

ENHANCE ESTIMATION IN HIDDEN DATABASES USING CACHE MEMORY WITH LEFT DEEP TREE

S.Nandagopal¹, S.Jegadeesan², N.Priyadarshan³, L.Sridhar⁴, K.A.Sujith Kumar⁵,

¹Professor, Department of Information Technology,
Nandha College Of Technology, Erode.

^{2,3,4,5} UG Scholar, Department of Information Technology,
Nandha College of Technology, Erode.

¹asnandu@yahoo.com, ²jegadeesan1506@gmail.com, ³priyadarshan211@gmail.com
⁴sridharraja805@gmail.com, ⁵sujithka1312@gmail.com

ABSTRACT

HIDDEN databases are data repositories “hidden behind” only accessible through a restrictive web search interface. Input capabilities provided by such a web interface range from a simple keyword-search textbox to a complex combination of textboxes, drop-down controls, checkboxes.

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. Most existing techniques assume a form-like web interface which consists solely of categorical attributes. Nonetheless, many real-world web interfaces (of hidden databases) also feature checkbox interfaces—e.g., the specification of a set of desired features, such as A/C, navigation, etc., for a car-search website like Yahoo! Autos. It is found that, for the purpose of data analytics, such checkbox-represented attributes differ fundamentally from the categorical/numerical ones that were traditionally studied.

In this study, the problem of data analytics over hidden databases with checkbox interfaces is analyzed. Extensive experiments on both synthetic and real datasets demonstrate the accuracy and efficiency of the proposed algorithms

INTRODUCTION

In this paper, we develop novel techniques for estimating and tracking various types of aggregate queries, e.g., COUNT and SUM, over dynamic web databases that are hidden behind proprietary search interfaces and frequently changed. Hidden Web Databases: Many web databases are “hidden” behind restrictive search interfaces that allow a user to specify the desired values for one or a few attributes (i.e., form a conjunctive search query), and return to the user a small number (bounded by a constant k which can be 50 or 100) of tuples that match the user-specified query, selected and ranked according to a proprietary scoring function. Examples of such databases include Yahoo! Autos, Amazon.com, eBay.com, CareerBuilder.com, etc. Problem Motivations: The problem we consider in this paper is how a third party can use the restrictive web interface to estimate and track aggregate query answers over a dynamic web database. Aggregate queries are the most common type of queries in decision support systems as they enable effective analysis to glean insights from the data.

Tracking the number of tuples in a web database is by itself an important problem. For example, the number of

active job postings at Monster.com or listings at realestate.com can provide an economist with real-time indicators of US economy. Similarly, tracking the number of apps in Apple’s App Store and Google Play provides us a continuous understanding of the growth of the two platforms. Note that while some web databases (e.g., App Store) periodically publish their sizes for advertisement purposes, such published size is not easily verifiable, and some times doubtful because of the clear incentive for database owners to exaggerate the number

- To apply a novel problem of aggregate estimations over the hidden web databases with checkbox interfaces.
- To produce unbiased aggregate estimations over the hidden databases with checkbox interfaces and develop the data structure of left-deep-tree.
- Define the concept of designated query to form an injective mapping from tuples to queries supported by the web interface.

- To reduce the variance of aggregate estimations, we develop the ideas of weighted sampling and special tuple-crawling.

The proposed system main contributions also include a comprehensive set of experiments which demonstrate the effectiveness of enhanced UNBIASED-WEIGHTED-RAWL algorithm on aggregate estimation over real world hidden data-bases with checkbox interface, as well as the effectiveness of each of three ideas on improving the performance of UNBIASED-WEIGHTED-CRAWL.

RELATED WORKS

CHENG SHENG [1] stated that 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. In other words, data acquisition from such a source is not by following static hyper-links. Instead, data are obtained by querying the interface, and reading the result page dynamically generated. This, with other facts such as the interface may answer a query only partially, has prevented hidden databases

from being crawled effectively by existing search engines.

This paper remedies the problem by giving algorithms to extract all the tuples from a hidden database. Our 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. The derivation of our upper and lower bound results reveals significant insight into the characteristics of the underlying problem. Extensive experiments confirm the proposed techniques work very well on all the real datasets examined.

A. DASGUPTA [2] stated about the growth interest in random sampling from online hidden databases. These databases reside behind form-like web interfaces which allow users to execute search queries by specifying the desired values for certain attributes and the system respond by returning a few tuples that satisfy the selection conditions, sorted by a suitable scoring function. In this paper, they considered the problem of uniform random sampling over such hidden databases. A key challenge is to eliminate

the skew of samples incurred by the selective return of highly ranked tuples. To address this challenge, all state-of-the-art samplers share a common approach: they do not use overflowing queries.

SRIRAM RAGHAVAN [3]

crawlers retrieve content only from the publicly index able Web, i.e., the set of Web pages reachable purely by following hypertext links, ignoring search forms and pages that require authorization or prior registration. In particular, they ignore the tremendous amount of high quality content “hidden” behind search forms, in large searchable electronic databases. They address the problem of designing a crawler capable of extracting content from this hidden Web

EXISTING SYSTEM

The existing presents a solution for a novel problem: Aggregate Estimation for the Hidden Database with Checkbox Interface. In the hidden database with checkbox interface, a checkbox attribute is represented as a checkbox in the web interface. For example, in the home search website, features (e.g., central air, basement) for a home are represented by checkboxes. The checkbox interface has

its specialty. By checking the checkbox corresponding to a value v_1 , it ensures that all returned tuples contain the value v_1 . But it is impossible to enforce that no returned tuple contains v_2 —because unchecking v_2 is interpreted as “do-not-care” instead of “not-containing- v ” in the interface.

- If one considers a feature as a Boolean attribute, then the checkbox interface places a limitation that only TRUE, not FALSE, can be specified for the attribute.
- As a result, it is impossible to apply the existing techniques which require all values of an attribute to be specifiable through the input web interface.
- Such databases also have the same limitations as the hidden databases with drop-down-list interface. i.e., (1) a top-k restriction on the number of returned tuples, and (2) a limit on the number of queries one can issue (e.g., per IP address per day) through the web interface.
- Cache results of previous queries are not maintained in web server space and so the

burden of database server is more.

PROPOSED SYSTEM

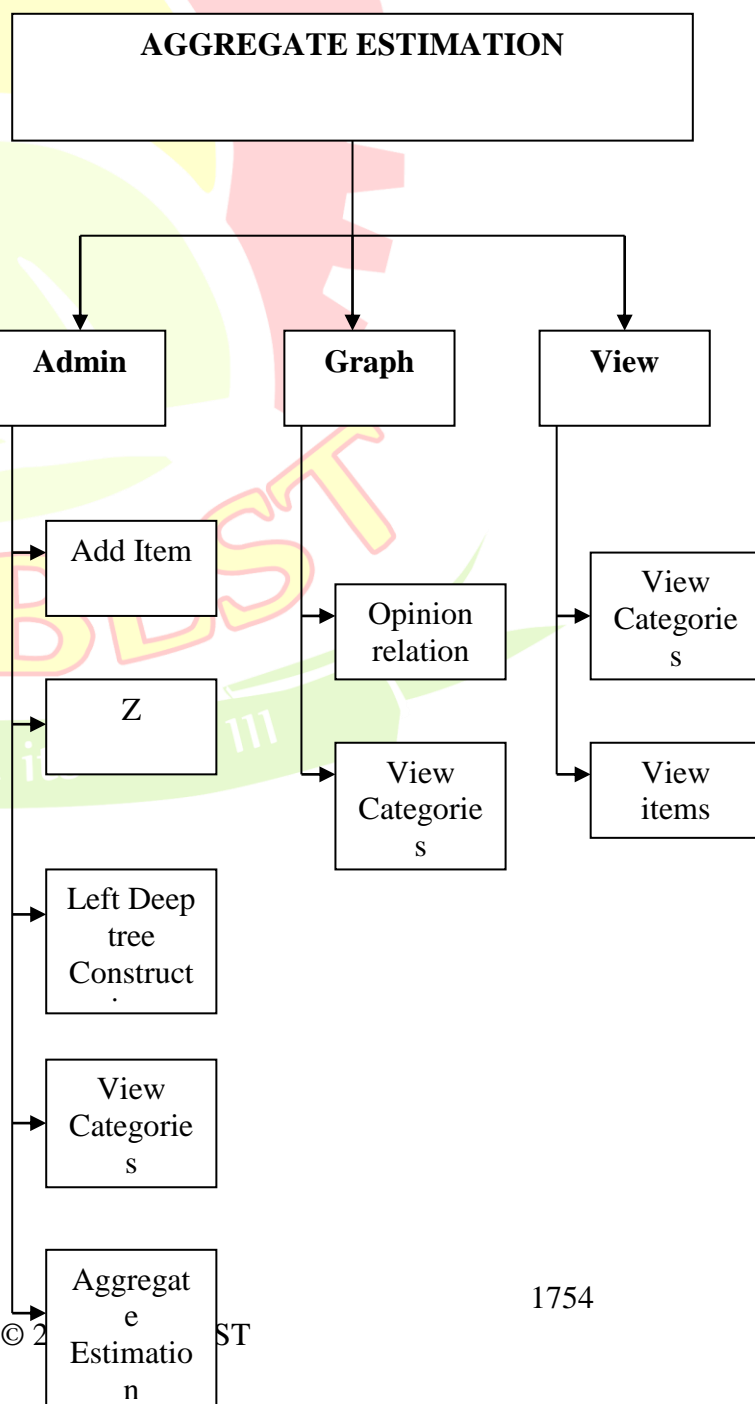
Like existing system, the first idea is to organize these overlapping queries in a left-deep-tree data structure which imposes an order of all queries. Based on the order, it is capable of mapping each tuple in the hidden database to exactly one query in the tree, which is referred as the designated 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 and SUM queries.

In addition, (1) a top-k restriction on the number of returned tuples, and (2) a limit on the number of queries one can issue (e.g., per IP address per day) through the web interface. Also, cache results of previous queries are maintained in web server space

- A top-k restriction on the number of returned tuples.
- A limit on the number of queries one can issue (e.g.,

per IP address per day) through the web interface.

- Cache results of previous queries are maintained in web server space and so eliminated the burden of database server.



S. No	Number Of Cluster	Cluster A	Cluster B	Cluster C	Cluster D	Cluster E	Cluster F
1	2	73	56	72	74	72	65
2	3	62	73	73	75	69	74
3	4	69	77	77	76	77	69
4	5	65	77	77	76	77	72
5	6	69	77	77	76	77	72
6	7	70	77	77	76	77	72
7	8	72	77	77	76	77	72
8	9	78	78	65	79	72	74
	No. of Aggregated data	558	574	570	542	566	563
	Average %	69.75	71.75	71.25	67.75	70.75	70.375

The following **Table 6.2** describes experimental result for existing system over all experimental result analysis. The table contains aggregated cluster, number of aggregated data cluster data and average aggregated data details are shown

Aggregate d Cluster	No. of. Aggregated Data	AVG % Aggregated
Cluster A	558	69.75
Cluster B	574	71.75
Cluster C	570	71.25
Cluster D	542	67.75
Cluster E	566	70.75
Cluster F	563	70.375

Table 6.2 Overall Experimental Result - Existing System

The following **Table 6.3** describes experimental result for proposed system performance rate analysis. The table contains number of cluster, cluster size and number of aggregated data and average aggregated data details are shown

The following **Table 6.4** describes experimental result for existing system over all experimental result analysis. The table contains aggregated cluster, number of aggregated data cluster data and average

Aggregated Cluster	No.Of. Aggregated Data	AVG % Aggregated
Cluster A	580	72.5
Cluster B	597	74.62
Cluster C	578	72.25
Cluster D	557	69.62
Cluster E	579	72.37
Cluster F	569	71.12

aggregated data details are shown

Table 6.4 Overall Experimental Result - Proposed System

The following **Fig 6.1** describes experimental result for existing system aggregated data analyses are shown.

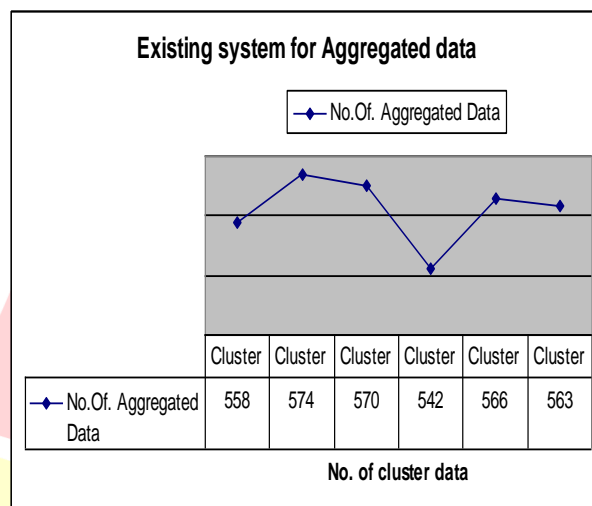


Fig 6.1 Existing System - Aggregated Data

The following **Fig 6.2** describes experimental result for proposed data transfers for hybrid method rate analyses are shown.

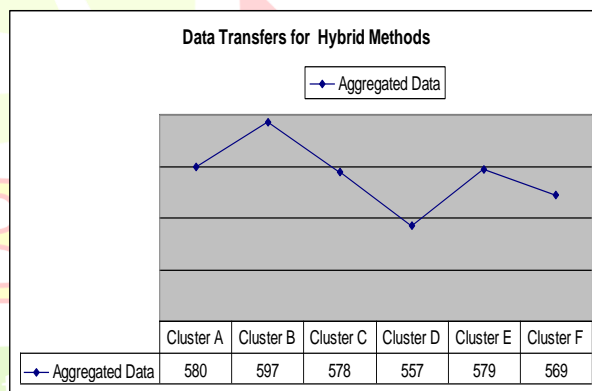


Fig 6.2 Proposed System - Aggregated Data

The following **Fig 6.2** describes experimental result for existing system aggregation scheme analyses are shown

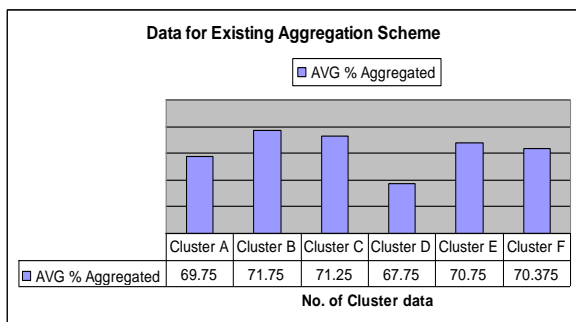


Fig 6.3 Aggregation Scheme- Existing System

The following Fig 6.4 describes experimental result for proposed system aggregation scheme analyses are shown

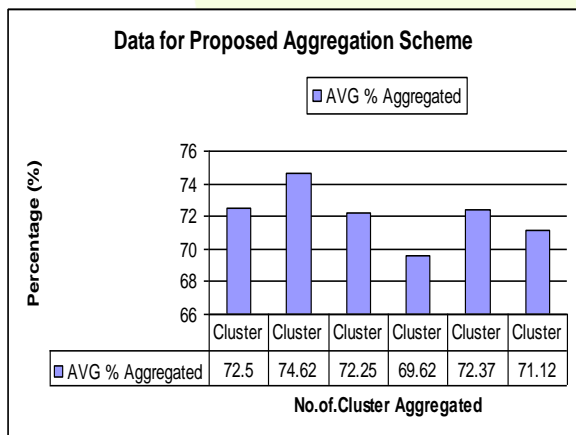


Fig 6.3 Aggregation Scheme- Proposed System

The following Table 6.4 describes experimental result for comparison between existing and proposed aggregation data scheme system. The table contains number of cluster, cluster size and number of aggregated data and average aggregated data in existing and proposed system detail are shown

CONCLUSION

This project addresses a novel problem where checkboxes exist in the web interface of a hidden database. To enable the approximation processing of aggregate queries and develops algorithm UNBIASED-WEIGHTED-CRAWL which performs random drill-downs on a novel structure of queries which we refer to as a left-deep tree and also propose weight adjustment and low probability crawl to improve estimation accuracy.

This project performed a comprehensive set of experiments on synthetic and real-world datasets with varying database sizes (from 5;000 to 100;000), number of attributes (from 20 to 50) and top-k restriction (from k = 10 to 30). We found that, as predicted by the theoretical analysis, the relative error decreases when the number of queries issued increases.

In addition, for the same query budget, the relative error is lower with a smaller number of attributes and/or a large k. In the worst-case scenario, achieve around 25 percent relative error with 300 queries issued for the synthetic dataset,

and less than 10 percent relative error with about 2,500 queries issued for the real-world dataset. The experimental results demonstrate the effectiveness of proposed algorithms.

SCOPE FOR FURTHER ENHANCEMENTS

The paper has the scope for probing the hidden databases since query probing techniques have been widely used in the hidden database. In this area, there are three key related subareas: (1) resource discovery, i.e., the discovery of hidden database URLs from the web, (2) 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, (3) crawling, sampling and data analytics over hidden web databases, which is the most related to our problem. The paper provides a best assistance in the network environment. It allows adding up the following facilities in future. The paper provides an enabling analytics on hidden web database which is a problem that has drawn much attention in proposed system. The application become useful if the below enhancements are made in future.

- If the application is designed as web service, it can be integrated in many network applications.
- The application is developed such that above said enhancements can be integrated with current modules.

REFERENCES

- [1] X. Jin, Y. Tao, C. Sheng, N. Zhang, and “Optimal algorithms for crawling a hidden database in the web,” Proc. VLDB Endowment, vol. 5, no. 11, pp. 1112–1123, 2012.
- [2] Monster, Job search page [Online]. Available: <http://jobsearch.monster.com/AdvancedSearch.aspx>, 2011.
- [3] Epicurious, Food search page [Online]. Available: <http://www.epicurious.com/recipesmenus/advancedsearch>, 2013.
- [4] Homefinder, Home finder page [Online]. Available: <http://www.homefinder.com/search>, 2013.
- [5] G. Das, A. Dasgupta, B. Jewell, N. Zhang, and “Unbiased estimation of size and other aggregates over hidden web

databases,” in Proc. Int. Conf.Manage.
Data, 2010, pp. 855–866.

[6] M. Benedikt, G. Gottlob, and P.
Senellart, “Determining relevance of
accesses at runtime,” in Proc. 30th ACM
SIGMOD-SIGACT-SIGART Symp.
Principles Database Syst., 2011, pp. 211–
222.

[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] Y. An, R. Khare, and I.-Y. Song,
“Understanding deep web search
interfaces: A survey,” ACM SIGMOD
Rec., vol. 39, no. 1, pp. 33–40, 2010

