using a Poisson distribution of duplicate records) and with a maximum of two modifications in a single attribute and in the full record.

I. CONCLUSION AND FUTURE WORK

Ending with the duplicate information is stored from different sources that require storage space for storing replicas or duplicated records. several events has been experimented on modern enterprise and data loss was common in all. Various methods have been tried to tackle with that problems identification and replica handling is important for guaranteeing the quality of information which is being available by data intensive system such as digital libraries and e-commerce brokers. This project help in presenting system for locating record duplication via genetic programming. Our approach is combination of several pieces of evidence that has been extracted from the data content for the production of duplication function that enables for identification whether two or more entries are in the repository or replicas or not. In future this approach will help in finding the complex matches in three different data repository scenarios where the repositories partially shares some data

some common data. The entire proposed work concept is concluded with the literature survey, design and modelling, project objective using system. For locating record duplication. Our evolutionary approach has implemented range of fuse for our future work. Experiments were carried out for evaluation our work with data from specific application domains with different datasets. Other investigation techniques were carried out. This paper suggests efficient implementation of kNN for betterment of results and advancement in technology.

ACKNOWLEDGMENT

I am thankful to Prof. Dr. V.H. Patil HOD and vice Principal, MCERC, Nashik for giving her precious time and guidelines during this paper also for her expert guidance and continuous encouragement throughout this paper. I would like to express deep appreciation towards Prof. Dr. G.K. Kharate, Principal MCERC, Nashik, and Prof. Ranjit Gawande (ME Coordinator) whose invaluable guidance supported me in completing this paper.

REFERENCES

- [1] Joises G. de Carvalho, Alberto H.F. Laender, Marcos Andre Goncalves, and Altigran S. da Silva, *A Genetic Programming Approach to Record Deduplication* IEEE Transaction on Knowledge and Data Engineering vol. 24, No. 3, March 2012
- [2] Agrawal, S. Chaudhuri, G. Das, A. Gionis, *Automated ranking of database query results*, in: *Proceedings of the First Biennial Conference on Innovative Data System Research*, 2003
- [3] Koudas, S. Sarawagi, and D. Srivastava, "Record Linkage: Similarity Measures and Algorithms," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 802-803, 2006.
- [4] Bhattacharya and L. Getoor, *Iterative Record Linkage for Cleaning and Integration*, *Proc. Ninth ACM SIGMOD Workshop Research Issues in Data Mining and Knowledge Discovery*, pp. 11-18, 2004.
- [5] P. Fellegi and A.B. Sunter, *A Theory for Record Linkage*, *J. Am. Statistical Assoc.*, vol. 66, no. 1, pp. 1183-1210, 1969.
- [6] S. Verykios, G.V. Moustakides, and M.G. Eljegy, "A Bayesian Decision Model for Cost Optimal Record Matching," *The Very Large Databases J.*, vol. 12, no. 1, pp. 28-40, 2003.
- [7] Bell and F. Dravis, "Is Your Data Dirty? and Does that Matter?," *Accenture Whiter Paper*, <http://www.accenture.com>, 2006
- [8] R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, 1992.
- [9] M. de Almeida, M.A. Goncalves, M. Cristo, and P. Calado, *A Combined Component Approach for Finding Collection-Adapted Ranking Functions Based on Genetic Programming*, *Proc. 30th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval*, pp. 399-406, 2007.
- [10] Zhang, Y. Chen, W. Fan, E.A. Fox, M. Goncalves, M. Cristo, and P. Calado, *Intelligent pp Fusion from Multiple Sources for Text Classification* *Proc. 14th ACM Int'l Conf. Information and Knowledge Management* pp. 477-484, 2005
- [11] G. de Carvalho, M.A. Goncalves, A.H.F. Laender, and A.S. da Silva *Learning to Deduplicate*, *Proc. Sixth ACM/IEEECS Joint Conf. Digital Libraries*, pp. 41-50, 2006.
- [12] Bilenko and R.J. Mooney, *Adaptive Duplicate Detection Using Learn-*

- ableStringSimilarityMeasures, Proc. Ninth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 39-48, 2003.*
- [13] Lawrence, L. Giles, and K. Bollacker, *Digital Libraries and Autonomous Citation Indexing, Computer, vol. 32, no. 6, pp. 67-71, June 1999.*
- [14]. K. Elmagarmid, P. G. Ipeirotis, and V. S. Verykios, *Duplicate Record Detection: A Survey, IEEE Trans. Knowledge and Data Eng., vol. 19, no. 1, pp. 1-16, Jan. 2007*
- [15] W. Cohen, *Data Integration Using Similarity Joins and a Word-Based Information Representation Language, ACM Trans. Information Systems, vol. 18, no. 3, pp. 288-321, 2000*
- [16] C. P. Carvalho and A. S. da Silva, "Finding Similar Identities among Objects from Multiple Web Sources," *Proc. Fifth ACM Int'l Workshop Web Information and Data Management, pp. 90-93, 2003.*
- [17] B. Newcombe, J. M. Kennedy, S. Axford, and A. James, *Automatic Linkage of Vital Records, Science, vol. 130, no. 3381, pp. 954-959, Oct. 1959.*
- [18] *Freely Extensible Biomedical Record Linkage, <http://sourceforge.net/projects/febrl>, 2011*
- [19]. W. Cohen and J. Richman, *Learning to Match and Cluster Large High-Dimensional Data Sets for Data Integration, Proc. Ninth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 475-480, 2002*

International Journal of Innovative Technology & Adaptive Management (IJITAM)
ISSN: 2347-3622, Volume-2, Issue-7, April-2015

www.ijitam.org

International Journal of Innovative Technology & Adaptive Management (IJTAM)
ISSN: 2347-3622, Volume-2, Issue-7, April-2015

www.ijtam.org

International Journal of Innovative Technology & Adaptive Management (IJITAM)
ISSN: 2347-3622, Volume-2, Issue-7, April-2015

www.ijitam.org

International Journal of Innovative Technology & Adaptive Management (IJITAM)
ISSN: 2347-3622, Volume-2, Issue-7, April-2015

www.ijitam.org

International Journal of Innovative Technology & Adaptive Management (IJITAM)
ISSN: 2347-3622, Volume-2, Issue-7, April-2015

www.ijitam.org