

HiDER: Query-Driven Entity Resolution for Historical Data

Bijan Ranjbar-Sahraei¹, Julia Efremova², Hossein Rahmani¹, Toon Calders³,
Karl Tuyls⁴, and Gerhard Weiss¹

¹ Maastricht University, Maastricht, The Netherlands

² Eindhoven University of Technology, Eindhoven, The Netherlands

³ Université Libre de Bruxelles. Brussels, Belgium

⁴ University of Liverpool, Liverpool, UK

Abstract. Entity Resolution (ER) is the task of finding references that refer to the same entity across different data sources. Cleaning a data warehouse and applying ER on it is a computationally demanding task, particularly for large data sets that change dynamically. Therefore, a query-driven approach which analyses a small subset of the entire data set and integrates the results in real-time is significantly beneficial. Here, we present an interactive tool, called HiDER, which allows for query-driven ER in large collections of uncertain dynamic historical data. The input data includes civil registers such as birth, marriage and death certificates in the form of structured data, and notarial acts such as estate tax and property transfers in the form of free text. The outputs are family networks and event timelines visualized in an integrated way. The HiDER is being used and tested at BHIC center¹; despite the uncertainties of the BHIC input data, the extracted entities have high certainty and are enriched by extra information.

1 Introduction

In the domain of historical research vast amount of historical data exists. Digitization and correction of data is an everyday process in historical centers. Additionally, some projects such as Ancestry.com² are using crowdsourcing and volunteering efforts to improve the quality of their database on census records and civil registers. This results in many dynamically changing large data corpora, requiring efficient ER.

This work develops, based on the work of [1], a query-driven tool for Historical Data Entity Resolution called HiDER. HiDER has the following advantages: (a) HiDER allows for ER across different data sources; (b) the changes in input data and ER algorithms can be incorporated in generating outcomes in real time; (c) by using *Lucene*'s inverted indexing, both structured and unstructured data are handled, and fuzzy search allows for compensating missing data and spelling variations, and (d) graph-based ER allows for detecting and visualizing “family networks”.

¹ Brabant Historical Information Center, <https://www.bhic.nl>

² Ancestry.com Inc., <http://www.ancestry.com>

2 The HiDER System

The HiDER system is developed on an Apache web server, equipped with Solr search platform. HiDER works as follows: a user gives a query which consists of at least a family name, but can also contain names of a couple, date and location and relatives' names. Subsequently, HiDER searches for relevant records existing in different sources and presents them in an integrated way. To do this, HiDER uses an inverted index data structure to retrieve a subset of records from multiple corpora, and applies an ER process, developed previously by the authors in [2, 3], on this subset on the fly. As such, the system is flexible in the sense that it adapts with minimal effort to changes in the corpus. In Fig. 1, different modules of HiDER are shown. Next we introduce each of these modules, in detail.

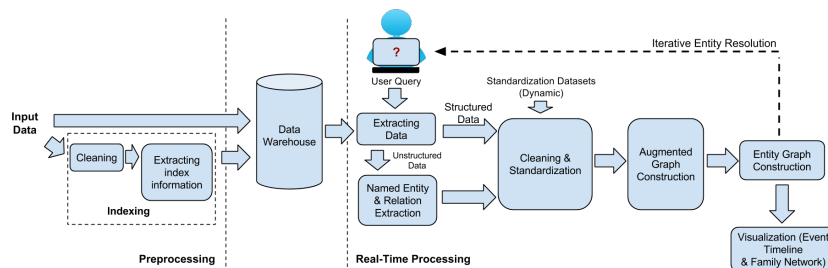


Fig. 1. The HiDER query-driven ER process.

Preprocessing: The **input data** consists of historical documents of the 18th and 19th centuries in the form of structured civil registers and unstructured notarial acts. We refer to each civil register or notarial act as a *record* and each person mentioned in a record as a *reference*. Upon arrival of a new record or when an existing record is updated, the important information of the record is **cleaned** and stored in an **inverted index**. For structured data, the names, locations, date and type of the record are the indexed information, and a *general text* field is used to generate an inverted index for every term which appears in the record. For the unstructured data the *general text* field is used to generate the inverted index for every term in the text of document. The indexing procedure is computationally light and still captures every information in the record.

Real-Time Processing: Real-Time processing is the main part of the HiDER system. Depending on the user query, HiDER uses the available indexes in the data warehouse for **Extracting Data**. The available *faceting feature* guides a user to drill into his/her target data (see left column in Fig. 2). Furthermore, the user can choose between strict and fuzzy search, where the latter one allows for compensation of spelling errors and missing data. The retrieved unstructured

data is then further processed for **Named Entity** and **Relation Extraction**; for more information we refer to [2]. Additional **Cleaning** and **Standardization** is applied to the outputs of previous modules. For instance, extra symbols are removed from names, and names with spelling variations are standardized. The standardization databases³ are continuously updated based on the user feedback and experts knowledge and updates are incorporated in answering future queries.

In the **Augmented Graph Construction** phase, the contextual information available in each record is translated to a graph component. The graph consisting of these components is then augmented by adding so called *block nodes* which capture the important features of each name such as its first and last few letters and its length (see [4] for examples). Once the augmented graph is constructed, a *random walk*-based entity detection approach is used to detect all references with highly similar contextual information (i.e., similar neighbor nodes in the augmented graph), indicating that they all refer to the same entity. Once the entities are detected, the **Entity Graph Construction** is accomplished by merging each set of references that correspond to the same entity (this technique is elaborated by the authors in [3]).

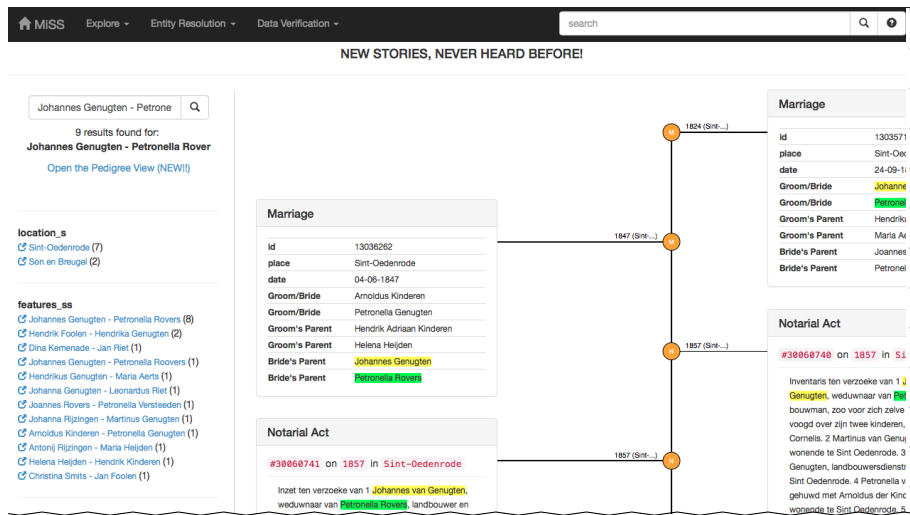


Fig. 2. Part of the HiDER interface upon arrival of a query: searching tool and faceting are shown on the left, and the event timeline is shown on the right.

Visualization serves as an indispensable tool to evaluate the entity graph manually, and is also a way to deliver the results to the user. HiDER is capable of visualizing the entity graph in the form of *event timelines*. In event timelines the information of each record is shown in the form of a floating card, while the important entities are highlighted, and the cards are sorted based on the date

³ e.g., <http://www.meertens.knaw.nl/cms/en/collections/databases>

of the records (see Fig. 2). To visualize the *family networks*, due to complexity of the generated entity graphs, we propose a novel visualization scheme for the genealogical data by combining every two individuals with marriage relations into single couple nodes, and use graph traversing algorithms to categorize nodes into different generations (see Fig. 3).



Fig. 3. The HiDER visualization of a family network: each link connects parents, on the left, to a child and his/her spouse, on the right. Users can interactively focus on the nodes and expand them. Coloring of nodes adapts to the mouse position.

HiDER allows for **Iterative ER**; the entity graph constructed in one round is used to extend the current query and as such to iteratively construct new entity graphs. Therefore, the user can retrieve the family network of farther relatives of specific entities, and also to manually compensate some of the missing links.

3 Concluding Remarks

The HiDER interactive tool targets different experts including data scientists, genealogists and demographers. Any individual who is interested in generating his/her family tree is among the main audience of HiDER, too. According to evaluations by the experts of BHIC center, using HiDER for searching the available 3,000,000 documents generates precise results (e.g., the precision of ER in [3] is 92%).

References

1. Hotham Altwaijry, Dmitri V Kalashnikov, and Sharad Mehrotra. Query-driven approach to entity resolution. *Proceedings of the VLDB Endowment*, 6(14):1846–1857, 2013.
2. Julia Efremova, Bijan Ranjbar-Sahraei, Hossein Rahmani, Frans A Oliehoek, Toon Calders, and Karl Tuyls. Multi-source entity resolution for genealogical data. In *Population Reconstruction*. Springer, 2015 (in press).
3. Hossein Rahmani, Bijan Ranjbar-Sahraei, Gerhard Weiss, and Karl Tuyls. Entity resolution in disjoint graphs: an application on genealogical data. *Intelligent Data Analysis*, 20(2), 2016 (in press).
4. Hossein Rahmani, Bijan Ranjbar-Sahraei, Gerhard Weiss, and Karl Tuyls. Contextual entity resolution approach for genealogical data. In *Workshop on Knowledge Discovery, Data Mining and Machine Learning*, Aachen, Germany, 2014.