



SportsTables: A New Corpus for Semantic Type Detection (Extended Version)

Sven Langenecker^{1,2} · Christoph Sturm¹ · Christian Schalles¹ · Carsten Binnig^{2,3}

Received: 7 June 2023 / Accepted: 22 September 2023 / Published online: 16 October 2023
© The Author(s) 2023

Abstract

Table corpora such as VizNet or TURL which contain annotated semantic types per column are important to build machine learning models for the task of automatic semantic type detection. However, there is a huge discrepancy between corpora and real-world data lakes since they contain a huge fraction of numerical data which are not present in existing corpora. Hence, in this paper, we introduce a new corpus that contains a much higher proportion of numerical columns than existing corpora. To reflect the distribution in real-world data lakes, our corpus SportsTables has on average approx. 86% numerical columns, posing new challenges to existing semantic type detection models which have mainly targeted non-numerical columns so far. To demonstrate this effect, we show in this extended version paper of [18] the results of an extensive study using four different state-of-the-art approaches for semantic type detection on our new corpus. Overall, the results demonstrate significant performance differences in predicting semantic types for textual and numerical data.

Keywords Semantic type detection · Column annotated corpora

1 Introduction

Semantic Type Detection Is Important for Data Lakes Semantic type detection of table columns is an important task to exploit the large and constantly changing data collections residing in data lakes. However, manually annotating tables in data lakes comes at a high cost. Hence, in the past a lot of approaches have been developed that automatically derive semantic types from table data [6, 15, 24, 26]. Many of the recent approaches use deep learning techniques to build semantic type detection models. As such, corpora containing

large amounts of table data with assigned semantic types are required for training and validating. Existing annotated table corpora (e.g. VizNet, TURL) primarily contain tables extracted from the web and therefore limit the capability to represent enterprise data lakes.

Existing Corpora and Models Fall Short on Real-World Data Lakes However, as we can see in Fig. 1, almost all existing corpora that provide annotated columns labeled with semantic types have a lack of table columns that contain numerical data, and tables in these datasets incorporate either only for a very high percentage of textual data. Only GitTables [17] contains a more balanced ratio of textual and numerical data. Nevertheless, compared to real enterprise data lakes, there is a significant discrepancy in the ratio of textual to numerical data. An inspection of a large real-world data lake at a company¹ has shown that on average approx. 20% textual data and 80% numerical data are present (see. Fig. 1 bars on the right). Moreover, semantic type detection models [6, 15, 24, 26] that are trained on the available corpora also mainly target non-numerical data.

Semantic Type Detection for Numerical Data Is Challenging Detecting semantic types of numerical columns is gener-

✉ Sven Langenecker
sven.langenecker@mosbach.dhbw.de

Christoph Sturm
christoph.sturm@mosbach.dhbw.de

Christian Schalles
christian.schalles@mosbach.dhbw.de

Carsten Binnig
carsten.binnig@cs.tu-darmstadt.de

¹ Duale Hochschule Baden Württemberg Mosbach, Mosbach, Germany

² TU Darmstadt, Darmstadt, Germany

³ DFKI Darmstadt, Darmstadt, Germany

¹ The analyses were done at the company LÄPPLE AG.

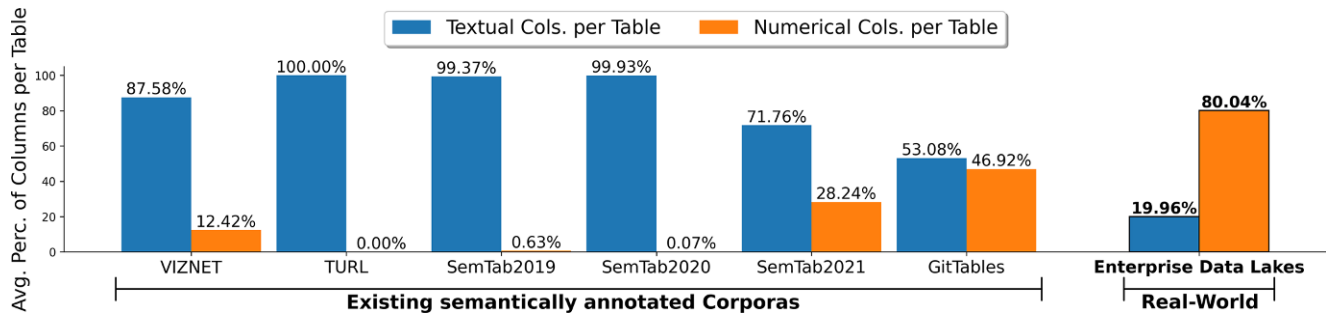


Fig. 1 Average percentage of textual and numerical based columns per table in existing semantically annotated corpora (Notice that for GitTables we only considered the tables and columns labeled by terms from DBpedia using the semantic annotation method as described in the GitTables paper. Therefore our reported ratios of textual and numerical data differ from those shown in the GitTables paper because they consider all data, whether annotated or not.) (left bars) compared to real-world data lakes (right bar). This shows the fact that there is a significant shift in the ratio of textual to numeric columns per table from existing corpora to real data lakes. Since all existing semantic type detection models were developed by using the existing corpora, shortcomings in validating the models on numerical data are present and it has not yet been studied in depth how well the models can perform on datasets containing a high proportion of numerical data.

ally harder than for textual columns. For example, for a textual column with the values {Germany, USA, Sweden, ...} a model can easily identify the semantic type *country*. Instead, for a numeric column with e.g. the values {20,22,30,34,...} it is not that straightforward and several possibilities for a matching semantic type exist such as *age*, *temperature*, *size*, *money*. The fundamental reason here is that numerical values can be encoded with much fewer bits than string values [23], resulting in a lower overall entropy and thus providing less information content that can be used by a machine learning model to infer the underlying semantic type. Due to the existing corpora providing annotated columns that have been used to create and validate semantic type detection models, we see several essential shortcomings that could not be addressed until now because of the absence of a sufficient dataset for this purpose.

Contributions In this paper, we thus contribute a new corpus containing tables with semantically annotated columns with numeric and non-numeric columns that reflect the distribution of real-world data lakes. We will make the corpus available which should stimulate research directions such as working on new model architectures that can reliably annotate types to numeric and non-numeric columns. In the following, we discuss the main contribution of this paper.

As a first contribution, we present and provide our new corpus SportsTables². To the best of our knowledge, SportsTables is the first corpus with annotated table columns, which contains a significantly larger proportion of numerical data than textual data. In total, the tables in our corpus have on average about 3 textual and 18 numerical columns. Moreover, the tables in our new corpus are much larger in both the number of columns and the number of rows than

in existing corpora which better reflects the characteristics of real-world tables.

As a second contribution that comes together with the corpus, we specify an ontology with semantic types for the sports baseball, basketball, football, hockey, and soccer. This ontology provides fine granular semantic types for all kinds of sports we considered to build SportsTables and allows us to semantically describe each occurring table column, which is not possible with the current ontologies (e.g. DBpedia) at this level of detail. Using a manually created dictionary, we assign a semantic type to each existing column in SportsTables.

As an extension of [18] and as a third contribution, we present in this paper results of extensive experimental analyses using our new corpus on four different state-of-the-art semantic type detection models. Overall, we can see that when trained on our new corpora, the models can improve the performance on numerical data types. However, one shortcoming that our analysis shows is that the current model architectures are not targeting numerical columns. To be more precise, our analysis demonstrates that textual data columns are mostly correctly semantically interpreted with the models (best F1-Score 0.98), but on numerical data columns, the models only achieve F1-Scores in the range of 0.31–0.7. This large difference indicates that new model architectures that take the characteristics of numerical columns into account are needed which is a direction that could be stimulated by the availability of our corpus.

Outline In Sect. 2, we first provide an overview of existing corpora which was used to build and validate semantic type prediction models and discuss their characteristics and statistics. Afterward, in Sect. 3, we then introduce our new corpus SportsTables and describe in detail how we created the corpus and labeled the table columns with semantic types. Section 4 first demonstrates the main characteristics

² Available on <https://github.com/DHBWMosbachWI/SportsTables.git>.

of our corpus before we then show the results of using our new corpus on the different semantic type detection models. Next, further research challenges are discussed in Sect. 5 before Sect. 6 concludes the paper.

2 Existing Corpora

In the following, we describe different existing corpora that contain annotated table columns and therefore can be used to build and validate semantic column type detection models. We summarized the main statistics for all corpora in Table 1.

VizNet [14] The original VizNet corpus [14] is a collection of data tables from diverse web sources [4, 20, 22, 25] which initially do not contain any semantic label annotation. The corpus we consider in this paper is a subset of the original VizNet corpus, which was annotated by a set of mapping rules from column headers to semantic types and then used to build and validate the Sherlock [15] and Sato [26] prediction models. The corpus contains in total 78,733 tables and 120,609 columns annotated with 78 unique semantic types. Overall, the tables in the corpus contain only 1.53 columns and 18.35 rows on average. Furthermore, the distribution of the column data types is 87.58% textual and 12.42% numerical and thus leads to the shortcomings as described before.

TURL [6] The TURL corpus uses the WikiTable corpus [3] as basis. To label each column they refer to the semantic types defined in the Freebase ontology [9] with a total number of 255 different semantic types. What distinguishes TURL from other corpora is that columns can have multiple semantic types assigned. In total, there are 406,706 tables resulting in 654,670 columns, and on average a table consists of 1.61 columns and 12.79 rows. Again, these are rather small dimensions. In addition, the Turl corpus includes no numerical data at all, which leads to the shortcomings mentioned above when using the corpora.

SemTab SemTab is a yearly challenge with the goal of benchmarking systems that match tabular data to knowledge graphs since 2019. The Challenge includes the tasks of assigning a semantic type to a column, matching a cell to an entity, and assigning a property to the relationship between columns. Every year, the challenge provides different datasets to validate the participating systems against each other. In this paper we observed the provided corpora for the years 2019 [11], 2020 [5, 12], and 2021 [1, 5, 13, 16, 21]. Statistic details of the corpora are shown in Table 1. In case more than one dataset was provided per year, we aggregated the statistics over all datasets included in the challenge. While SemTab2019 consists of 13,765 tables and 21,682 columns in total, there are 131,253 tables and 190,494 columns in SemTab2020. In both corpora, the dimensions of the included tables are rather small (on average 1.58 columns and 35.61 rows in 2019 and 1.45 columns and 9.19 rows in 2020). In SemTab2021, the contained tables are the largest in terms of rows with almost 875 on average. However, the number of columns (3.86 on average) is only moderate and the corpus in general is the smallest with a total of 795 tables and 3,072 columns. Numerical data is almost nonexistent in the first two years (0.63% in 2019/0.07% in 2020), increasing to 28.24% numeric columns per table on average in 2021, which is still not comparable to the number of numeric data in real world data lakes.

GitTables [17] GitTables is a large-scale corpus of relational tables created by extracting CSV files from GitHub repositories. Table columns are labeled with semantic types from Schema.org [10] and DBpedia [2] using two different automated annotation methods (syntactically/semantically similarity matching from semantic type to column header). In this paper, we have focused on the annotations origin from DBpedia and the results of the semantic annotations method as described in the GitTables paper [17]. This leads to a corpus containing over 1.37M tables and 9.3M columns in total. Although this is by far the largest collection of data tables, the dimensions of the tables are on average only

Table 1 Corpus statistics about the number and sizes of tables. Additionally, we see the average number of textual and numerical columns per table for each existing annotated corpora and our new SportsTables corpus. This shows the absence of numerical data columns per table in most existing corpora and the dominance of textual data columns per table in all existing corpora. Instead, our new corpus SportsTables contains on average over 6 times more numerical columns than textual columns

Corpus	Tables	Cols	Cols/Table	Text. Cols/Table	Num. Cols/Table	Rows/Table
VIZNET	78,733	120,609	1.53	1.34	0.19	18.35
TURL	406,706	654,670	1.61	1.61	0	12.79
SemTab2019	13,765	21,682	1.58	1.57	0.01	35.61
SemTab2020	131,253	190,494	1.45	1.45	0.001	9.19
SemTab2021	795	3072	3.86	2.77	1.09	874.6
GitTables	1.37M	9.3M	6.82	3.62	3.2	184.66
SportsTables	1,187	24,838	20.93	2.83	18.1	246.72

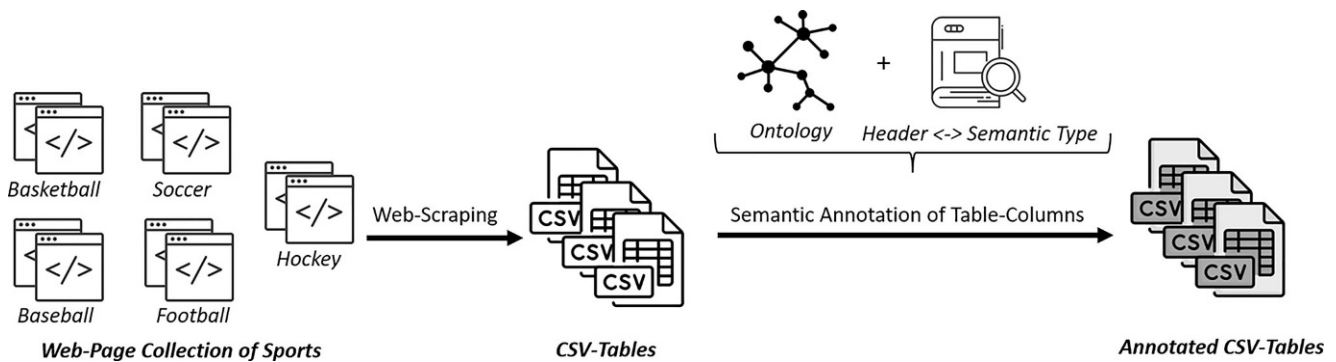


Fig. 2 Overview of the implemented pipeline to build SportsTables. We use web-scraping techniques to extract HTML tables from a manually defined web page collection for each selected sport and convert the tables to CSV files. With the help of a defined ontology and a manually created dictionary that maps column headers to semantic types, we annotate each table column with an appropriate semantic type

moderate with 6.82 columns and 184.66 rows. Overall, GitTables incorporates the most numeric data with an almost balanced ratio of 53.08% textual and 46.92% numerical columns per table.

Discussion The overview in Table 1 and the discussion before shows that most existing corpora contain no or only a minimal fraction of numerical data types which is very different from real-world data lakes. An exception is GitTables which has a much higher ratio of numerical columns. However, as we show in Sect. 4, GitTables still lacks a good coverage of different numeric semantic types which is one important aspect that we tackle with our new corpus SportsTables which covers a wide variety of different numerical semantic types. Moreover, another important (but orthogonal) aspect is that existing corpora include a large number of tables. However, on average the tables are very small in terms of the number of columns and the number of rows. Instead, our new corpus SportsTables contains fewer tables, but on average a significantly higher number of columns and rows per table to better reflect the characteristics of real-world data lakes

3 The SportsTables Corpus

In the following, we will introduce our new corpus and describe in detail the implemented construction pipeline to build SportsTables.

Methodology to Generate the Corpus Figure 2 gives an overview of our implemented pipeline to generate the new corpus. The main idea was to collect data tables from different sports domains such as soccer, basketball, baseball, etc. since data tables coming from such kinds of sources are rich in numerical columns. For example, a soccer player statistic table of a soccer season contains typically 3 textual columns (e.g., player name, team name, field position) and

18 numerical columns (e.g., goals, games played, assists). Hence, building a collection of such tables will lead to a corpus that contains many numerical columns which are in addition semantically interpretable. As a result, the corpus will enable performance analysis of semantic type prediction models in a much more rigorous manner regarding numerical data.

Scraping Data from the Web [8] A vast amount of data covering information about player statistics, team statistics, coach statistics, or season rankings of different sports are available on various web pages. Therefore, for collecting the data, we built a data collection pipeline based on web scraping technology [8]. In the first step, we manually searched and defined a set of different web pages for each of the selected sports of which we want to scrape contained data tables (left side of Fig. 2). We first converted each HTML table on the web pages to Pandas-Dataframes using Python and then saved them as CSV files (center of Fig. 2), since this file format is most known and used to store raw structured data [19]. During the scrape process, we kept the respective column headers from the original HTML table and used them as headers in the CSV file.

Annotating Columns with Semantic Types Due to the low granularity of existing ontologies (e.g. DBpedia) regarding semantics of a given sport, we manually created an ontology-like set of valid semantic types for all sports. For example, in DBpedia there is the type *Person.Athlete.BasketballPlayer*, but semantic labels in the particular that would match individual numerical columns such as *NumberOfGoals* are not defined. Next, we annotated all table columns with semantic types using a manually created dictionary that maps column headers to matching semantic types from our created set. Since the column headings were in many cases identical if the semantic content was the same, this procedure significantly reduces the manual labeling effort. In addition, to ensure

that the labels are of very high quality in terms of correctness, we manually checked each assignment based on the content of the columns.

4 Analysis of the Corpus

This section describes the characteristics of SportsTables in detail and then demonstrates the significant impact of these characteristics on semantic type prediction frameworks in a study where we apply the corpus to several state-of-the-art semantic type detection models.

4.1 Corpus Characteristics

In the following, we discuss the statistics of the SportsTables corpus and compare them to the existing corpora.

Data Statistics (Table 1) Using the described pipeline for creating SportsTables, a total of 1,187 tables which comprises 24,838 columns (approx. 86% numeric and 14% textual) are scraped from the web resulting in 20.93 columns (2.83 textual and 18.1 numerical) per table on average. This ratio of textual to numerical columns, as well as the total average number of columns in a table, differs significantly from existing corpora.

In Table 1 we can also see a comparison of the average number of textual and numerical columns per table of SportsTables versus that of the existing corpora. Here we can see that numerical columns only exist in the corpora VizNet with 0.33, SemTab2021 with 1.09, and GitTables with 3.2 columns per Table. Compared to GitTables, in SportsTables there are thus on average over 6 times more numeric columns per table. Moreover, as we discuss below, our corpus uses a much richer set of numerical data types that better reflects the characteristics in real-world data lakes which is very different from GitTables. For example, when looking at the semantic types that are assigned to numerical columns in GitTables, more than half (393,925) of the columns are labeled with just a single type *Id*.

In terms of the total number of columns, the tables in SportsTables (20.93 columns per table) are on average about 3 times wider than in GitTables (6.82 columns per table), which contains the widest tables among the existing corpora. As such, the number of columns in tables of SportsTables are reflecting better the width when comparing this to the characteristics of the tables in real-world data lakes which we analyzed. Moreover, considering the average number of rows per table, it can be seen that the tables in SportsTables have on average 246.72 rows. In comparison, tables in SportsTables are larger on average than in many other corpora where tables have typically fewer rows.

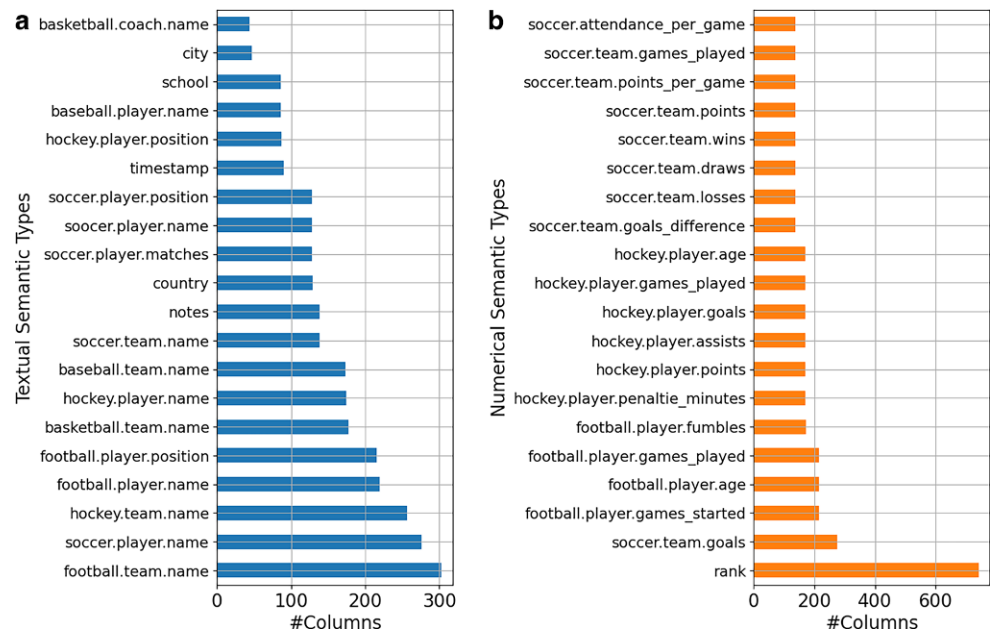
Annotation Statistics Semantic type annotation follows a two step process. First, we establish a directory with manually defined mappings from column header to semantic type for each existing header. Second, we label each column with the semantic type listed in the directory for its header. As a result, 56 textual and 419 numerical semantic types are present in the corpus. Thereby textual semantic types are those which specify textual columns and numerical types are those which specify columns containing numeric values. To compare the annotation statistics, we also counted the number of textual and numerical semantic types in an analysis of the existing corpora. The results of these analyses can be seen in Table 2. Different from our corpus, the sets of textual and numerical types are not disjoint in all other corpora (except TURL where no numeric values are present). This indicates that individual semantic types were assigned to both textual and numerical columns which is problematic if semantic type detection models should be trained and tested on these corpora. In particular, GitTables has a very large overlap and almost all semantic types are used in both column data types. To give an example, in GitTables the semantic types *comment*, *name* and *description* are assigned to both column data types. Next, we take a closer look into the semantic types of our corpus.

Figure 3a, b show the top 20 semantic types (textual and numerical) in regards to how often they were assigned to a table column. It can be seen that the most common textual types across all sports are *player.name* and *team.name*. These are types that occur in almost every table. Other types such as *country* or *city* are also common, describing, for example, the player’s origin or the team’s hometown. Among numeric semantic types, *rank* is by far the most common and is present in almost all tables. The type describes a column containing the placement of e.g., a team in a “seasons standing” table or a player in a “top scorer” table. All other numeric semantic types show mainly an equal distribution of the frequency, which is a good precondition for training machine learning models.

Table 2 Statistics about the number of unique semantic types. Showing that our new corpus has a higher proportion of numerical semantic types than textual semantic types in contrast to the existing corpora. In addition, there is a large overlap of semantic types used for textual and numeric columns in the existing corpora. In comparison, the semantic types in SportsTables are disjoint for the two column data types

Corpus	#Textual Sem. Types	#Numerical Sem. Type	#Total Sem. Types
VIZNET	78	44	78
TURL	255	0	255
SemTab2019	360	19	360
SemTab2020	5804	32	5832
SemTab2021	177	93	251
GitTables	2646	2426	2693
SportsTables	56	419	475

Fig. 3 Semantic type annotation statistics of SportsTables. **a** Shows column annotation counts of the top 20 textual semantic types. Across all kinds of sports, *player.name* and *team.name* are the most common. **b** Shows column annotation counts of the top 20 numerical semantic types. A dominant type here is *rank*, which describes a column containing the placements of e.g. a team in a season standings table



SportsTables vs. GitTables Since GitTables is the largest corpus with the most tables, one could argue that a subset of GitTables would result in a new corpus with similar characteristics as SportsTables. To analyze this, we executed a small experiment in which we filtered out only tables from GitTables where the number of textual and numerical columns (min. 3 textual and 18 numerical columns) is at least the same as it is in SportsTables. The result was a corpus containing a total of 16,909 tables and 743,432 columns. On average a table has 12.53 textual columns, 31.43 numerical columns, and 17.35 rows. However, looking at the semantic types that are assigned to numerical columns, more than half (393,925) of the columns are labeled with the type *Id*. In terms of training and validating semantic type detection models, this is rather an unfavorable type representing no semantically meaning. Moreover, the next 5 most common numerically based semantic types are *parent*, *max*, *comment*, *created* and *story editor*, constituting a large proportion of the columns. The assignment of these types to numerical data is slightly less understandable and indicates a lack of quality in the automatically generated labels for table columns.

4.2 Study of Using SportsTables

In the following, we report on the results of using four different state-of-the-art semantic type detection models on our new corpus. With this, we want to measure how well the semantic types in our corpus can be inferred by the models with a special focus on how each performs on textual and numerical columns.

Models As state-of-the-art models, we used Sherlock [15], Sato [26], Dosolo [24] and Doduo [24] in our experiments, which all use deep learning techniques to build the model. In the following, we describe the fundamental functionalities of each model to explain the differences between them.

Sherlock: Sherlock utilizes multiple feature sets, such as character distributions, word embeddings, paragraph embeddings, and column statistics (e.g., mean of numerical values), in its single-column prediction model. Each columnwise feature set, except for the column statistics, is processed by a multi-layer feature specific subnetwork to generate compact dense vectors. The resulting outputs from the subnetworks, along with the column statistics features, are then inputted into the primary network, which comprises two fully connected layers.

Sato: Building upon Sherlock, Sato is a multi-column prediction model that incorporates LDA features to capture table context and integrates a CRF layer to account for column type dependency in its predictions. With this, Sato's prediction quality improves over Sherlock on the VizNet data corpus.

Dosolo & Doduo: Dosolo & Doduo are both models from [24] and use pre-trained language models (LM) (e.g. BERT) combined with an attached output layer to implement a model for the semantic type classification task. Given that LMs receive token sequences (i.e. text) as input, it is essential to convert a table into a token sequence so that the LM can process it. What distinguishes Dosolo and Doduo is the way in which a table is serialized into a token sequence. Dosolo implements a columnwise serialization where each column C and its values v_1, \dots, v_m of a table is separately serialized as follows: $serialize(C) ::= [CLS]v_1, \dots, v_m[SEP]$. In con-

trast, Doduo is a tablewise model designed to process an entire table as input. To accomplish this, Doduo serializes the complete table and its entries as follows: for each table that has n columns $T = (c_i)_{i=1}^n$, where each column has N_m column values $c_i = (v_i^j)_{j=1}^{N_m}$, they let $serialize(T) ::= [CLS]v_1^n \dots [CLS]v_1^n \dots v_m^n [SEP]$. In both sequences, the special token [CLS] marks the beginning of a new table column and [SEP] the end of a token sequence. The major difference between the two approaches and their serialization techniques is that with Dosolo a column type is predicted independently of other data in the table (e.g. neighboring column values), whereas Doduo model captures the table context to make a prediction of a column type. In summary, we can conclude that Sherlock and Dosolo are single-column (columnwise) prediction models that only take into account the individual column values for the prediction. In contrast, Sato and Doduo are multi-column (or tablewise) prediction models, which consider table contexts for predicting the semantic type of an individual column.

Experiment Setup For the experiments, we split the SportsTables corpus into training, validation, and test set. While creating the splits, we first extracted 20% of the data for the test set and then another 20% of the remaining 80% for the validation split. The rest of the data was used as the training set. We used the four pre-trained models as described above and re-trained them with the training data set. During the re-training, we replaced the last layer of the different models to support the number of semantic types that occur in SportsTables and then re-trained the entire neural network. In order to optimize the hyperparameters, we measured the performance of the respective re-trained models against the validation split. To report the final performance, we applied the re-trained models to the 20% test data set. For obtaining statistically reliable results, we ran each experiment with five different random seeds and report the mean and standard deviation over multiple runs.

Results of Study Figure 4 shows the results of the experiments reporting the support weighted and macro average F1-Scores in individual subplots for all four models. For each model, we plot the F1-Score across all semantic types (numerical & non-numerical) to show the total performance, but also the separate average F1-Score for only textually and numerically based semantic types, respectively. In the following we want to discuss the main aspects of the results in detail.

Non-numeric vs. numeric: As we can see in the figure, there is a significant performance difference between predicting textual and numerical semantic types for all models. While textual columns can be predicted with performances in a very promising range of 0.82–0.98, the performances

for numerical columns are rather moderate ranging from 0.31–0.7. On average, the difference in F1-Score between textual and numeric types is 0.35 across all models. These results demonstrate that the models can better handle textual data and determine its associated semantic types more accurately than numerical data. Looking at the total performances over all types for each model, we see that they are rather moderate in the range of 0.38–0.74, but these insufficient results are primarily caused by poor prediction performances on the numerical based types.

Columnwise vs. tablewise: Looking and comparing the results of the columnwise models Sherlock & Dosolo and the results of the tablewise models Sato & Doduo, we observe that the tablewise models outperform the columnwise models. The results underline the known importance of considering not only individual column values for the task of semantic type detection of table columns but also to involve the table context. In particular, what we can see from the comparison of Dosolo and Doduo is how important it is, especially for numerical based columns, to include table context data for semantic type detection. As described above, numerical values provide less information content that can be used by a machine learning model to identify the type and therefore Doduo doubles the performance of Dosolo by considering the complete table context. However, the resulting performance of 0.62 is rather moderate and

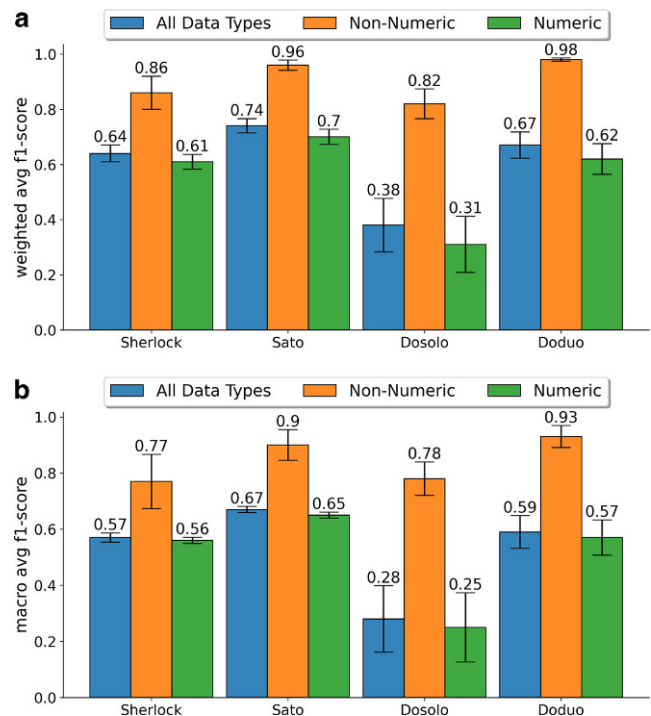


Fig. 4 Results using different state-of-the-art semantic type detection models on our new SportsTables corpus. The overall differences in F1-Scores for predicting textual and numeric columns indicate that the models can handle textual data more effectively than numeric data

demonstrates the shortcomings of the model on numerical semantic types. Comparing Sherlock and Sato also reflects the advantages of a tablewise semantic column type detection, whereas the performance improvement on just numerical columns is not as significant as in Dosolo vs. Doduo. We will discuss the reasons for this in the following.

Sherlock & Sato vs. Dosolo & Doduo on numeric: As described above, Sherlock & Sato (same feature set) as well as Dosolo & Doduo (same LM model) are models with an identical foundation. Focusing only on the F1-Scores on the numerical types, one can see that Sherlock & Sato outperform Dosolo & Doduo. We think that this aspect is due to the fact that Sherlock & Sato extract features of numerical columns that better address numerical data (e.g. mean of individual digits occurring in a column), while in Dosolo & Doduo a LM model is used as the basis to encode the numerical column. LM models are optimized for text and can therefore not provide a representative encoding to infer the semantic type of numerical columns. Therefore, Dosolo & Doduo predictions on numerical columns are inferior to Sherlock & Sato.

5 Future Challenges

Detecting semantic types in real-world data lakes comes with many challenges that need to be addressed. In particular, based on our findings of the analysis using the different models in Sect. 4, we think that new model architectures are needed for detecting numerical data types which have different characteristics from non-numerical data. In the following, we list some of the challenges we think are important to be addressed. We hope that our corpus enables research on those challenges.

Embedding numerical data: Most state-of-the-art models like Dosolo&Doduo apply LMs like BERT [7] to encode literals to infer the semantic type of a table column. Since such approaches are optimized for textual data, using them on numerical data is not sufficient, as our experimental results show. Therefore, we need improved embeddings for numerical data, which can now be studied with SportsTables.

Leveraging numerical context: To improve the semantic type prediction of a table column, recent approaches like Sato [26], TURL [6] and Doduo [24] incorporate also context information like the table-topic or values from neighboring columns of the same table as described above. Given that tables in existing corpora contain almost entirely textual columns, the contexts (e.g. values from neighboring columns) used are rich in information and therefore also lead to performance improvements. However, it is unclear how effective this approach is in case the tables contain many numerical columns and only a few textual columns

since the context information provided is reduced due to the lower entropy of numeric values as described before. Our first results using SportsTables show that in principle adding context information leads to an improvement on non-numeric and numeric types. However, we believe that specifically for the prediction of the numerical types the leverage of contextual informations needs to be researched in more depth.

Supporting wide tables: Existing datasets for semantic type detection consist of tables with small numbers of columns and rows. In nearly all corpora, the existing tables contain on average less than two columns and less than 40 rows (see Table 1). Therefore, at the current state, it has not been analyzed how state-of-the-art models can handle such large tables. To give an example of why large tables could be a problem for recent models, we will briefly discuss this aspect on the Doduo model. As described above, Doduo uses pre-trained LMs (e.g., BERT) and hence serializes the entire table into token sequences with a fixed tensor length of 512 elements, which is given by the LM model. With this methodology of serialization and the fixed given tensor length, increasing the number of table columns means that decreasing number of values of each column can be included for serialization. For example, a table with 512 columns would allow only one value per column to be considered and this would most likely result in an insufficient semantic representation of the column based on that one value.

6 Conclusion

Existing corpora for training and validating semantic type detection models mainly contain tables with either only or a very high proportion of textual data columns and either no or just a limited number of numerical data columns. Therefore, it has not been studied precisely how well state-of-the-art models perform on a dataset with a very high percentage of numerical columns as it occurs in real-world data lakes. Moreover, tables in existing corpora are very small regarding the total number of columns and rows. To tackle these shortcomings, we built a new corpus called SportsTables which contains tables that have on average approx. 3 textual columns, 18 numerical columns, and 250 rows. With our new corpus, semantic type detection models for table columns can now be holistically validated against numerical data. We show results by using the different state-of-the-art semantic type detection approaches Sherlock, Sato, Dosolo, and Doduo on our new corpus and report significant differences in the performance of predicting semantic types of textual data and numerical data on all models. Finally, we think that the corpus is just a first step to stimulate more research on new model architectures that can better

deal with numerical and non-numerical data types. The corpus is available on <https://github.com/DHBWMosbachWI/SportsTables.git>.

Acknowledgements We thank the reviewers for their feedback. This research is funded by the BMBF projects AICoM and KompAKI (grant numbers 02P20A064 and 02L19C150) by the state of Hesse as part of the NHR Program, as well as the HMWK cluster project 3AI (The Third Wave of AI). Finally, we want to thank DHBW Mosbach, hes-dian.AI, TU Darmstadt as well as DFKI Darmstadt for their support.

Funding Open Access funding enabled and organized by Projekt DEAL.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Abdelmageed N, Schindler S, König-Ries B (2021) fusion-jena/biodivtab <https://doi.org/10.5281/zenodo.5584180>
- Auer S, Bizer C, Kobilarov G et al (2007) Dbpedia: A nucleus for a web of open data. In: Aberer K, Choi KS, Noy N, al (eds) The Semantic Web. Springer Berlin Heidelberg, Berlin, Heidelberg, pp 722–735
- Bhagavatula CS, Noraset T, Downey D (2015) Tabel: Entity linking in web tables
- Cafarella MJ, Halevy A, Wang DZ et al (2008) Webtables: Exploring the power of tables on the web <https://doi.org/10.14778/1453856.1453916>
- Cutrona V, Bianchi F, Jiménez-Ruiz E et al (2020) Tough Tables: Carefully Evaluating Entity Linking for Tabular Data <https://doi.org/10.5281/zenodo.4246370>
- Deng X, Sun H, Lees A et al (2021) TURL: Table Understanding through Representation Learning <https://doi.org/10.14778/3430915.3430921> (<https://github.com/sunlab-osu/TURL>)
- Devlin J, Chang MW, Lee K et al (2019) BERT: Pre-training of deep bidirectional transformers for language understanding <https://doi.org/10.18653/v1/N19-1423> (<https://aclanthology.org/N19-1423>)
- Diouf R, Sarr EN, Sall O et al (2019) Web scraping: State-of-the-art and areas of application <https://doi.org/10.1109/BigData47090.2019.9005594>
- Google (2022) Freebase data dumps. <https://developers.google.com/freebase>
- Guha RV, Brickley D, Macbeth S (2016) Schema.org: Evolution of structured data on the web. *Commun ACM* 59(2):44–51. <https://doi.org/10.1145/2844544>
- Hassanzadeh O, Efthymiou V, Chen J et al (2019) SemTab 2019: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching Data Sets <https://doi.org/10.5281/zenodo.3518539>
- Hassanzadeh O, Efthymiou V, Chen J et al (2020) SemTab 2020: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching Data Sets <https://doi.org/10.5281/zenodo.4282879>
- Hassanzadeh O, Efthymiou V, Chen J et al (2021) SemTab 2021: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching Data Sets <https://doi.org/10.5281/zenodo.6154708>
- Hu K, Gaikwad SNS, Hulsebos M et al (2019) Viznet: Towards a large-scale visualization learning and benchmarking repository <https://doi.org/10.1145/3290605.3300892>
- Hulsebos M, Hu K, Bakker M et al (2019) Sherlock: A deep learning approach to semantic data type detection <https://doi.org/10.1145/3292500.3330993>
- Hulsebos M, Demiralp C, Demiralp P (2021a) Gittables benchmark - column type detection <https://doi.org/10.5281/zenodo.5706316>
- Hulsebos M, Demiralp Ç, Groth P (2021b) Gittables: A large-scale corpus of relational tables. *CoRR* abs/2106.07258. <https://arxiv.org/abs/2106.07258>. Accessed 24 April 2023.
- Langenecker S, Sturm C, Schalles C et al (2023) Sportstables: A new corpus for semantic type detection. In: König-Ries B, Scherzinger S, Lehner W et al (eds) BTW 2023. Gesellschaft für Informatik e.V., <https://doi.org/10.18420/BTW2023-68>
- Mitlöhner J, Neumaier S, Umbrich J et al (2016) Characteristics of open data csv files. In: IEEE (ed) International Conference on Open and Big Data (OBD). IEEE, pp 72–79 <https://doi.org/10.1109/OBD.2016.18>
- Neumaier S, Umbrich J, Polleres A (2016) Automated quality assessment of metadata across open data portals. *J Data Inf Qual.* <https://doi.org/10.1145/2964909>
- Oliveira D, Pesquita C (2021) Semtab 2021 biotable dataset <https://doi.org/10.5281/zenodo.5606585>
- Plotly (2018) Plotly. <https://chart-studio.plotly.com/feed/>. Accessed 24 April 2023.
- Shannon CE (1948) A mathematical theory of communication. *Bell Syst Tech J* 27:379–423 (<http://plan9.bell-labs.com/cm/ms/what/shannonday/shannon1948.pdf>)
- Suhara Y, Li J, Li Y et al (2022) Annotating columns with pre-trained language models
- Viegas FB, Wattenberg M, van Ham F et al (2007) Manyeyes: A site for visualization at internet scale. *IEEE Trans Visual Comput Graphics* 13(6):1121–1128. <https://doi.org/10.1109/TVCG.2007.70577>
- Zhang D, Hulsebos M, Suhara Y et al (2020) Sato: Contextual semantic type detection in tables <https://doi.org/10.14778/3407790.3407793>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.