

Multiple Parameters Clustering for Web Services

S.Rajendra Kumar, K.Chandra Shekar

Abstract— In the Web database scenario, the records to match are highly query-dependent. Since they can only be obtained through online queries. Moreover, they are only a partial and biased portion of all the data in the source Web databases. Consequently, hand-coding or offline-learning approaches are not appropriate for two reasons. First, the full data set is not available beforehand and therefore, good representative data for training are hard to obtain. Second, most importantly, even if good representative data are found and labeled for learning, the rules learned on the representatives of a full data set may not work well on a partial and biased part of that data set.

Index Terms— SOA, Web Services, Networks



I INTRODUCTION

Today, more and more databases that dynamically generate web pages in response to user queries are available on the Web. These Web databases compose the deep or hidden Web, which is estimated to contain a much larger amount of high quality, usually structured information and to have a faster growth rate than the static Web. Most Web databases are only accessible via a query interface through which users can submit queries. Once a query is received, the Web server will retrieve the corresponding results from the back-end database and return them to the user.

To build a system that helps users integrate and more importantly, compare the query results returned from multiple Web databases, a crucial task is to match the different source's records that refer to the same real-world entity. The problem of identifying duplicates, that is, two (or more) records describing the same entity, has attracted much attention from many research fields, including Databases, Data Mining, Artificial Intelligence and Natural Language Processing. Most previous work is based on predefined matching rules hand-coded by domain experts or matching rules learned offline by some learning method from a set of training examples. Such approaches work well in a traditional database environment, where all instances of the target databases can be readily accessed, as long as a set of high-quality representative records can be examined by experts or selected for the user to label.

II PREVIOUS WORK

Data integration is the problem of combining information from multiple heterogeneous databases. One step of data integration is relating the primitive objects that appear in the different databases specifically, determining which sets of identifiers refer to the same real-world entities. A number of recent research papers have addressed this problem by exploiting similarities in the textual names used for objects in different databases. (For example one might suspect that two objects from different databases named "USAMA FAYYAD" and "Usama M. Fayyad" respectively might refer to the same person). Integration techniques based on textual similarity are especially useful for databases found on the Web or obtained by extracting information from text, where descriptive names generally exist but global object identifiers are rare. Previous publications in using textual similarity for data integration have considered a number of related tasks. Although the terminology is not completely standardized, in this paper we define entity-name matching as the task of taking two lists of entity names from two different sources and determining which pairs of names are co-referent (i.e. refer to the same real-world entity).

We define entity-name clustering as the task of taking a single list of entity names and assigning entity names to clusters such that all names in a cluster are co-referent. Matching is important in attempting to join information pair of relations from different databases and clustering is important in removing duplicates from a relation that has been drawn from the union of many different information across of sources. Previous work in this area includes work in distance functions for matching and scalable matching and clustering algorithms. Work in record linkage is similar but does not rely as heavily on textual similarities. [1]

Important business decision; therefore, accuracy of such analysis is crucial. However, data received at the data warehouse from external sources usually contains errors: spelling mistakes, inconsistent conventions etc. Hence, significant amount of time and money are spent on data cleaning, the task of detecting and correcting errors in data. The problem of detecting and eliminating duplicated data is one of the major problems in the broad area of data cleaning and data quality. [2]

Many times, the same logical real world entity may have multiple representations in the data warehouse. For example, When Lisa purchases products from Super Mart twice, she might be entered as two different customers due to data entry errors. Such duplicated information can significantly increase direct mailing costs because several customers like Lisa may be sent multiple catalogs. Duplication elimination is hard because it is caused by several types of errors like type-graphical and equivalence error (non-unique and non standard) representation of the same logical value. Moreover, such duplicates can cause incorrect results in analysis queries (say, the number of Super Mart customers in customers in Seattle) and erroneous data mining models to be built. We refer to this problem of detecting and eliminating multiple distinct records representing the same real world entity as the fuzzy duplicate elimination problem, which is sometimes also called merge/purge, record linkage problems. This problem is different from the standard duplicate elimination problem, say for answering "select distinct" queries in relational database systems, which considers two tuples to be duplicates if they match exactly on all attributes. However elimination, which use duplicate elimination to mean fuzzy duplicate elimination. [3]

III. PROPOSED SYSTEM

A. Weighted Component Similarity Summing Classifier

In our algorithms, classifier plays a vital role. At the beginning, it is used to identify some duplicate vectors when there are no positive examples available. Then, after iteration begins, it is used again to cooperate with other classifier to identify new duplicate vectors. Because no duplicates vectors are available initially, classifiers that need class information to train, such as decision tree, cannot be used. An intuitive method to identify duplicate vectors is to identify duplicate vectors is to assume that two records are duplicates if most of their fields that are under consideration are similar. On the other hand, if all corresponding fields of the two records are dissimilar, it is unlikely that the two records are duplicates. To evaluate the similarity scenario between two records, we combine the values of each component in the similarity vector for the two records.

Different fields may have different importance, when we decide whether two records are duplicates. The importance is usually data-dependent, which, in turn depends on the query in the Web database.

B. Component Weight Assignment

In this classifier, we assign a weight to a component to indicate the importance of its corresponding field. The similarity between two duplicate records should be close to 1. For a duplicate vector that is formed by a pair of duplicate records r_1 and r_2 , we need to assign large weights to be components with large similarity values and small weights to the components with small similarity values. The similarity for two non-duplicate records should be close to 0. Hence, for a non-duplicate vector that is formed by a pair of non-duplicate records r_1 and r_2 , we need to assign small weights to the components with large similarity values and large weight to the components with small similarity values. The component will be assigned a small weight, if it usually has a small similarity value in the duplicate vectors.

C. Component Weight Assignment

Duplicate are defined based only some of the fields in a record. Using a Straight forward pre-processing step exact matching can merge those records that are exactly the same in all relevant matching fields. Similarity vector formed by a duplicate record pair a duplicate vector. After, we assign a weight for each component; the duplicate vector detection is rather intuitive. Two records r_1 and r_2 are duplicates if they are similar. i.e, if their similarity value is equal to or greater than a similarity threshold. In general, the similarity threshold should be close to 1 to ensure that the identified duplicates are correct. Increasing the value of similarity will reduce the number of duplicate vectors identified while, at the same time, the identified duplicates will be more precise.

D. Similarity Calculation

The similarity calculation quantifies the similarity between a pair of record fields. As the query results to match are extracted from HTML pages, namely, text files, we only consider string similarity. Given a pair of strings S_a and S_b , which must be between 0 and 1. Since the similarity function is orthogonal to the iterative duplicate detection, any kind of similarity calculation method can be employed. Domain knowledge or user preference can also be incorporated into the similarity function. In particular, the similarity function can be learned if training data is available.

IV. RESULT

The concepts of these papers are implemented and different results are show below. The proposed paper is implemented in Java technology on a Pentium-IV PC with 20 GB hard-disk and 512 MB RAM with apache web server. The propose paper's concepts shows efficient results and has been efficiently tested on different Datasets. The Fig. 1, Fig. 2, Fig. 3 and Fig. 4 shows the real time results compared.

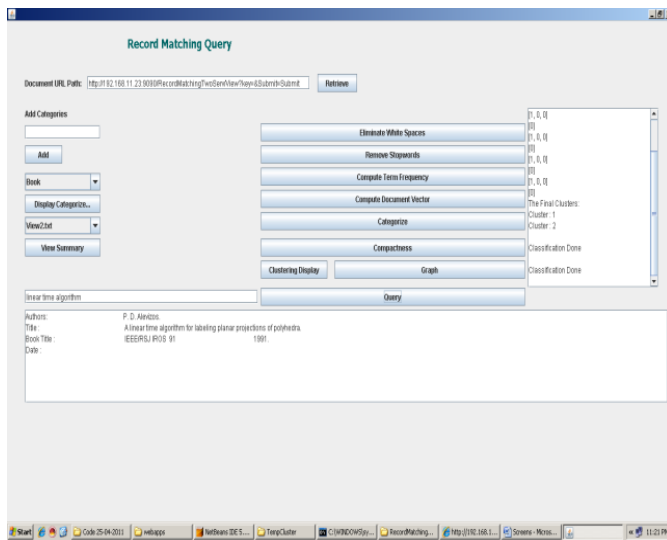


Fig. 1: Proposed system performing clustering on the dataset

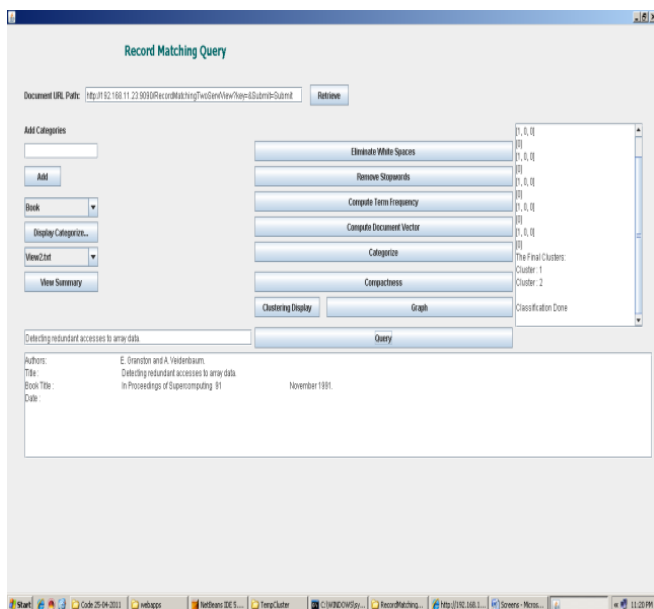


Fig. 2: Proposed system performing clustering on different dataset

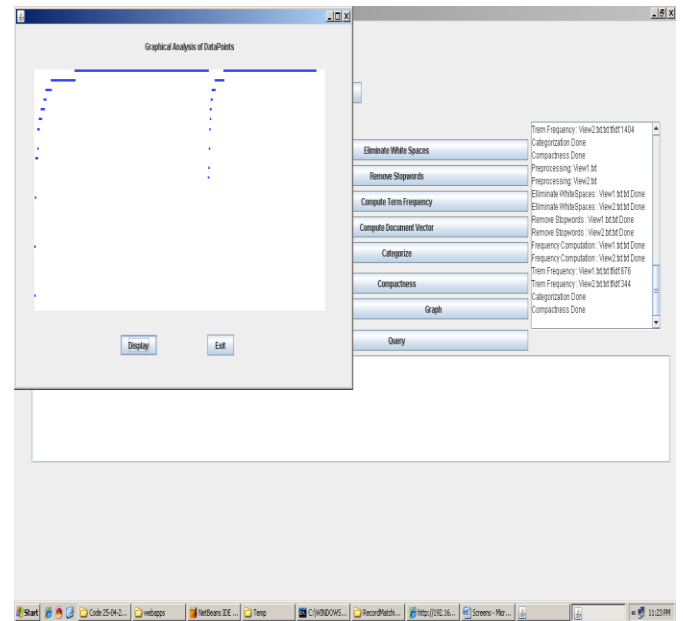


Fig. 3: Comparative Graphs of clustering

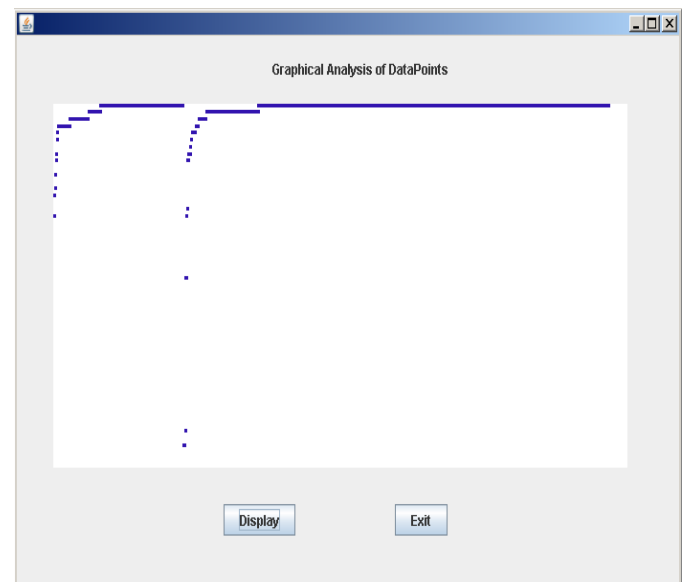


Fig. 4: Comparative Graphs of clustering

V. CONCLUSION

Duplicate detection is an important step in data integration and most state-of-the-art methods are based on offline learning techniques, which require training data. In the Web database scenario, where records to match are greatly query-dependent, a pre-trained approach is not applicable as the set of records in each query's results is a biased subset of the full data set. To overcome this problem, we presented an unsupervised, online approach, UDD for detecting duplicates over the query results of multiple Web database. Two classifiers, WCSS and SVM are used cooperatively in the convergence step of record matching to identify the duplicate pairs from all potential duplicate pair iteratively.

VI. ACKNOWLEDGMENT

We would to thank the anonymous referee for helpful comments.

VII. ACKNOWLEDGEMENT

- 1] W.W. Cohen and J.Richman, "Leaning to Match and Cluster Large High-Dimensional Datasets for Data Integration," Proc. ACM SIGKDD, pp. 475-480, 2002.
- 2] R. Ananthakrishna, S. Chaudhuri, and V. Ganti "Eliminating Fuzzy Duplicates in Data Warehouses," Proc. 28th Int'l Conf. Very Large Data Bases, pp.586-597, 2002.
- 3] S. Chaudhuri, K. Ganjam, V. Ganti ad R. Motwani, "Robust and Efficient Fuzzy Match for Online Data Cleaning." Proc. ACM SIGMOD, pp. 313-324, 2003.
- 4] L. Gravano, P.G. Iperiotis, H.V. Jagadish, N. Koudas, S. Muthukrishnan and D. Srivastava, "Approximate String Joins in a Database (Almost) for Free," Proc. 27th Int'l Conf. Very Large Data Bases, pp. 491-500, 2001.
- 5] X. Dong, A. Halevy, and J. Madhavan, "Reference Reconciliation in Complex Information Spaces," Proc. ACM SIGMOD, pp.85-96, 2005.
- 6] W.W.Cohen, H. Kautz, and D. McAllester, "Hardening Soft Information Source." Proc. ACM SIGKDD, pp.255-259, 2000.
- 7] P. Christen, T. Churches, and M. Hegianed, "Febrl-A Parallel Open Source Data Linkage System," Advances in Knowledge Discovery and Data Mining, pp 638-647, Springer, 2004.
- 8] P. Christen and K. Goiser, "Quality and Complexity Measures for Data Linkage and Deduplication Quality Measures in Data Mining. F. Guillet and H.Hamilton, eds., vol. 43, pp. 127-151, Springer, 2007.

AUTHOR SHORT BIOGRAPHIES:



Mr. Sanku Rajendra Kumar received his B.E. in Computer Science and Engineering from B.L.D.E's College of Engineering and Technology, Bijapur, Karnataka University and Pursuing M.Tech in Computer Science (Software Engineering) from Aurora's Technological And Research Institute, JNTU Hyderabad, India.

E-mail: rajendhar.sanku@gmail.com



Mr. K. Chandra Shekar working as Associate Professor in the Department of Computer Science and Engineering in Aurora's Technological and Research Institute with a teaching experience of 9 Year and worked as Project Trainee. He had received his Master Degree in Software Engineering. His areas of interest include Data mining and Information Security and his Specialization in Associative Classification Mining in Non-Binary Data.

E-mail: chandhra2k7@gmail.com