

DAGOBDAH: An End-to-End Context-Free Tabular Data Semantic Annotation System

Yoan Chabot¹[0000-0001-5639-1504], Thomas Labbé¹[0000-0001-9295-7675],
Jixiong Liu¹, and Raphaël Troncy²[0000-0003-0457-1436]

¹ Orange Labs, France

yoan.chabot|thomas.labbe@orange.com

² EURECOM, Sophia Antipolis, France

raphael.troncy@eurecom.fr

Abstract. In this paper, we present the DAGOBDAH system which tackles the Tabular Data to Knowledge Graph Matching (TDKGM [6]) challenge. DAGOBDAH aims to semantically annotate tables with Wikidata and DBpedia entities, and more precisely performs cell and column annotation and relationship identification, via a pipeline starting from pre-processing to enriching an existing knowledge graph using the table information. This paper presents techniques for typing columns with fine-grained concepts while ensuring good coverage, and for valuing these types when disambiguating the cell content. This system obtains promising results in the CEA and CTA tasks on the challenge test datasets.

Keywords: Tabular Data · Knowledge Graph · Entity Linking · Embedding · DAGOBDAH · TDKGM Challenge

1 Introduction

The annotation of tables using knowledge graphs is an important problem as large parts of both companies internal repositories and Web data are represented in tabular formats. This type of data is difficult to interpret by machines because of the limited context available to resolve semantic ambiguities and the layout of tables that can be difficult to handle. The ability to annotate tables automatically using knowledge graphs (encyclopedia graphs such as DBpedia and Wikidata or enterprise-oriented graphs) allows to support new semantic-based services. For example, it opens the way to more efficient solutions to query (e.g. “moving beyond keyword” for dataset search [2]), manipulate and process heterogeneous table corpus [1].

In this paper, we propose a complete system to annotate tabular data, from pre-processing to entities and relations extracted from knowledge bases, without using any tables context. The main contributions of our system are:

Copyright © 2019 for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

- A new pre-processing tool chain improving the results of the DWTC framework³.
- A three-step annotation process (cell-entity annotation, column-type annotation and disambiguation) leveraging on Wikidata and DBpedia.
- An alternative approach based on clustering operations using Wikidata embedding.

2 DAGOB AH: an End-to-End System for Annotating Tables

DAGOB AH is implemented as a set of tools, used in sequence, to provide three main functionalities.

1. The identification of mapping relationships between tabular data and knowledge graphs.
2. The enrichment of knowledge graphs by transforming into triples the knowledge contained in table. DAGOB AH’s target knowledge graph is Wikidata for several reasons including its dynamics, coverage and the quality of data [4]. Thus, adaptations had to be made for the challenge to support DBpedia.
3. The production of metadata that can be used for datasets referencing, search and recommendation processes [1].

To provide these features, DAGOB AH is structured in the following way. The pre-processing modules (Section 2.1) perform the tables cleaning and a first high-level characterisation of their shape and content. The entity linking modules accomplish the tasks of the challenge itself, namely the cell-entity annotation (CEA), the column-type annotation (CTA) and the columns-property annotation (CPA). Two methods have been studied to carry out these tasks: a baseline exploiting lookup and voting mechanisms (Section 2.2) and a geometric approach based on clustering applied to Wikidata embeddings (Section 2.3).

2.1 Tables Pre-Processing

In order to correctly process the information contained in the tables, it is first necessary to infer several characteristics on the shape of the table. In a context of real exploitation in which there is sometimes little or no knowledge on the tables, the information produced by this chain is decisive for the quality of the annotations. The pre-processing toolbox of DAGOB AH, partly based on the DWTC-Extractor, generates four different types of information. The precision of the toolbox was evaluated on round 1 dataset (Table 1) and compared to a modified version of the DWTC (which does not use HTML tags and thus supports more formats).

Table Orientation Detection. The initial algorithm proposed in the DWTC-Extractor is based on the length of the strings and the assumption that the cells

³ <https://github.com/JulianEberius/dwtc-extractor>

in the same column have a similar size. However, the robustness of the DWTC algorithm can be improved. Indeed, for example, two strings representing very different elements may have the same length (for example "Paris" and "10cm²"). DAGOBDAH introduces a new algorithm based on a primitive cell typing system with eleven types (string, floating numbers, date, etc.). Based on these types t_i , an homogeneity score is computed on each row and each column x (Equation 1). The mean of all rows and all columns is then compared, and depending on the ratio, the table is said "HORIZONTAL" or "VERTICAL".

$$Hom(x) = \left[\frac{1}{len(x)} \sum_{t_i \in x} \left(1 - \left(1 - 2 * \frac{count(t_i)}{len(x)} \right)^2 \right) \right]^2 \quad (1)$$

Header Extraction. The algorithm used in DAGOBDAH is based on the primitive types defined above. The header extractor assumes that the header of a column contain mainly strings and, in most cases, does not share the type of the column cells. The use of these two heuristics allows to identify if a table contains or not a header with a good accuracy (see Table 1). It should be noted that the DWTC framework also offers a header detection tool but it contains several bugs that make impossible its evaluation.

Key Column Detection. A first step aims to identify a first low level type for each column among five given types (Object (mentions that are potentially lookup-able in a knowledge base), Number, Date, Unit (e.g. "12 km") and Unknown (containing all the unclassified columns)). The key column is an Object column containing a large number of unique values and located on the left side of the table. This pre-processing toolbox was especially useful during the first

Table 1. Precision of pre-processing tasks

Task/Tool	DWTC	DAGOBDAH
Orientation Detection	0.9	0.957
Header Extraction	Not evaluated	1.0
Key Column Detection	0.857	0.986

round to automatically identify the information contained in the header as well as the column to be annotated in each table. In addition, when the goal is to enrich a knowledge graph with information contained in tables, key column detection is a critical element in determining the subject of the generated RDF triples. The availability of targets during the second round made the use of this pre-processing chain obsolete. However, this does not affect the usefulness of such tools in real world applications.

2.2 Baseline Approach

Entities Lookup. A preliminary cleaning process is first applied in order to optimize the lookups. The intent is not to correct every string issues, but to

have a macroscopic transformation process covering the most known artefacts: encoding homogenization and special characters deletion (parenthesis, square bracket and non alphanumeric characters). Five lookup services are simultaneously queried to retrieve entity candidates from cell mention: Wikidata API, Wikidata Cirrus Search Engine, DBpedia API, Wikipedia API and an internal ElasticSearch index where Wikidata has been ingested, that contains labels, aliases and types associated to Wikidata QIDs. The benefit of having such source is also to manage the indexes, thus the search strategy. An occurrence count is performed at the output of the lookups to select the most popular candidates, and their corresponding types, among the five services. As DBpedia is the target knowledge base for the challenge, QIDs and Wikidata types are translated to equivalent DBpedia entities (using SPARQL query and following *owl:sameAs* and *owl:equivalentClass* predicates). Ancestors are retrieved for each class in the resulting list through SPARQL query to the DBpedia endpoint in order to have an extended list of candidates types.

Column Typing. Before proceeding to counting, it is necessary to remove non-relevant types. A basic type coverage of the cells criteria is not relevant as right types might be more specific but not frequent enough to be consider as the target ones. To solve this issue, a threshold based on relative scores is used. To each type t in $\{t_i\}$ (list of all types in a given column), a score S_t is first associated, from which a relative score R_t is computed (Equation 2).

$$S_t = \frac{\text{count}(t)}{\text{sum}(t_i)} \quad \text{and} \quad R_t = \frac{S_t}{\max(S_{t_i})} \quad (2)$$

Only types with $R_t > 0.7$ (configurable threshold) are considered in the next steps. In order to select the target type from the short-listed ones, a TF-IDF-like method is used to compute the specificity. It reflects the importance of a type t in $\{t_i\}$ (Equation 3).

$$S_{\text{spec}}(t) = S_t * \log\left(\frac{N_L}{\text{count}(t)}\right) \quad (3)$$

where N_L is the number of lookup candidates and $\text{count}(t)$ is the occurrences of type t within the given column. The inherent advantage of this method is to be independent from the target knowledge base.

Entities Disambiguation. In order to enhance the CEA results, the previous ordered types are used to disambiguate the candidates entities. If the first candidate of the lookup has the target type, it is selected as the target entity; if not, we select the entity associated to this type in the lookup list (if the list is empty, no annotation is produced).

2.3 Embedding Approach

The intuition behind this approach is that entities in the same column should be closed in the embedding space as they share semantic characteristics, and thus could form coherent clusters. In order to compare this approach to the baseline,

DAGOBDAH: Tabular Data Semantic Annotation System

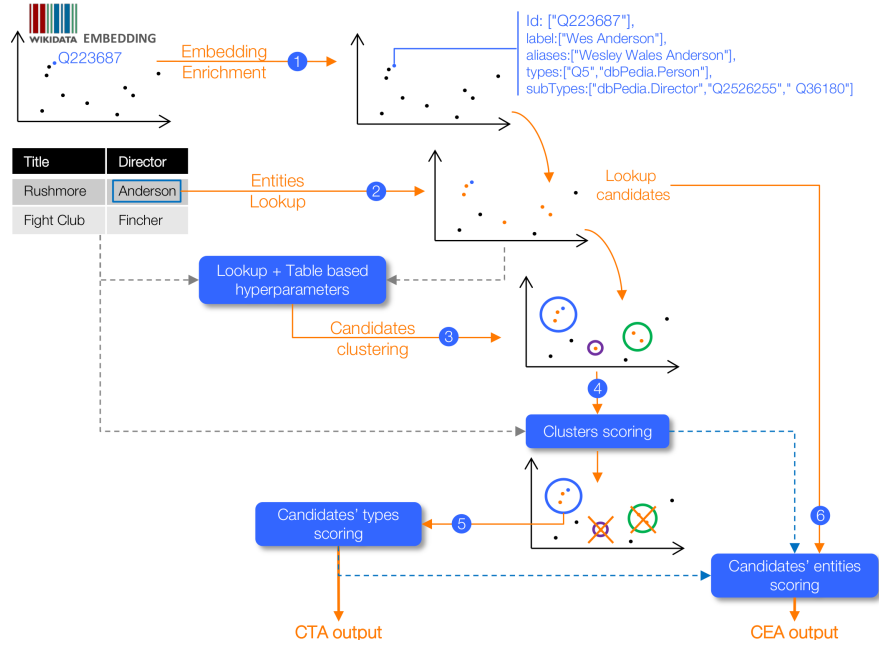


Fig. 1. Workflow of the embedding approach

a syntactic lookup to Wikidata embedding is done.

Embedding enrichment and Lookup. The pre-trained Wikidata embedding [5] only contains Wikidata QIDs. Labels, aliases as well as types are added to each entity through an internal Elasticsearch server (step 1 in Figure 1). Lookups are then used to find candidates matching the content of each cell (step 2 in Figure 1). Both regex and Levenshtein distance strategies have been implemented.

1. Entity candidate label or aliases should include all the words of the original mention (with no order).
2. Levenshtein ratio [7] between an entity string and mention ($item_{sim}(i)$) should be larger or equal than 0.75.

If an entity (label or aliases) satisfies one of the previous conditions, it is accepted as candidate.

Candidates Clustering and Scoring. The challenge is to have most of the expected candidates for a given column in one of the clusters. A grid search strategy measuring precision of correct candidates clustering with different algorithms and K values was implemented to determine the best algorithm (K-means with hyperparameter K equals to number of lookup candidates divided by number of rows)(step 3 in Figure 1). In order to select the relevant cluster, a scoring algorithm has been defined (step 4 in Figure 1). The clusters with the highest rows coverage (i.e. number of rows having at least one candidate in the cluster)

Chabot et al.

are selected. Then, a confidence score is associated to each candidate within these clusters (Equation 4).

$$S_c(i) = item_{conf}(i) * item_{sim}(i)^x \quad (4)$$

where $item_{conf}(i)$ is the co-occurrence score given by Equation 5 and $x \in \mathbb{N}^+$ allows to give more importance to the textual similarity.

$$item_{conf} = FE + 0.5 * FH \quad (5)$$

where FE and FH are defined in Equation 6.

$$FE = \frac{e + 1}{N_C}, \quad FH = \frac{h + 1}{N_C} \quad (6)$$

where e = number of entities in other columns matching with Wikidata properties values of the candidate; h = number of headers in other columns matching with Wikidata properties labels of the candidate; N_C = number of table columns. The normalized confidence score for a given cluster n is then computed using:

$$S_k(n) = \frac{\sum_{i \in n} S_c(i)}{len(n)} \quad (7)$$

where $len(n)$ is the number of elements within cluster n .

From all candidates in the selected clusters, a counting for every existing type is computed, each type inheriting the confidence score of its corresponding candidate (step 5 in Figure 1). All types with a score higher than a threshold ($Max(score) * 0.75$) is selected. Thus, the output type is the one having the highest specificity within DBpedia hierarchy (using subclasses count and distance to *owl:Thing*). Finally, the candidates of each cell resulting from the lookup operations are examined according to the selected type (step 6 in Figure 1). The following score is computed for each lookup entity i belonging to cluster n :

$$S_e(i) = R_t * (0.2 * S_k(n) + 0.5 * S_c(i)) \quad (8)$$

where $R_t = 1.5$ if the entity belongs to the type T produced in CTA, 1.2 if it belongs to a parent of T and else 1. From a given row candidates, the output entity is the one with the highest score.

3 Results

The baseline was used during the four rounds of the challenge (Table 2). The CEA's results are satisfactory, but the baseline has difficulty in producing the CTA results expected by the evaluator. In round 1, the predicted type was often too generic or too specific. In addition, the baseline showed two important limitations: a high dependency on lookup services (over which DAGOBAN has little control) and difficulties in correctly setting up algorithms (in particular finding the right compromise between specificity and representativeness of types in the

Table 2. Results of the baseline and embedding approaches: rounds 2-3-4 of the challenge evaluated by AICrowd on 30 November 2019 and round 1 evaluated by the same scorer after the organizers have made it available

	Task	CTA		CEA		CPA	
	Criteria	Prim. Score	Sec. Score	F1 Score	Precision	F1 Score	Precision
Round 1	Baseline	0.479	0.242	0.883	0.892	0.415	0.347
	Embedding	1.212	0.336	0.841	0.853	-	-
Round 2	Baseline	0.641	0.247	0.713	0.816	0.533	0.919
Round 3	Baseline	0.745	0.161	0.725	0.745	0.519	0.826
Round 4	Baseline	0.684	0.206	0.578	0.599	0.398	0.874

case of CTA). Concerning the CPA, a simple lookup technique on the header was used during round 1 explaining the low accuracy.

For the next rounds, a search over relationships between each pair of instances (output of CEA) of the two candidate columns followed by a majority vote was used. This second technique has significantly improved the accuracy of the CPA. To show the contribution of the embedding approach, an evaluation was carried out on the corpus of round 1⁴ with an enrichment of the CTA GT consisting in the addition of the parent types (except *owl:Thing*) of the perfect type in order to allow the correct evaluation of the primary and secondary scores. CEA performances are slightly poorer because of basic lookup strategies used (compared to the fully-optimised lookups used in the baseline) and the absence of expected candidates in the selected clusters. But the embedding approach proves to be highly proficient to determine the type of a column which is the core of more reliable annotations. In addition, the results are particularly interesting in cases where string mentions in the original table are incomplete. In table 54719588_0.8417197176086756912 for instance, one column references movie directors only by their family names (in that case, our baseline performance was poor). Doing K-means clustering on a subset of this table (four rows) performs pretty well even with very few data, as illustrated in Figure 2 (2 clusters shown among 12). Indeed, the green cluster gives the expected candidates, even if disambiguation still has to be done for some cells (using S_e).

4 Conclusion and Future Works

In this paper, we have presented a baseline and an embedding approach to carry out the task of generating semantic annotations of tables. The embedding approach shows very encouraging CTA results and an ability to infer candidates from incomplete information and with no tailored data cleaning. However, optimizing the hyper parameters still remains challenging. The way DAGOBAB computes the number of clusters for K-means (based on lookup and table properties) gives good results but might not be robust with all datasets (this is why

⁴ For timing and performance reasons, this approach could not be tested on other rounds in time.



Fig. 2. Result of K-means clustering applied to Wikidata embedding

we have decided to consider the high-scored clusters instead of relying entirely on the first one). Other clustering algorithms shall be tested that may be more accurate to find the best compromise between having all candidates in the same cluster and enough clusters for good discrimination between candidates.

Combining Wikipedia and Wikidata in a joint embedding space and use Fasttext-like embeddings, as well as exploiting semantic lookup in addition to syntactic one could significantly enhance the entities mapping. Finally, a full vectorial approach [3] consisting of learning table rows embedding (possibly Poincaré-like) combined with geometric constraints derived from column mentions, and then find mapping (general or local) transformation(s) with a Wikidata/Wikipedia embedding space could enhance the results.

References

1. Chabot, Y., Grohan, P., Le Calvez, G., Tarnec, C.: Dataforum: Data exchange, discovery and valorisation through semantics. In: Extraction et Gestion des Connaissances (EGC). Metz, France (2019)
2. Chapman, A., Simperl, E., Koesten, L., Konstantinidis, G., Ibáñez, L.D., Kacprzak, E., Groth, P.: Dataset search: a survey. *The VLDB Journal* pp. 1–22 (2019)
3. Chen, J., Jimenez-Ruiz, E., Horrocks, I., Sutton, C.: ColNet: Embedding the Semantics of Web Tables for Column Type Prediction. In: 33rd AAAI International Conference on Artificial Intelligence (2018)
4. Färber, M., Ell, B., Menne, C., Rettinger, A.: A Comparative Survey of DBPedia, Freebase, OpenCyc, Wikidata, and YAGO. *Semantic Web Journal* pp. 1–5 (2015)
5. Han, X., Cao, S., Lv, X., Lin, Y., Liu, Z., Sun, M., Li, J.: Openke: An open toolkit for knowledge embedding. In: International Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 139–144 (2018)
6. Hassanzadeh, O., Efthymiou, V., Chen, J., Jiménez-Ruiz, E., Srinivas, K.: SemTab2019: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching - Data Sets. Zenodo (2019), <https://doi.org/10.5281/zenodo.3518539>
7. Sarkar, S., Pakray, P., Das, D., Gelbukh, A.: JUNITMZ at SemEval-2016 Task 1: Identifying semantic similarity using levenshtein ratio. In: 10th International Workshop on Semantic Evaluation (SemEval). pp. 702–705 (2016)