



International Journal of Intellectual Advancements and Research in Engineering Computations

Enhanced aggregate estimation using left deep tree

Dr. P. Manimaran ME.,Ph.D¹, S.Vinitha², V. Sruthi³, M. Selva kumar⁴, K.S. Rangasamy College of Technology, Tiruchengode, Namakkal, Tamilnadu, India ^{2,3,4} Professor, B.E.-CSE, K.S. Rangasamy College of Technology, Tiruchengode, Namakkal, Tamilnadu, India¹ sivavinitha07@gmail.com

Abstract- A large number of web data repositories are hidden behind restrictive web interfaces, making it an important challenge to enable data analytics over these hidden web databases. In this propose system novel techniques which use a small number of queries to produce unbiased estimates with small variance. This paper presents real datasets demonstrate the accuracy and efficiency algorithms. Left-deep-tree data structure which imposes an order of all queries. Based on the order, it is capable of mapping each tuple in the hidden database to exactly one query in the tree, which is referred as the designated query. In this propose system novel techniques which use a small number of queries to produce unbiased estimates with small variance. These techniques can also be used for approximate query processing over hidden databases. Present theoretical analysis and extensive experiments to illustrate the effectiveness of proposed approach. In this paper initiate a study of estimating, without bias, the size and other aggregates over a hidden database.

1 INTRODUCTION

Data mining or knowledge discovery, is the computer-assisted process of digging through and analysing enormous sets of data and then extracting the meaning of the data. Data mining tools predict behaviours and future trends, allowing businesses to make proactive, knowledge-driven decisions. They scour databases for hidden patterns, finding

predictive information that experts may miss because it lies outside their expectations.

Data mining derives its name from the similarities between searching for valuable information in a large database and mining a mountain for a vein of valuable ore. Both processes require either sifting through an immense amount of material, or intelligently probing it to find where the value resides.

For businesses, data mining is used to discover patterns and relationships in the data in order to help make better business decisions. Data mining can help spot sales trends, develop smarter marketing campaigns, and accurately predict customer loyalty.

1.1 Enhanced aggregate estimation using left deep tree

The proposed methodology overlapping the queries in a left-deep-tree data structure which imposes an order of all queries. Based on the order, it is capable of mapping each tuple in the hidden database to exactly one query in the tree, which is referred as the designated query. By performing a drill-down based sampling process over the tree and testing whether a sample query is the designated one for its returned tuple(s), it develops an aggregate estimation algorithm that provides completely unbiased estimates for **COUNT**, **SUM** and **AVG** queries. In addition through the research analysis of a top-k restriction on the number of returned tuples, and a limit on the number of queries one can issue through the web interface. Also, cache results of previous queries are maintained in web server space per IP address per day.

1.2 OUTLINE OF TECHNICAL RESULTS

In this paper, we develop three main ideas for aggregate estimation over Enhanced Aggregate Estimation Using Left Deep Tree

UNBIASED-SET We start by showing a unique challenge imposed by the databases using left deep tree. For the databases using left deep tree, however, it is easily construct such a query tree because, unlike the hidden databases with drop-down-list interfaces, it is possible to pre-compute a set of overlapping queries which guarantee to return all tuples in this kind of hidden database. As a result, one has to rely on a set of overlapping queries to support aggregate estimation (e.g., through a query-sampling process)—which may lead to biased results because different tuples may be returned by different numbers of queries. Our first idea is to organize these overlapping queries in a left-deep-tree data structure which imposes an order of all queries. Based on the order, we are capable of mapping each tuple in the hidden database to exactly one query in the tree, which we refer to as the designated query.

WEIGHTED-SET-The error of an aggregate estimation consists of two components: bias and variance. We develop idea of weighted sampling to minimize variance. Specifically, we dynamically adjust the probability of sampling a query based on the query answers we receive so far, in order to “align” the sampling process to both the data distribution and the aggregate to be estimated, and thereby reduce the variance of our aggregate estimations.

CRAWL-Finally, we find that certain tuples in a hidden database, lowly ranked tuples that can only be returned by queries with a large number of conjunctive predicates—can cause a significant increase in the variance of aggregation estimations. To address the problem, we develop a special-case handling procedure which crawls such tuples to significantly reduce our final estimation error.

2 RELATED WORK

Cheng Sheng et al [1] describe a hidden database refers to a dataset that an organization makes accessible on the web by allowing users to issue queries through a search interface. Instead, data are obtained by querying the interface, and reading the result page dynamically generated. To extract all the tuples from a hidden database and algorithms are provably efficient, namely, they accomplish the task by performing only a small number of queries, even in the worst case. They also establish theoretical results indicating that these algorithms are asymptotically optimal i.e., it is impossible to improve their efficiency by more than a constant factor.

Michael Benedikt [3] verified about systems whose transitions consist of accesses to a Web-based data-source. An access is a lookup on a relation within a relational database, fixing values for a set of positions in the relation. AccLTL, is based on a first-order extension of linear-time temporal logic, interpreting access paths as sequences of relational structures. We also present a lower-level automaton model, they show that AccLTL and A-automata can express static analysis problems related to “querying with limited access patterns” that have been studied in the database literature in the past, such as whether an access is relevant to answering a query, and whether two queries are equivalent in the accessible data they can return.

Jayant Madhavan et al [4] describe The Deep Web refers to content hidden behind HTML forms. In order to get to such content, a user has to perform a form submission with valid input values. The Deep Web is also believed to be the biggest source of structured data on the Web and hence accessing its contents has been a long standing challenge in the data management community. The results of our surfacing are now shown in over 1000 web-search queries per-second, and the content surfaced is in over 45 languages and in hundreds of domains.

3 ESTIMATION ALGORITHM

A. Unbiased Estimation Algorithm

The Unbiased Estimation Algorithm (UEA) presents a solution for a novel problem: Aggregate Estimation for the Hidden Database with Checkbox Interface. The unbiased weighted-crawl which performs only random drill-downs on a novel structure of queries which refer to as a left-deep tree. The weight adjustment and low probability crawl for estimation accuracy.

Weight adjustment and low probability crawl is used for random drill-downs for estimation accuracy. The relative error occurred when the number of queries issued increases through the checkbox interfaces. The search results for the designated query are not controlled by the top-k restriction and it overflowing the search results.

B. Fast Moving Item Analysis (FMIA)

The proposed methodology overlapping the queries in a left-deep-tree data structure which imposes an order of all queries. Based on the order, it is capable of mapping each tuple in the hidden database to exactly one query in the tree, which is referred as the designated query. In addition through the research analysis of a top-k restriction on the number of returned tuples, and a limit on the number of queries one can issue through the web interface. In this paper also concentrates in finding the most favorable Fast Moving Item Analysis (FMIA) that got sale and also got the order for the purchase of the particular item in the market. With these particular details it is easier for the producers in finding the more frequent item sale in particular area and the particular fast moving item easily from particular time to time.

To combine the three techniques to produce an efficient and unbiased estimator for the hidden database size.

Algorithm: Left Deep Tree

1. $q \leftarrow$ root node. $p \leftarrow 1$. $i \leftarrow 1$.
2. Randomly generate $v \in \{0, 1\}$.
3. Issue $q \ 0 \leftarrow q \wedge (A_i = v)$.
4. if $q \ 0$ underflows then
5. $q \leftarrow q \wedge (A_i = 1 - v)$. Goto 2.
6. else if $q \ 0$ overflows then
7. Issue $q \wedge (A_i = 1 - v)$.
8. if $q \wedge (A_i = 1 - v) \neq 0$
9. $p \leftarrow p/2$. . Update $p(q)$
10. end if
11. $q \leftarrow q \ 0$. $i \leftarrow i + 1$. Goto 2.
12. end if
13. return $m \sim \leftarrow |q| / p$
14. Return an estimation database size

A. Left Deep Tree Construction

This output is used to enabling the aggregate queries over a hidden database with checkbox interface by issuing a small number of queries (sampling) through its web interface. In this output form, an aggregate estimation algorithm is implemented and used that provides completely unbiased estimates for COUNT and SUM queries.

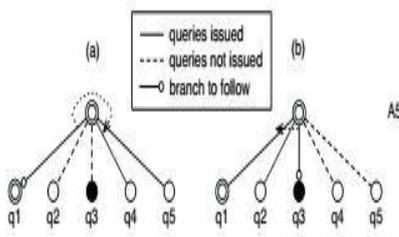
B. Aggregate Estimation with Fast Moving Item

In this aggregate estimation with fast moving item, in this form every node is corresponding to a query and a directed edge from a node to a child node indicates that the query corresponding to this child node includes all attributes item in the parent query and one additional attribute item.

4 BACKTRACKING FOR CATEGORICAL DATA

First, Boolean attributes ensure that the sibling of an under flowing node always overflow. There is no such guarantee for categorical databases. Thus, to successfully backtrack from an under flowing branch, we must find one of its sibling branches that returns non-empty and count the number of such non-empty siblings (in order to compute $p(q)$). A Novel problem a non-empty branch always exists given an overflowing parent node.

The other change required is the computation of $p(q)$. If the above-mentioned simple backtracking is used.



Backtracking diagram

computation of $p(q)$ becomes $p(q) = 1/Q_{hi} - 1 - c_i$, where c_i is the number of non-underflowing branches for the i th predicate then route to the top-valid query q . The two changes for categorical databases do not affect the unbiasedness of the estimation as the proof of Theorem 1 remains unchanged. These changes, however, do affect the query cost of the drill-down process. In particular, if the above simple backtracking technique is used, we must issue queries corresponding to all branches to find the number of non-empty ones, leading to a high query cost for large-fan-out attributes.

To reduce such a high query cost, to develop smart backtracking which aims to avoid testing all branches of a high fan-out attribute. Consider a categorical attribute A_i with possible values v_1, \dots, v_w ($w = |\text{Dom}(A_i)|$). Assume a total order of the values which can be arbitrarily assigned for cardinal or ordinal attributes. In the following, we describe the random drill-down and the computation of $p(q)$ with smart backtracking, respectively.

5 RESULTS AND DISCUSSION

The following table 6.1 shows the experimental results for Left Deep Tree and Fast moving Item base Left Deep Tree model. The table contains number of dataset query, average of left Deep Tree query and Average of Fast Moving item tree query.

Find the aggregate estimation query is following formula estimated,

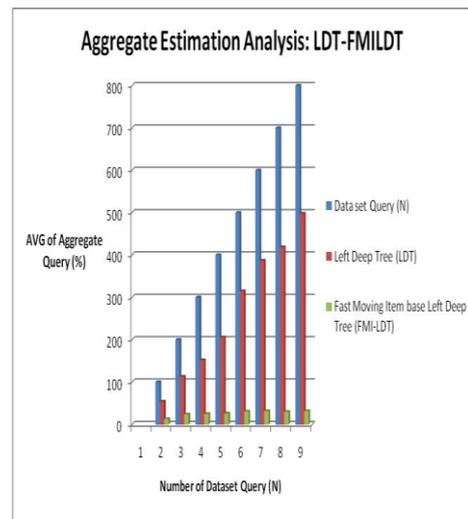
$$LDT = qn \sum_{i=3} n/i$$

$$FMILDT = qn \sum 2 * n/3$$

| s.no | Data set query | Left deep tree | Fast moving item using left deep tree |
|------|----------------|----------------|---------------------------------------|
| 1 | 100 | 54 | 13 |
| 2 | 200 | 113 | 24 |
| 3 | 300 | 152 | 25 |
| 4 | 400 | 206 | 26 |
| 5 | 500 | 315 | 31 |
| 6 | 600 | 387 | 32 |
| 7 | 700 | 419 | 30 |
| 8 | 800 | 498 | 32 |

Table 6.1 Aggregate Estimation Analysis: LDT-FMILDT

The following figure 6.2 shows the experimental results for Left Deep Tree and Fast moving Item base Left Deep Tree model. The figure contains number of dataset query, average of left Deep Tree query and Average of Fast Moving item tree query



6:2Fast moving Item base Left Deep Tree model

6 CONCLUSION

In this paper we have initiated an investigation of the unbiased estimation of the size and other aggregates over hidden web databases through its restrictive web interface. We proposed backtrack-enabled random walk schemes over the query space to produce unbiased estimates for SUM, COUNT and AVG queries, including the database size. In this fast moving item analysis system we also described the two ideas, weight adjustment and capture & recapture, to reduce the estimation variance. The proposed system provides theoretical analysis for estimation accuracy and query cost of the proposed ideas. We described a comprehensive set of experiments that demonstrate the effectiveness of the proposed approach over real-world hidden databases.

7 FUTURE ENHANCEMENTS

The paper has the scope for probing the hidden databases since query probing techniques have been widely used in the hidden database. The application becomes useful if the below enhancements are made in the future. In this area, there are three key related subareas: Resource discovery, i.e., the discovery of hidden database URLs from the web, Interface understanding, i.e., the proper understanding of how to issue (supported) search queries through a web interface and to extract query answers from the returned web pages,

REFERENCE

- [1] C. Sheng, N. Zhang, Y. Tao, and X. Jin, —Optimal algorithms for crawling a hidden database
- [2] Monster, Job search page
- [3] Epicurious, Food search page [Online]. recipesmenus/advancedsearch, 2013.
- [4] R. Khare, Y. An, and I.-Y. Song, “Understanding deep web search interfaces: A survey,” ACM SIGMOD Rec., vol. 39,
- [5] A. Dasgupta, X. Jin, B. Jewell, N. Zhang, and G. Das, Unbiased estimation of size and other aggregates over hidden web databases, in Proc.
- [6] M. Benedikt, G. Gottlob, and P. Senellart, —Determining relevance of accesses at runtime, in Proc.
- [7] M. Benedikt, P. Bourhis, and C. Ley, —Querying schemas with access restrictions, Proc. VLDB Endowment, vol. 5, no. 7, pp. 634–645, 2012
- [8] R. Khare, Y. An, and I.-Y. Song, —Understanding deep web search interfaces: A survey, ACM SIGMOD Rec.
- [9] B. He, M. Patel, Z. Zhang, and K. C.-C. Chang, “Accessing the deep web,” Commun. ACM, vol. 50, no. 5, pp. 94–101, 2007.
- [10] J. Madhavan, L. Afanasiev, L. Antova, and A. Halevy, “Harnessing the deep web: Present and future,” CoRR, vol. abs/0909.1785, 2009.