



ENTITY IDENTIFICATION PROBLEM ESTIMATION HEART DATA INTEGRATION

Lalitha Kumari G^a, Surekha Y^{b,*}

Article DOI:

[10.55434/CBI.2024.10101](https://doi.org/10.55434/CBI.2024.10101)

Author's Affiliations

^aDept. of CSE,
PVP Siddhartha Institute of
Technology, Vijayawada, AP-India

^bDept. of CSE,
PVP Siddhartha Institute of
Technology,
Vijayawada, AP-India

Corresponding Authors

glalitha@pvpsiddhartha.ac.in,
ORCID : /0000-0002-5583-8471

ysurekha@pvpsiddhartha.ac.in,
ORCID: 0000-0002-9470-1638

Received- 01-12-2023

Accepted- 31-03-2024

©2024 The Authors Published under
Caribbean Journal of Science and
Technology

ISSN 0799-3757

<https://caribjscitech.com/>

Abstract

The incorporation of healthcare data from various sources is necessary for knowledge finding from various health data repositories. An essential research topic is preserving record linkage during the integration of medical data. In developed nations where patients' electronic health records are kept with identifiers like their social security number (SSN), universal patient identifier (UPI), health insurance number, etc., researchers have offered a variety of solutions to this issue. Due to missing information, uncertainty in patient identification, and a high level of noise in patient information, these methods cannot be properly used for record linking of health data from developing nations. People in developing nations lack medical ID cards with personalized health information. Health care facilities do not keep Social Security information or National ID numbers on file(SSN). The irony is that despite how many times a patient receives treatment or a diagnosis at the same hospital, his or her records will constantly be documented as being associated with a different patient with a different ID or serial number. In PATIENT IDENTIFICATION TECHNIQUE BASED ON SECURED RECORD LINKAGE(PITSRL) approach, NAMEVALUE algorithm is used to identify the name of the patient among diversity repositories.

Keywords:

Data Security; Health Data Warehouse; Privacy Preserved Record Linkage; Data Mining.

I. Introduction

The experiments with our testing facility are presented in this article. For both small and large training set sizes, the clustered data sets with decision trees process data significantly faster than the un-clustered data sets with decision trees. On the other hand, the accuracy obtained through clustering is always lower than the accuracy obtained without clustering, but the gap between the two diminishes as the quantity of unique entities rises. Using the clustered method over the non-clustered approach includes a trade-off between speed and accuracy. With only a slight reduction in correctness compared to un-clustered accuracy, clustering offers a significant reduction in processing time comparison to processing for un-clustered data sets.

We suggest a technique designed to enhance the effectiveness and efficiency of a current decision-tree-based method for conducting entity recognition. Prior to conducting entity recognition, our method performs data pre-processing. The data is clustered during pre-processing, and EI is then run on each cluster. To conduct EI on the clusters, we combine the decision tree method with an extra classification technique called k-NN. We set up a testbed and create appropriate experiments to investigate how well the classification methods work in different contexts. The trials change various factors such as training set size and number of unique entities. We create metrics to assess how well classification methods work both with and without pre-processing.

Except when applied to data sets with a limited number of unique entities for the smaller training set size, the decision tree method outperforms the k-NN technique in every situation. In this project the author made an attempt to resolve entity identification problem for estimation heart data integration.

2. Literature Survey

Rosario *et al*¹ proposed cross lingual named entity recognition for clinical de-identification applied to a COVID-19 Italian data set. In the proposed system it is observed that The results show that clinical de-identification is preferable in this case given the limited resources and language issue, but they also show that there is still space for improvement. The quantity of the clinical de-identification data sets accessible continues to be a limitation of this study area; therefore, it would be necessary to expand the availability of de-identification data sets.

McFarland *et al*² proposed an Improved estimation of cancer dependencies from large-scale RNAi screens using model-based normalization and data integration. In the proposed system it is observed that Absolute dependency scores from D2 are now available, allowing for the development of more specialised and reliable methods for locating genes with particular dependency patterns.

Previous models that addressed RNAi off-target effects^{1,8} had the drawback of only providing estimates of the relative variations in gene dependency across cell lines, which prevented the discovery of shared critical genes and direct comparisons of dependency scores across genes.

In the proposed system Zhang *et al*³ observed that the maintainability, ease of extension, and capacity to track the entire data history are benefits of the DataPile structure. According to a strongly desired requirement, all present applications used at all branches stay intact and operational without any change.

1. Efficiency is poor, particularly during export and matching
2. It is challenging to create direct queries because of the central repository's structure.

Srah *et al*⁴ Proposed an accuracy of claim data in the identification and classification of adults with congenital heart diseases in electronic medical records. In the proposed system it is observed that A helpful technique for identifying ACHD is administrative data using ICD- 10 codes, which can also be utilized to create a nationwide cohort.

The capacity to correctly describe the CHD populations in terms of CHD subtypes may be hampered by the lack of accuracy in the description of the CHD spectrum.

Table 1: Existing system analysis

S.No	Title	Authors	Methodology	Merits	Accuracy	De-Merits
1	Cross-lingual named entity recognition for clinical identification applied to a COVID-19 Italian dataset	Rosario Catelli, Francesco Gargiulo, Valentina Casola	IT strategy, ENIT strategy	The results obtained leave further room for improvement, although they have allowed to highlight thow, in this situation, it is desirable to proceed with	a higher accuracy, especially finer-grained and therefore at subcategory level, of the EN-IT system than the IT	limitation of this research sa remains the size of the data sets available for clinical de-identification: it would be appropriate to

				clinical de-identification given the low resources language problem.	system.	increase the availability of de-identification datasets
2	Improved estimation of cancer dependencies from large-scale RNAi screens using model-based normalization and data integration	James M. McFarland, Zandra V. Ho, Guillaume Kugener	SAP Integration	With the availability of absolute dependency scores from D2, more refined and stable approaches can be used to identify genes showing dependency pattern too interests.	D2 greatly improves identification of common-essential genes compared with existing approaches	Limitation of previous models designed to address RNAi off-target effects ^{1,8} is that they only provide estimates of the relative differences in gene dependency across cell lines, precluding identification of common essential genes and direct comparisons of dependency scores across genes
3	Data Integration Using Data Pile Structure.	David bednarek, Jakub Yaghob, Filip Zavoral	global-as-view and local-as-view	advantage of the Data Pile structure is its maintainability easy extensibility and ability to keep the track of the whole data history. All current applications used at all branches remain preserved and functional without any change according to a strongly desired requirement	All current applications used at all branches remain preserved and functional without any change according to a strongly desired requirement .	1. Efficiency, especially during export and matching, is low 2. the structure of the central repository makes constructing direct queries difficult.
4	Accuracy of claim data in the identification and classification of adults with congenital heart diseases in electronic medical records	Sarah Cohena, Anne-Sophie Jannot, Laurence Iserin c	regex-based filtering method	Administrative data using ICD-10 codes is a useful tool for detecting ACHD, and may be used to establish a national cohort.	when limited to those with moderate or complex lesions, accuracy reached 77%	the lack of accuracy in describing the spectrum of CHD may affect the ability to precisely describe the CHD populations in terms of CHD subtypes.
	How far have we come with	The odora Katsila & Minos-	QSAR	Data integration has	No matter the context,	there are vast discrepancies

5	contextual data integration in drug discovery	Timotheos Matsoukas	methodology	been considered as the panacea and road-map towards data interpretation, opinion-mining, and decision-making in data-intensive and cognitively complex settings	contextual data yields highly accurate predictions, as it is based on several sources of information – adding more contextual data to a prediction, the more precise the latter becomes	or even limitations when population-specific thinking is applied, calling for cost-effectiveness and sustainability in diverse setting
---	---	---------------------	-------------	---	---	--

3. Methodology

Semantic integration, which determines which characteristics are similar between databases, is a component of schema integration. Use neural networks to automatically conduct semantic synthesis⁵. They autonomously pull metadata from databases, including property titles and descriptions, schema details, and data contents. Then, they employ neural network methods to find comparable characteristics and discover metadata similarities. To determine whether any two provided attributes are synonyms and, if they are, to determine the connection between them, we use correlation and regression analysis methods. The first stage of database integration, model integration, is automated in all of the aforementioned methods using data mining techniques. This project assumes that schema integration has already been done and focuses on the application of data mining techniques to automate the data integration process, the second phase of database integration⁶.

The role-set method is founded on the observation that many contradictory data values for the real world entity are not contradictions but rather values that correlate to the same real world entity acting in different roles⁷. The response to a user's question is displayed as a collection of relations known as the role-set, which represents the unique interactions between the relations pertaining to various roles. As a result of the users' ability to define how role-sets should be generated, the data merging is dynamic.

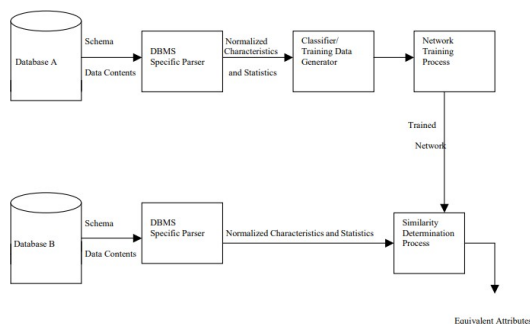


Fig. 1: Methodology of Entity Identification

This essay makes use of The EI procedure may be impacted as the number of occurrences of each unique entity declines as the number of unique entities rises. The remaining instances of the data collection are spread randomly among the unique entities, each of which has at least one occurrence. To obtain precise values for A1C, blood glucose, blood pressure, hyper lipidemia, LDL and cholesterol, obesity BMI and waist size, and date, regular expressions were created. For each risk event annotation that was followed by a measurement, a collection of training data samples was created, and regular expressions were written to extract each measurement. After a potential event had been identified by the NER-based event recognition algorithm previously discussed, regular expressions were used. A list of the regular expressions used is provided⁸.

Semantic integration, which determines which characteristics are similar between databases, is a component of schema integration. Use neural networks to automatically conduct semantic synthesis. They autonomously pull metadata from

databases, including property titles and descriptions, schema details, and data contents. Then, they employ neural network methods to find comparable characteristics and discover metadata similarities. To determine whether any two provided attributes are synonyms and, if they are, to determine the connection between them, we use correlation and regression analysis methods⁹. The first stage of database integration, model integration, is automated in all of the aforementioned methods using data mining techniques.

Data structure program architecture, operational specifics (algorithms, etc.), and interaction between modules are the main areas of emphasis in the multi-step design process. Before any coding is done, the design process turns the specifications into a display of software that can be tested for quality¹⁰. The architecture of computer software is constantly changing as new techniques, improved analysis, and greater knowledge are developed. The change in software design is still in its early stages. As a result, the depth, flexibility, and quantitative character typically linked with more traditional engineering fields are absent from the Software Design approach. However, there are methods for creating software designs, as well as standards for design attributes and design language that can be used.

4. Results and discussion

Our approach provides a representation of how closely-related records in particular homogenous sets are to those professionals in charge of merging comparable data.

One of the key objectives of our solution is detection: preventing duplicate entry by warning the user that some nearby entries already exist. In addition, they assist in decision-making by providing proximity values between these comparable records, which helps. This real-time usage would be connected to identity generation or multi-criteria identity searches. However, this prospect can only be realized with a straightforward front-end algorithm and a quick reaction time. Response times are directly correlated with algorithm optimization, particularly in the blocking phase. The change makes it possible to cut down on potential duplication.

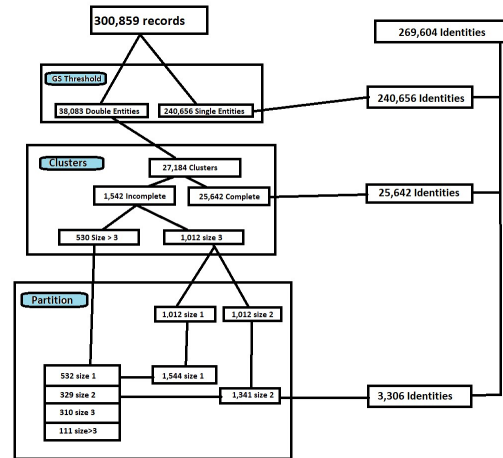


Fig. 2: Summary of Entire Linkage Procedure

Additionally, we only take into account couples that were truly positive matches rather than those that were found through traditional methods. However, there are alternative approaches that may also be taken into account. For instance, the CART approach (classification tree-based models) might be used to define clusters of data that need to be seen as "n-plicates" rather than duplicates when mapping duplicates. Another option is to propose a logistic model with the edit-distance as the independent variable and a dummy variable indicating whether or not a duplicate exists as the dependent variable. You might also utilise some additional independent variables.

5. Conclusions

In this study, we address the instance level issue of entity identification while merging existing autonomous databases. Most current approaches presume that the original relations have at least one common candidate key and that key equivalence is a legal identification condition in order to combine instances from several autonomous databases. For example "A combination of radiation therapy and chemotherapy was the first successful treatment for the patient, a 62-year-old male with squamous cell lung cancer."

As this literature demonstrates, a non-medical individual cannot comprehend the many medical terminologies. To clarify the issue and provide context for the terms "cancer" and "chemotherapy," we have picked a straightforward sentence from a set of medical data.

The Continuous Random Field (CRF) model is trained using the train dataset, and the model is

then tested using the test dataset. Further, the authors plan to (i) assess a bigger sample and suggest modifications for improved performance; (ii) incorporate more NE kinds; and (iii) obtain the needed obtaining the necessary ethics bodies' clearance before sharing some of the data for additional study in future.

6. References

- [1] Rosario, C.; Francesco, G.; Valentina, C.; Giuseppe, De P.; Hamido, F.; Massimo, E. *Applied Soft Computing*, **2020**, 97(Part A), , 106779. <https://doi.org/10.1016/j.asoc.2020.106779>
- [2] McFarland, J.M.; Zandra V. Ho.; Guillaume, K.; Joshua, M.; Dempster, P. G.; Montgomery, J. G.; Bryan, J.M.; Krill-Burger, T. M.; Green, Fr. V.; Jesse, S.; Boehm, T. R.; Golub, W. C.; Hahn, D. E.; Root & Aviad T.. *Nature Commun.*, **2018**, 9, 4610. <https://doi.org/10.1038/s41467-018-06916-5>
- [3] Zhang, H.; Guo, Y.; Li, Q.; Thomas, J.; George, El. S.; François Modave & Jiang B. *BMC Med Inform Decis Mak.*, **2018**, 18 (Suppl 2), 41. <https://doi.org/10.1186/s12911-018-0636-4>
- [4] Sarah, C.; Anne-Sophie, J.; Laurence, I.; Damien, B.; Anita, B.; Jean-Baptiste, E. *Archives of cardiovascular diseases*, **2019**, 112, 31. DOI: 10.16/j.acvd.2018.07.002.hal-03486373
- [5] Katsila, T.; Matsoukas, M. T. *Expert Opinion on Drug Discovery*, **2018**, 13(9), 791. <https://doi.org/10.1080/17460441.2018.1504767>
- [6] Dhayne, H.; Haque, R.; Kilany, R.; Taher, Y.; *IEEE Access*, **2019**, 7, 91265. Doi: 10.1109/ACCESS.2019.2927491.
- [7] Prasad, G.L.V.; Kollu, V.N.; Sailaja, M.; Radhakrishnan, S.; Jagan Mohan, K.; Kishore Reddy, A., Rajesh Chandra, G. *Trans. Electr. Electron. Mater.*, **2024**, 25, 89. <https://doi.org/10.1007/s42341-023-00487-z>
- [8] Kumar, C. N. S.; Sailaja, M.; Hussain, M. A.; Rahman, S. Z. *Intl. J. Rec. Innov. Trends Comp. Commun.*, **2022**, 10(2s), 146. <https://doi.org/10.17762/ijritcc.v10i2s.5921>
- [9] Prasad B.D.C.N.; Sailaja, M.; Suryanarayana, V. **2022**, *ECS Trans.*, 107(1), 19777. DOI: 10.1149/10701.19777ecst
- [10] Sailaja, M.; Ahad, A.; Sivaramakrishna, K.; Hussain, A. *J. Phys.: Conf. Ser.*, **2021**, 2089, 012082. DOI 10.1088/1742-6596/2089/1/012082.