

Efficient and Effective On-line Learning-Based Instance Matching over Heterogeneous Data

Samur Araujo ^{#1}, Duc Thanh Tran ^{*2}, Arjen de Vries ^{#3}

[#]*Delft University of Technology, PO Box 5031, 2600 GA Delft, the Netherlands*

¹*s.f.cardosodearaujo@tudelft.nl*

³*a.p.devries@tudelft.nl*

^{*}*Karlsruher Institute of Technology, Germany*

²*ducthanh.tran@kit.edu*

Abstract—This document gives formatting instructions for authors preparing papers for publication in the Proceedings of an IEEE conference. The authors must follow the instructions given in the document for the papers to be published. You can use this document as both an instruction set and as a template into which you can type your own text.

I. OVERVIEW

In this section, we introduce the data model, discuss the problem, and finally, present a brief overview of existing solutions as well as our approach.

A. Data

We focus on heterogeneous Web data including relational data, XML, RDF and other types of data that can be modeled as graphs. Closely resembling the RDF data model, we employ a graph-structured data model where every dataset is conceived as a graph $G \in \mathbb{G}$ comprising a set of triples:

Definition 1 (Data Model): A dataset set is a graph G formed by a set of triples of the form (s, p, o) where $s \in U$ (called subject), $p \in U$ (predicate) and $o \in U \cup L$ (object). Here, U denotes the set of Uniform Resource Identifiers (URIs) and L the set of literals. Every literal $l \in L$ is a bag of tokens, $l = \{t_1, \dots, t_i, \dots, t_n\}$, drawn from the vocabulary V , i.e. $t_i \in V$.

With respect to this model, *instances* are resources that appear at the subject position of triples. An instance representation can be obtained from the data graph as follows:

Definition 2 (Instance Representation): The instance representation $IR : U \times \mathbb{G} \rightarrow \mathbb{G}$ is a function, which given an instance $s \in U$ and a graph $G \in \mathbb{G}$, maps s to a set of triples in which s appears as the subject, i.e. $IR(s) = \{(s, p, o) | (s, p, o) \in G\}$.

Thus, an instance is basically represented through a set of *predicate-value* pairs. **@TODO: add a graph to provide example, also, provide one example for instance representation** We will use the terms instance and instance representation interchangeably from now on. Note that for the sake of presentation, only the outgoing edges (s, p, o) of an instance s are considered while incoming edges (triples where s appears as the object) can also be added to the representation of s .

B. Problem

Instance matching is about finding instances that refer to the same real-world object based on their representations extracted from the data. In this paper, we tackle the problem of *instance matching across multiple datasets*: given instances of a *source dataset* G_s , the problem is to find instances in other datasets, collectively referred to as the *target dataset* G_t , which represents the same real world object. Compared to the single dataset setting, this problem entails additional challenges especially when the data is heterogeneous not only at the data but also schema level. In this regard, *data-level heterogeneity* means that for the same property, instances referring to the same object may have different values, or different syntactical representations of the same value. For instance, the values "Michael Jackson" or "Jackson, Michael" can be both used as the value of the property name of an instance representing Michael Jackson. Schema-level heterogeneity arises when there is only little or no schema overlaps between the datasets. That is, instances referring to the same object may be represented by different predicates, or different representations of the same predicates. Finding different representations of the same predicate is part of a problem also known as schema matching.

C. Existing Solutions

One main class of solutions for this problem comprises learning-based approaches, which relying on training examples, find a (weighted) combination of similarity function predicates, $\sum_i w_i s_i(p_m, p_n) > \alpha$, to decide if a given pair of instances match (when the similarity exceeds the threshold α), or not. This learning entails the problem of finding the comparable predicates p_m and p_n (schema matching), choosing (and weighting) these predicate pairs, as well as determining the threshold α . The second category of solutions comprises unsupervised learning-based approaches, which without training example, try to predict with two instances match. Usually those solutions are clustering-based approaches or probabilistic-based models, that try to predict the comparable predicates and the threshold α , by only inspecting the distribution of the terms in data and schema. For instance, [1] can unsupervisely learn the source and target comparable predicates used during the matching process; however, the

mapping between those predicates is manually defined. As well as, the similarity threshold α . [2] describe an schema agnostic unsupervised approach. In this approach, the learning of comparable predicate is unnecessary. All instances are presented as a bag of tokens and they match when they share similar tokens. **@TODO: class-based disambiguation: make clear that learning/finding comparable predicates is not possible when there are no overlaps between schemas, i.e. predicates in different datasets do not matches.**

D. Our Solution

In this work, we perform learning but... **@TODO: differences to existing learning-based approaches: not a (weighted) combination of similarity function predicates to be applied to all pairs of instances but one by one: for every source instance, we learn optimal match selection queries:**

Definition 3 (Template Query): A template query Q is conjunction of n tuples (p, k) , where p is a predicate in G_t and k is a token, defined as $\bigwedge_i^n (p_i, k_i) = \{t | (t, p_i, o) \in G_t \wedge o \sim k_i\}$. The similarity function \sim is given. Without loss of generality we can also assume that can exist tuples where p are undefined, acting as a wild card.

Definition 4 (Candidate Set): A candidate set of an instance s in G_s is a instantiation of a template query Q where we have a tuple (p, k) where $(s, x, k) \in IR(s)$

In this work, we discuss how query templates can be derived. Also, executing all these queries is expensive, hence we propose to reduce the number of queries that has to be executed for retrieving candidates for every instance in the source.

For a given set of source instances S , we can model a candidate selection query as a template that can be instantiate for each source instance $s \in S$.

Candidate selection queries refer to the problem of finding queries, for each instance s in a source set G_s , that retrieve a set of possible candidate matches for s in the target dataset G_t . Ideally, assuming that there is a 1-to-1 mapping between the instances, the best query produces candidates set with cardinality 1, because we can not have less than one candidate. We pose the problem as an optimization problem, where the goal is to find the most selective conjunction of query clauses that produces minimal candidate sets.

Solutions. Aiming to approximate our solution to the optimal solution, we approach this problem in an iterative fashion, where at each iteration we use information from the previous iteration to refine and increase the selectivity of the query templates used on next iterations. Notice that as more selective a query is, as more efficient and effective it is, because, it selects less elements and it takes less time; and it selects less incorrect matches.

Without any knowledge of the source or target schemas, we start the process with queries templates that are less selective but easy to build. As the process moves on, information obtained from the candidate sets produced in previous iterations are used to refine the query templates for next iterations. Basically, this refining process adds two types of clauses in the

query templates, aiming to make them more selective: attribute and class clauses. Attributes clauses are composed of highly selective target predicates with values similar to the values of at least one highly selective source predicate. The value of this source predicate is used as the object value of the attribute clause. To build attribute clauses we apply the Algorithm 2 described in our previous work [?] over the source instances and instances in the candidate sets already generated. Class clauses are predicate/value pairs that represent the class of interest of the target instances (e.g. `rdfs:type=geo:country`). To build class clauses, we apply a matcher over the already generated candidate sets obtaining positive and negative examples. Then, those positive and negatives examples are input to an algorithm that output the set of predicate/value pairs that select all positives and avoid the negatives examples; namely, the class of interest. Finally, those pairs are used to compose the class clauses. The matcher can be any approach that uses a more complex similarity measure to select the correct matches (positive examples) among the possible candidates. In this work, we assume that the source instances belong to a specific class of interest (e.g. countries), therefore we use the class-based disambiguation paradigm as the matcher. We do so, because this is the only complex matcher that can lead with datasets with non-overlapping schemas.

The set of query templates generated in this process can be large and their results for a specific instance may overlap. Hence, a set of heuristic is applied to decide whether or not to evaluate a template, aiming to avoid evaluating templates that will produce overlapping candidate sets. Those heuristics are embedded in a branch-and-bound optimization framework that on-line learns to efficiently evaluate the most effective query templates. As result, this framework produces the minimal candidate sets passing only once over each source instance.

II. CONCLUSIONS

REFERENCES

- [1] D. Song and J. Heflin, "Automatically generating data linkages using a domain-independent candidate selection approach," in *International Semantic Web Conference (I)*, 2011, pp. 649–664.
- [2] G. Papadakis, E. Ioannou, C. Niederée, and P. Fankhauser, "Efficient entity resolution for large heterogeneous information spaces," in *WSDM*, 2011, pp. 535–544.