

Retrieving Missing Data Based on Key Levels

J. Catherine Princy, V. S. Priyanga, P. Sailaja

Information Technology, Velammal Institute of Technology, Panchetti, Tamilnadu, India

ABSTRACT

Data imputation method is to fill the missing values from a group of datasets. Existing imputation approaches to non-quantities string knowledge may be roughly placed into two categories: 1. Inferring based approaches and 2. Retrieving primarily based approaches. Specifically, the inferring-based approaches notice substitutes or estimations for the missing ones from the entire part to the information set. However, they usually come short in filling in distinctive missing attribute values that don't exist in the complete part of the information set. During this project we tend to investigate the interaction between the inferring based methods and also the retrieving based approaches. We tend to show that retrieving a tiny low variety of selected missing values will highly improve the imputation recall of the inferring based ways. With this institution, we tend to propose associate interactive Retrieving-Inferring knowledge imputation approach, that performs retrieving and inferring alternately in filling missing attribute in an exceedingly datasets to confirm the high recall at the minimum values. This approach faces a challenge of choosing the smallest amount of variety of missing values for retrieving to maximize the amount of inferable values.

Keywords: Data Imputation, Data Repairing, Interactive Retrieving-Inferring.

I. INTRODUCTION

Data incompleteness is a data quality problem in all kinds of databases. Data Imputation is the process of filling in missing attribute values. So far, plenty of imputation techniques have been developed for missing quantitative data, which is either continuous data such as temperature, salary, age, etc. Data with a relatively small scope of values such as weather, gender, country, etc. Only limited attention has been paid to nonquantitative data which is pure string data with a large scope of values, such as email, phone, company, address, etc. However, pure string data takes up a large part of the missing data in many databases.

Specifically, the inferring-based approaches find substitutes or estimations for the missing ones from the complete part of the data set. However, they typically fall short in filling in unique missing attribute values which do not exist in the complete part of the data set. The retrieving based approaches resort to external resources for help. Based on the premise that the missing data might be available at some external data sources over the World Wide Web, some work has been conducted to harvest missing values from web lists and web tables.

II. METHODS AND MATERIAL

A. Existing System Principle

Existing imputation approaches to nonquantitative string data can be roughly put into two categories: (1) inferring- approaches and (2) retrieving-based

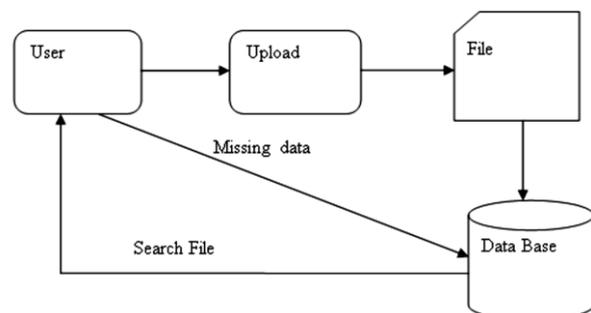


Figure 1. Existing System

B. Drawbacks

1. Low imputation recall.
2. Will not retrieve data

C. Proposed System

Retrieving-based methods are proposed to retrieve missing values from external resources such as the World Wide Web, which tend to reach a much higher imputation recall, but inevitably bring a large overhead by issuing a large number of search queries. In this paper, we investigate the interaction between the inferring based methods and the retrieving-based methods. We show that retrieving a small number of selected missing values can greatly improve the imputation recall of the inferring based methods. With this intuition, we propose an interactive Retrieving-Infering data imputation approach (TRIP), which performs retrieving and inferring alternately in filling in missing attribute values in a dataset. To ensure the high recall at the minimum cost, TRIP faces a challenge of selecting the least number of missing values for retrieving to maximize the number of inferable values. Our proposed solution is able to identify an optimal retrieving inferring scheduling scheme in Deterministic Data Imputation (DDI), and the optimality of the generated scheme is theoretic retrieving-based methods are proposed to retrieve missing values from external resources such as the optimal scheme is not feasible to be achieved analyzed with proofs. We also analyze with an example that the optimal scheme is not feasible to be achieved.

D. Advantages

1. High imputation recall.
2. Both insert and retrieve data

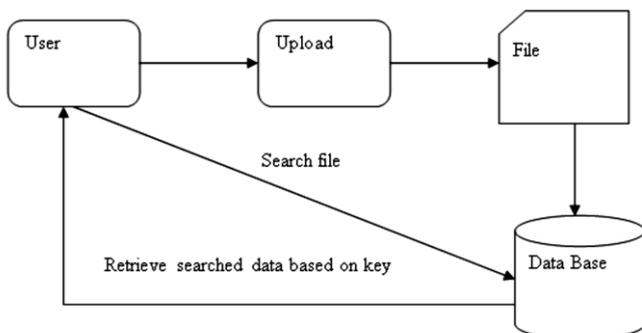
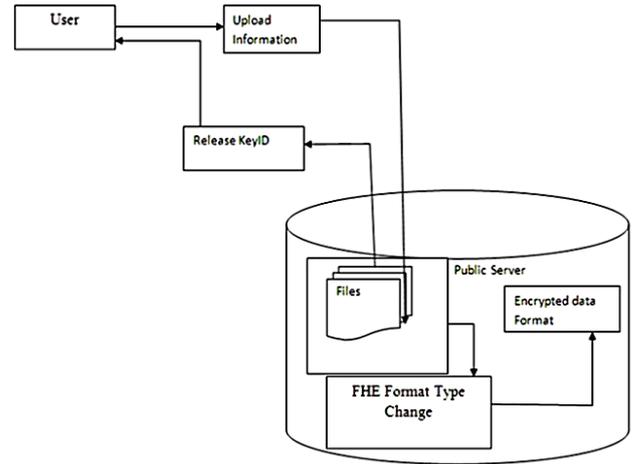


Figure 2. Proposed System

E. Architecture Diagram



F. Algorithm

Algorithm : Identifying an Optimal Scheme in DDI

Input : An incomplete table with missing value set \mathcal{O}

Output: An imputation scheme $\mathcal{S} = (\mathcal{I}_0, \mathcal{R}_1, \mathcal{I}_1, \dots, \mathcal{R}_n, \mathcal{I}_n)$

```

Set  $i = 0$ ;
while  $\mathcal{O} \neq \emptyset$  do
  1.  $\mathcal{I}_i \leftarrow$  All inferable values that can be inferred now;
  2.  $\mathcal{O} = \mathcal{O} - \mathcal{I}_i$ ;
  3. Infer all missing values in  $\mathcal{I}_i$ ;
  4.  $i++$ ;
  5. Build an inference dependency graph;
  6.  $\mathcal{R}_i \leftarrow$  Values in identified un-inferable nodes;
  7.  $\mathcal{R}_i \leftarrow \mathcal{R}_i \cup$  Values in identified minimum unlocking nodes;
  8.  $\mathcal{O} = \mathcal{O} - \mathcal{R}_i$ ;
  9. Retrieve all missing values in  $\mathcal{R}_i$ ;
return  $(\mathcal{I}_0, \mathcal{R}_1, \mathcal{I}_1, \dots, \mathcal{R}_n, \mathcal{I}_n)$ ;
  
```

G. Future Enhancement

However, pure inferring-based imputation method often fails to fill in some missing values as not all missing values are inferable based on the inference rules corresponding to those constraints. We say a missing value is inferable if there is at least one way to infer its value from the other existing or inferable values.

Advantages:

It will unlock some uncertain inferable deadlocks, which are harder to get unlocked.

H. Modules

- Filter Generation
- Threshold Queries
- Public DB Management
- Semantic Security

• Filter Generation

User interface design to create window for this application. User has to login to the system by using User ID and password. You can always see the result of the mix of all these components in the one file. Continuous aggregation queries over dynamic data are used for real time decision making and timely business intelligence. In this paper we consider queries where a client wants to be notified over distributed data crosses a specified file. Get an overview of the process of creating portlets, learn about the concepts of the APIs used to develop portlets, and view the samples to get you started. Also, learn about integrating features such as single sign-on, cooperative sharing of information using the property broker.

- **Threshold Queries**

The performance comparison of our threshold query protocols can be summarized in Complexity client and Complexity server, where enc. And dec. stand for encryption and decryption of bit, add. Andmulti. denote the holomorphic addition and multiplication of bits, and ADD. Represents the holomorphic. The performance of our generic construction depends on the performance of the underlying basic constructions.

- **Public Database Management**

Databases systems are central to most organizations' information systems strategies. At any organizational level, users can expect to have frequent contact with database systems.

Therefore, skill in using such systems understanding their capabilities and limitations, knowing how to access data directly or through technical specialists, knowing how to effectively use the information such systems can provide, and skills in designing new systems and related applications is a distinct advantage and necessity today. In public DB management always our data stored in encrypt format because attackers don't understand this data.

- **Semantic Security**

Semantic security provides measures for preventing, detaining or minimizing effects of semantic attacks. Traditional approaches to information system security focused on protecting systems and the information stored, processed and distributed on them. The goal of this project is to develop techniques to detect inconsistencies or irregularities (Behavior that breaches the rule, custom or morality) in online information, identify false information and evaluate the reliability of information sources and track those sources. A semantic attack is one in which the attacker modifies electronic

information in such a way that the result is incorrect, but looks correct to the casual or perhaps even the attentive viewer.

III. CONCLUSION

We propose a hybrid retrieving-inferring data imputation approach TRIP to alternately perform retrieving and inferring in imputing missing values in a database. Our proposed solution in TRIP is able to identify an optimal retrieving inferring scheduling scheme in DDI, and an expected optimal scheme in SDI. Extensive experimental results based on several data collections demonstrate that TRIP retrieves on average 20% missing values and the same high recall that was reached by the retrieving based approach.

IV. REFERENCES

- [1] S. Abiteboul, R.Hull, and V.Vianu Foundations of databases. Addison-Wesley Reading, 1995.
- [2] E.Agichtein and L.Gravano.Snowball: Extracting relations from large plaintextcollections. In ACM DL, pages 85–94, 2000.
- [3] J. Barnard and D. Rubin. Small-sample degrees of freedom with multiple imputation. *Biometrika*, 86(4):948–955, 1999.
- [4] G. Batista and M. Monard. An analysis of four missing data treatment methods for supervised learning. *Applied Artificial Intelligence*, 17(5-6):519–533, 2003.
- [5] P. Bohannon, W. Fan, M. Flaster, and R. Rastogi.A cost-based model and effective heuristic for repairing constraints by value modification. In SIGMOD, pages 143–154, 2005.
- [6] P. Bohannon, W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis.Conditional functional dependencies for data cleaning. In ICDE, pages 746–755, 2007.
- [7] S. Brin. Extracting patterns and relations from the world wide web. *The World Wide Web and Databases*, pages 172–183, 1999.
- [8] M. Cafarella, A. Halevy, and N. Khoussainova.Data integration for the relational web. *PVLDB*, 2(1):1090–1101, 2009.
- [9] M. Cafarella, A. Halevy, D. Wang, E. Wu, and Y. Zhang.Webtables: exploring the power of tables on the web. *PVLDB*, 1(1):538–549, 2008.
- [10] H. Elmeleegy, J. Madhavan, and A. Halevy.Harvesting relational tables from lists on the web. *PVLDB*, 2(1):1078–1089, 2009.