



## Mining E- Commerce for Customer Behaviour Analysis through Clustering

F. Gh. Kharrat<sup>1\*</sup>, Monaco F. J<sup>2</sup>, R. Amini. Niaki<sup>3</sup><sup>1</sup> USP - Interunidades Bioengenharia (EESC/FMRP/IQSC), Universidade de São Paulo, São Carlos, SP, Brazil<sup>2</sup> USP- Institute of Mathematics and Computer Sciences, Dept. of Computer Systemde São Paulo, São Carlos, SP, Brazil<sup>3</sup> USP-Mechanical Engineering, Escola de Engenharia de São Paulo , Sao Carlos, SP, Brazil

**Abstract-** In the Web mining 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 minings. To overcome this problem, we presented an unsupervised, online approach, Unsupervised Duplicate Detection (UDD), for detecting duplicates over the query results of multiple Web minings. Two classifiers, Weighted Component Similarity Summing Classifier (WCSS) and support vector machine (SVM), are used cooperatively in the convergence step of record matching to identify the duplicate pairs from all potential duplicate pairs iteratively with Weighted Component Similarity Summing Classifier, Component weight assignment and Duplicate Identification to achieve the efficient of the proposed system.

**Keywords-** SOA, Web Services, Networks, E-commerce

### I. INTRODUCTION

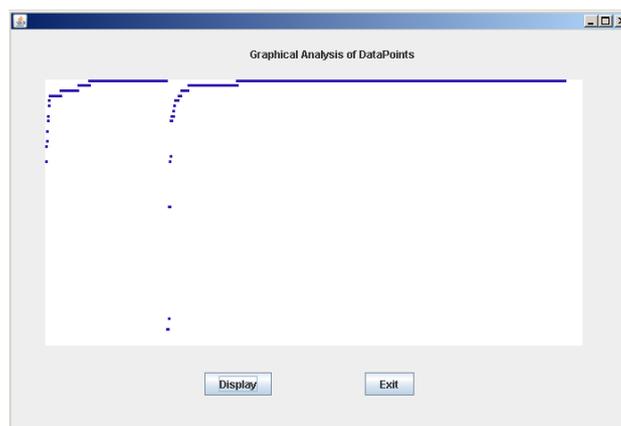
Today, more and more mining's that dynamically generate Web pages in response to user queries are available on the Web. This Web mining's compose the deep or hidden Web, which is estimated to contain a much larger amount of high quality, usually structured information and to have a faster growth rate than the static Web. Most Web mining's 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 mining 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 mining's, a crucial task is to match the different sources 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 mining's, 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 mining environment, where all instances of the target mining's 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 minings. One step of data integration is relating the primitive objects that appear in the different minings 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 minings. (For example one might suspect that two objects from different minings named "USAMA FAYYAD" and "Usama M. Fayyad" respectively might refer to the same person.) Integration techniques based on textual similarity are especially useful for minings 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 across of pair of relations from different minings, and clustering is important in removing duplicates from a relation that has been drawn from the union of many different information 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 decisions; 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





## V. PERFORMANCE ANALYSIS

The proposed paper is implemented in Java technology on a Pentium-III PC with 20 GB hard-disk and 256 MB RAM with apache web server. The propose paper's concepts shows efficient results and has been efficiently tested on different Messages.

## VI. 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 mining scenario, where records to match are greatly query-dependent, and a pertained 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 mining's. Two classifiers, WCSS and SVM, are used cooperatively in the convergence step of record matching to identify the duplicate pairs from all potential duplicate pairs iteratively.

## REFERENCES

- [1] W.W. Cohen and J. Richman, learning 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, and R. Motwani, "Robust and Efficient Fuzzy Match for Online Data Cleaning," Proc. ACM SIGMOD, pp. 313-324, 2003.
- [4] L. Gravano, P.G. Ipeirotis, H.V. Jagadish, N. Koudas, S. Muthukrishnan, and D. Srivastava, "Approximate String Joins in a Mining (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 Sources," Proc. ACM SIGKDD, pp. 255-259, 2000.
- [7] P. Christen, T. Churches, and M. Hegland, "Febri—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.