

A Survey on Duplicate Reduction in Graph Mining: Approaches, Analysis, and Evaluation

S. Raja¹, G. D. Nimalan², B. Sai Lakshman³, S. Sree Vishnu⁴

¹Assistant Professor, Department of CSE, Panimalar Institute of Technology, Chennai, India

^{2,3,4}Student, Department of CSE, Panimalar Institute of Technology, Chennai, India

Abstract: Utilizing on the web shopper surveys as electronic verbal exchange to help buy basic leadership has turned out to be progressively prevalent. The Web gives a broad wellspring of shopper audits, however one can scarcely read all surveys to acquire a reasonable assessment of an item or administration. A content preparing system that can condense surveys, would hence be attractive. A subtask to be performed by such a structure is locate the general angle classes tended to in survey sentences, for which this paper presents two techniques. As opposed to most existing methodologies, the principal strategy exhibited is an unsupervised technique that applies affiliation control mining on co-event recurrence information acquired from a corpus to discover these angle classifications. While not comparable to best in class directed strategies, the proposed unsupervised strategy performs superior to a few straightforward baselines, a Cloud computing provides flexible data management and ubiquitous data access. This paper takes such an approach in identifying the effect of duplicates on the performance of graph mining algorithms. Based on that observation, it proposes a number of heuristics to reduce the number of duplicates generated to significantly improve the performance of these algorithms. Further, we establish their correctness as well as their performance analysis for a number of graph characteristics. Based on these analysis, we show that it is possible to choose the best heuristic whether we have additional information about the graphs or not.

Keywords: Cloud Computing, Keyword frequency, Graph mining

1. Introduction

Data mining could be a method employed by corporations to show data into helpful data. By exploitation software system to appear for patterns in giant batches of knowledge, businesses will learn a lot of regarding their customers to develop simpler selling methods, increase sales and reduce prices. Substructure revelation is the route toward finding substructure(s) (a related subgraph) in a graph (or a woodlands) that best depicts a thought embedded in that graph reliant on a couple of criteria (repeat, compressibility et cetera.)

Many systems for substructure exposure have been proposed in the composition. Essential memory based circle based and database-arranged philosophies address substructure exposure on a single machine. With graphs that overwhelm essential memory, partition based approaches to manage substructure disclosure have substructures of growing sizes (starting from

substructure of size one that has one edge), checks the amount of specific undefined (or then again equivalent) substructures and applies a metric (e.g., repeat, Minimum Description Length, minimum help ET cetera.) to rank them. In each accentuation, either all-inclusive substructures or a subset (using the rank) are passed on forward to restrict the look space. This method is repeated until ensured substructure measure is come to or there are no more substructures to make.

In this paper, to propose associate economical giant universe regular language searchable coding theme for the cloud, that is privacy-preserving and secure against the off-line keyword idea attack (KGA). A notable highlight of the proposal over other existing schemes is that it supports the regular language encryption and deterministic finite automata (DFA) based data retrieval. The large universe construction ensures the extend ability of the system, during which the image set doesn't got to be predefined. Multiple users as supported within the system, and also the user might generate a DFA token exploitation his own personal key while not interacting with the key generation center. Furthermore, the concrete scheme is efficient and formally proved secure in standard model. Extensive comparison and simulation show that this theme has operate and performance superior than different schemes.

2. Literature survey

- *Yang Liu, Xiao Hong Jiang, Haunt Chen, Jun Ma, and Xiangyu Zhang:* Network motifs square measure basic building blocks in advanced networks. Motif detection has recently attracted a lot of attention as a subject to uncover structural style principles of advanced networks. Pattern finding is that most computationally costly step within the method of motif detection. during this paper, they style a pattern finding rule supported Google Map cut back to enhance the potency. Performance analysis shows our rule will facilitate the detection of giant r motifs in large size networks and has sensible measurability. They apply it within the prescription network and realize some usually used prescription network motifs that give the chance to any discover the law of prescription compatibility.
- *Saber Aridhi:* Big graph mining is an important

research area and it has attracted considerable attention. It allows to process, analyze, and extract meaningful information from large amounts of graph data. Big graph mining has been highly motivated not only by the tremendously increasing size of graphs but also by its huge number of applications. Such applications include bioinformatics, chemo informatics and social networks. One of the most challenging tasks in big graph mining is pattern mining in big graphs. This task consists on using data mining algorithms to discover interesting, unexpected and useful patterns in large amounts of graph data. It aims also to provide deeper understanding of graph data. In this context, several graph processing frameworks and scaling data mining/pattern mining techniques have been proposed to deal with very big graphs. This paper gives an overview of existing data mining and graph processing frameworks that deal with very big graphs. Then it presents a survey of current researches in the field of data mining / pattern mining in big graphs and discusses the main research issues related to this field. It also gives a categorization of both distributed data mining and machine learning techniques, graph processing frameworks and large scale pattern mining approaches.

- *Siddharth Suri Sergei Vassilvitskii*: The clustering coefficient of a node in a social network is a fundamental measure that quantifies how tightly-knit the community is around the node. Its computation can be reduced to counting the number of triangles incident on the particular node in the network. In case the graph is too big to fit into memory, this is a non-trivial task, and previous researchers showed how to estimate the clustering coefficient in this scenario. A different avenue of research is to perform the computation in parallel, spreading it across many machines. In recent years Map Reduce has emerged as a de facto programming paradigm for parallel computation on massive data sets. The main focus of this work is to give Map Reduce algorithms for counting triangles which we use to compute clustering coefficients. Our contributions are twofold. First, we describe a sequential triangle counting algorithm and show how to adapt it to the Map Reduce setting. This algorithm achieves a factor of 10-100 speed up over the naive approach. Second, we present a new algorithm designed specifically for the Map Reduce framework. A key feature of this approach is that it allows for a smooth tradeoff between the memory available on each individual machine and the total memory available to the algorithm, while keeping the total work done constant. Moreover, this algorithm can use any triangle counting algorithm as a black box and distribute the computation across many machines. We

validate our algorithms on real world datasets comprising of millions of nodes and over a billion edges. Our results show both algorithms effectively deal with skew in the degree distribution and lead to dramatic speed ups over the naive implementation.

- *Illinois Univ., Urbana, IL, USA*: We investigate new approaches for frequent graph-based pattern mining in graph datasets and propose a novel algorithm called gSpan (graph-based substructure pattern mining), which discovers frequent substructures without candidate generation. gSpan builds a new lexicographic order among graphs, and maps each graph to a unique minimum DFS code as its canonical label. Based on this lexicographic order gSpan adopts the depth-first search strategy to mine frequent connected subgraphs efficiently. Our performance study shows that gSpan substantially outperforms previous algorithms, sometimes by an order of magnitude.
- *Mansurul A. Bhuiyan and Mohammad Al Hasan*: Frequent sub graph mining (FSM) is an important task for exploratory data analysis on graph data. Over the years, many algorithms have been proposed to solve this task. These algorithms assume that the data structure of the mining task is small enough to fit in the main memory of a computer. However, as the real-world graph data grows, both in size and quantity, such an assumption does not hold any longer. To overcome this, some graph database-centric methods have been proposed in recent years for solving FSM; however, a distributed solution using Map Reduce paradigm has not been explored extensively. Since Map Reduce is becoming the de-facto paradigm for computation on massive data, an efficient FSM algorithm on this paradigm is of huge demand. In this work, we propose a frequent sub graph mining algorithm called FSM-H which uses an iterative Map Reduce based framework. FSM-H is complete as it returns all the frequent sub graphs for a given user-defined support, and it is efficient as it applies all the optimizations that the latest FSM algorithms adopt. Our experiments with real life and large synthetic datasets validate the effectiveness of FSM-H for mining frequent sub graphs from large graph datasets.
- *Cong Wang, Ning Cao, Jin Li, Kui Ren*: As Cloud Computing becomes prevalent, sensitive information are being increasingly centralized into the cloud. For the protection of data privacy, sensitive data has to be encrypted before outsourcing, which makes effective data utilization a very challenging task. Although traditional searchable encryption schemes allow users to securely search over encrypted data through keywords, these techniques support only Boolean search, without capturing any relevance of data files.

This approach suffers from two main drawbacks when directly applied in the context of Cloud Computing. On the one hand, users, who do not necessarily have pre- knowledge of the encrypted cloud data, have to post process every retrieved file in order to find ones most matching their interest; On the other hand, invariably retrieving all files containing the queried keyword further incurs unnecessary network traffic, which is absolutely undesirable in today’s pay-as-you-use cloud paradigm. In this paper, for the first time we define and solve the problem of effective yet secure ranked keyword search over encrypted cloud data. Ranked search greatly enhances system usability by returning the matching files in a ranked order regarding to certain relevance criteria (e.g., keyword frequency), thus making one step closer towards practical deployment of privacy-preserving data hosting services in Cloud Computing. We first give a straightforward yet ideal construction of ranked keyword search under the state-of-the-art searchable symmetric encryption (SSE) security definition, and demonstrate its inefficiency

- **Suresh Mariam Varghese:** Using Cloud Storage, users can remotely store their data and enjoy on-demand high quality applications and services. To ensure safety of stored data, it becomes must to encrypt data before storing in the global space. In cloud data, search arises only with plain data. But it is essential to invoke search with encrypted data. The specialty of cloud data storage is that it allows copious keywords in a solitary query and sorts the resultant data documents in relevance order. The proposed multi keyword search based on ranking over an encrypted cloud data uses feature of similarity and inner product similarity matching. The vector space model helps to provide sufficient search accuracy and homomorphic encryption enables users to involve in ranking while majority of computing work is done on server side by operations only on cipher text. Thus in this method for Top-K retrieval user gets an interested/used link in top. To selectively share documents fine-grained attribute-based access control policies can be used.

3. Proposed work

In this project we have to secure the file is the main motivation. In this, there is two parts are there one is user side and another one is admin side. In user side, only they will upload the data in the form of file. After that in an admin side, there are four admins are there. If the first user wants the file, they need acknowledgements of the other three members then only they will use the file else they are not accepting the file. The main motive is that, if the first user wants the file the other three member’s acknowledgement is very important then only the requester will use the file. Graphs are common data

structures used to represent / model real-world systems. Graph Mining is one of the arms of Data mining in which voluminous complex data are represented in the form of graphs and mining is done to infer knowledge from them.

A. Architecture diagram

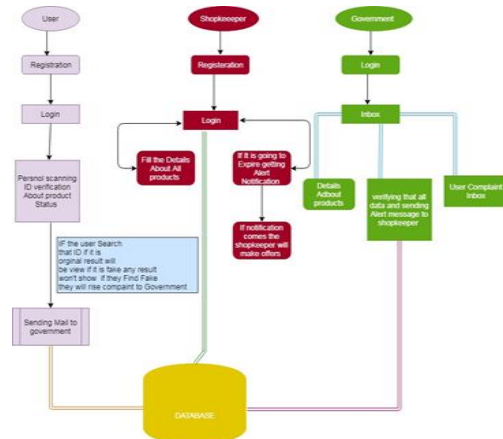


Fig. 1. Architecture diagram

B. Architecture explanation

Here in the above the concept each one of those issue first client need to keep up every one of the items with id. presently after login the businessperson account they need to transfer every one of the insights regarding items and they need to keep up make item and terminate date all they need to keep up in the wake of transferring all that these all data will goes to administrator group (carefulness group) now administrator group will deal with that all data and they can investigate and they will give all the data about the item lapsing date if the item will lapse they will send a notice to retailer before 15days of item will terminate. At that point businessperson will make offer for that specific id items then just it won't be squander capable that items. What's more, here client must be enroll one record they can login with id they can check with QR code whether that item is original item or not and on the off chance that it is original item, it will demonstrate the fabricate date and terminate date. in the event that it was phony it won't demonstrate any outcome. if like that any client discover like that they can send a mail. To administrator they can make a move on that specific shop.

4. Conclusion

In this paper, we studied the problem of record normalization over a set of matching records that refer to the same real-world entity. We presented three levels of normalization granularities (record-level, field-level and value component level) and two forms of normalization (typical normalization and complete normalization). For each form of normalization, we proposed a computational framework that includes both single-strategy and multi-strategy approaches. We proposed four single-strategy approaches: frequency, length, centroid, and feature-based to

select the normalized record or the normalized field value. We analyzed the record and field level normalization in the typical normalization. In the complete normalization, we focused on field values and proposed algorithms for acronym expansion and value component mining to produce much improved normalized field values. We implemented a prototype and tested it on a real-world dataset. The experimental results demonstrate the feasibility and effectiveness of our approach. Our method outperforms the state-of-the-art by a significant margin.

References

- [1] <http://ailab.wsu.edu/subdue>.
- [2] <http://snap.stanford.edu/data/com-LiveJournal.html>.
- [3] <http://snap.stanford.edu/data/ca-CondMat.html>.
- [4] <http://snap.stanford.edu/data/p2p-Gnutella04.html>.
- [5] Foto N. Afrati, Dimitris Fotakis, and Jeffrey D. Ullman. Enumerating subgraph instances using map-reduce. Technical report, Stanford University, December 2011.
- [6] Rakesh Agrawal and Ramakrishnan Srikant. Fast Algorithms for Mining Association Rules. In *Very Large Data Bases*, pages 487–499, 1994.
- [7] Md. Maksudul Alam, Maleq Khan, and Madhav V. Marathe. Distributed-memory parallel algorithms for generating massive scale-free networks using preferential attachment model. In *SC*, page 91, 2013.
- [8] Sofia Alexaki, Vassilis Christophides, Gregory Karvounarakis, and Dimitris Plexousakis. On Storing Voluminous RDF Descriptions: The Case of Web Portal Catalogs. In *International Workshop on the Web and Databases*, pages 43–48, 2001.
- [9] Björn Bringmann and Siegfried Nijssen. What is frequent in a single graph? In *PAKDD 2008*, Osaka, Japan, May 20–23, 2008 Proceedings, pages 858–863, 2008.
- [10] Aydın Buluc, Henning Meyerhenke, Ilya Safro, Peter Sanders, and Christian Schulz. Recent advances in graph partitioning. In *Algorithm Engineering*, pages 117–158. Springer, 2016.
- [11] Horst Bunke and Kim Shearer. A graph distance metric based on the maximal common subgraph. *Pattern Recognition Letters*, 19:255–259, 1998.
- [12] Deepayan Chakrabarti, Yiping Zhan, and Christos Faloutsos. R-MAT: A recursive model for graph mining. In *SIAM*, Florida, USA, April 22–24, 2004, pages 442–446, 2004.
- [13] Sharma Chakravarthy and Subhesh Pradhan. DB-FSG: An SQL-Based Approach for Frequent Subgraph Mining. In *DEXA*, pages 684–692, 2008.
- [14] Soumyava Das. Divide and Conquer Approach to Scalable Substructure Discovery: Partitioning Schemes, Algorithms, Optimization and Performance Analysis using Map/Reduce Paradigm. PhD thesis, The University of Texas at Arlington, May 2017.
- [15] Soumyava Das and Sharma Chakravarthy. Challenges and approaches for large graph analysis using map/reduce paradigm. In *BDA*, pages 116–132, 2013.
- [16] Soumyava Das and Sharma Chakravarthy. Partition and conquer: Map/reduce way of substructure discovery. In *DaWaK 2015*, Valencia, Spain, September 1–4, 2015, pages 365–378, 2015.
- [17] Soumyava Das, Ankur Goyal, and Sharma Chakravarthy. Plan before you execute: A cost-based query optimizer for attributed graph databases. In *DaWaK 2016*, Porto, Portugal, September 6–8, 2016, pages 314–328, 2016.
- [18] Mukund Deshpande, Michihiro Kuramochi, and George Karypis. Frequent Sub-Structure-Based Approaches for Classifying Chemical Compounds. In *IEEE International Conference on Data Mining*, pages 35–42, 2003.