

ABSTAT 1.0: Compute, Manage and Share Semantic Profiles of RDF Knowledge Graphs

Renzo Arturo Alva Principe, Blerina Spahiu, Matteo Palmonari, Anisa Rula,
Flavio De Paoli, Andrea Maurino

University of Milano-Bicocca

{renzo.alvaprincipe|blerina.spahiu|matteo.palmonari|anisa.rula|
flavio.depaoli|andrea.maurino}@unimib.it

Abstract. As Linked Data available on the Web continue to grow, understanding their structure and content remains a challenging task making such the bottleneck for their reuse. ABSTAT is an online profiling tool which helps data consumers in better understanding the data by extracting ontology-driven patterns and statistics about the data. This demo paper presents the capabilities of the new added feature of ABSTAT.

1 Introduction

Knowledge Graphs (KGs) in the Linked Open Data cloud¹ define possible classes and relations in a schema or ontology, and mainly describe instances and inter-link entities through relations. KGs cover different domains and are widespread, for example, in the EuBusinessGraph project², several parties contribute their data into the KG of the company. Despite the gross amount of data available on the Web, the selection of the data suitable for a given task is not straightforward as many data discovery steps have to be performed in order to understand data set's content and their characteristics. Thus, in order to use a data set, one needs to know which classes and properties are most commonly used, which predicates are generally associated with an instance of a given class, the potential domain and range of a given predicate, the cardinality of a predicate, etc. ABSTAT is an ontology-driven linked data summarization model which helps users in an effortless understanding of the data [5]. Given a RDF data set and, optionally, an ontology (used in the data set), ABSTAT computes a semantic profile which consists of a summary and statistics. ABSTAT's summary is a collection of patterns known as Abstract Knowledge Patterns (AKPs) of the form `<subjectType, pred, objectType>`, which represent the occurrence of triples `<sub, pred, obj>` in the data, such that `subjectType` is a minimal type of the subject and `objectType` is a minimal type of the object. With the term type we refer to either an ontology class (e.g., `foaf:Person`) or a datatype (e.g., `xsd:DateTime`). By considering only minimal types of resources, computed with

¹ <http://lod-cloud.net/>

² <http://eubusinessgraph.eu/>

the help of the data ontology, we exclude several redundant AKPs from the summary making them compact and complete. Summaries are published and made accessible via web interfaces, in such a way that the information that they contain can be consumed by users and machines (via APIs). The user interface is available and can be used to explore summarized datasets³. Several approaches to profile RDF data have been proposed, we refer to our research papers [5, 1] for a detailed discussion of state-of-the-art. While many of these approaches publish and make accessible the computed profiles, only a few are open source and, to the best of our knowledge, none of them provide support for the summarization process to the user. Based on requirements collected in the two industry-driven innovation projects EW-Shopp⁴ and EuBusinessGraph we have built ABSTAT 1.0, a tool to compute, manage and make accessible to humans and machines semantic profiles of RDF graphs. Compared to the ABSTAT research prototype [2], ABSTAT 1.0 not only provides more features, which are used in different applications scenarios [5, 1, 3] but it has also developed into a tool that lays on a more scalable modular and effective architecture, and is endowed with a user interface to help the management of the profiling process. ABSTAT 1.0 is released as open source⁵ under the GNU Affero General Public License v3.0⁶.

In this paper, we make the following contributions: (i) Minimalization over properties; (ii) AKPs inference and instance count; (iii) Cardinality extraction; (iv) Configuration and launch of the summarization via GUI; (v) Indexing of summaries via GUI; (vi) Browsing and full-text search; (vii) Access to summaries via APIs (viii) Autocomplete service over arbitrary strings.

2 Exploring and understanding a data set with ABSTAT

ABSTAT controller⁷ is designed to be modular and decoupled as in Fig. 1. The modules of ABSTAT 1.0 are the following:

- **ABSTAT Viewer** provides a graphic user interface to serve different types of tasks such as summary exploration, execution of the summarization process using a wizard and summaries indexing. Summary exploration can be performed using constrained queries (a desired subject and/or predicate and/or object) and full-text search. The summarization wizard provides a GUI to let users select datasets/ontologies from a populated list or using an upload module, configure and execute the summarization process. After the semantic profile is computed, the user can load/index it on a persistent storage/search engine in order to support its access through APIs or GUI.
- **ABSTAT Builder** is the module that executes the summarization algorithms and produces the profiles. The Summarizator component requires as

³ <http://abstat.disco.unimib.it>

⁴ <http://www.ew-shopp.eu>

⁵ https://bitbucket.org/disco_unimib/abstat

⁶ <https://www.gnu.org/licenses/>

⁷ <http://backend.abstat.disco.unimib.it>

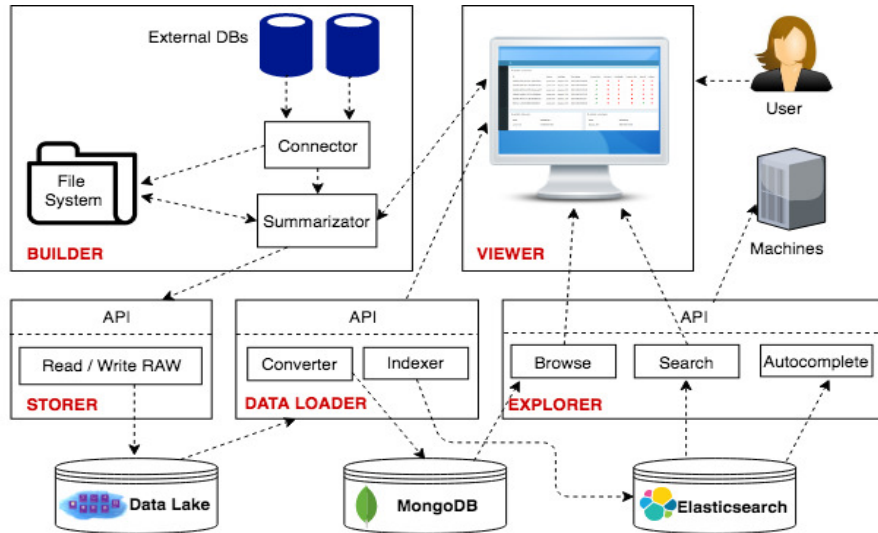


Fig. 1. ABSTAT Architecture

input a dataset (in N3 format) and an ontology (in OWL format) along with the configuration chosen by the user. If the data are in an external DB, the Connector component allows extracting a dump and storing it in the correct file to serve as input to the Summarizator.

- **ABSTAT Storer** component feeds a data lake storage with the raw data produced by the Builder. It also receives download requests from users who want to get raw summaries.
- **ABSTAT Loader** contains the Converter component, which converts the data formats in the Data Lake in a format suitable for the Explorer module. The Indexer component indexes summaries in a search engine. Note that the Loader component receives the control input from the Viewer.
- **ABSTAT Explorer** is organized as a set of APIs to satisfy profile exploration requests from Viewer or users who want to use them directly.

3 Demonstration

ABSTAT is a framework that computes and provides access to semantic profiles that consist in an RDF summary and statistics. The summary of a data set describes its content by listing every schema-level pattern that occur in the data. In addition, semantic profiles provide several statistics about the occurrence of patterns, types and properties and cardinality statistics. During the summarization process if the user specifies the main pay-level domain of the data set (e.g., dbpedia.org for DBpedia), ABSTAT can distinguish between resources (patterns, types and properties) that are internal (resources having the specified pay-level domain) and external (resources having a pay-level domain different from the

one specified by the user). This distinction has the only purpose of letting users filter out patterns that include some external resource (e.g., hide all patterns that contain the type foaf:Person when looking at patterns extracted from DBpedia).

browser the abstract knowledge patterns found into: <http://id.summaries.org/dbpedia/2014>

subject type (occurrences)	predicate (occurrences)	object type (occurrences)	frequency	instances	Max subj-obj	Avg subj-obj	Min subj-obj	Max subj-objs	Avg subj-objs	Min subj-objs
dbo:Actor	predicate	object								
dbo:Actor (4695)	DTP foaf:name (500230)	rdfs:Literal (1244446)	12072	15989	1159	1	1	8	3	1
dbo:Actor (4695)	DTP dbo:description (236198)	rdfs:Literal (1244446)	8382	12098	1006	7	1	2	2	1
dbo:Actor (4695)	DTP foaf:givenName (113449)	rdfs:Literal (1244446)	4036	5698	44	2	1	1	1	1
dbo:Actor (4695)	DTP foaf:surname (1104316)	rdfs:Literal (1244446)	4036	5698	40	1	1	1	1	1
dbo:Actor (4695)	DTP dbo:birthYear (1047864)	xmld:Year (3689462)	3977	5754	83	19	1	2	1	1
dbo:Actor (4695)	DTP dbo:birthDate (101416)	xmld:Date (184226)	3515	5007	3	1	1	2	1	1
dbo:Actor (4695)	QP dbo:country (26876)	dbo:Country (2684)	2640	2685	1072	710	472	1	1	1
dbo:Actor (4695)	DTP dbo:activeYearStartYear (236146)	xmld:Year (3689462)	2151	2298	76	18	1	1	1	1
dbo:Actor (4695)	QP dbo:occupation (313790)	owl:Thing	2050	4390	1231	16	1	11	2	1
dbo:Actor (4695)	QP dbo:birthPlace (116849)	dbo:Country (2684)	1929	2958	361	24	1	2	1	1
dbo:Actor (4695)	QP dbo:occupation (313790)	dbo:PersonFunction (18355)	1698	1993	1	1	1	5	1	1
dbo:Actor (4695)	QP dbo:birthPlace (116849)	dbo:Settlement (238436)	1530	3523	246	3	1	4	1	1
dbo:Actor (4695)	DTP dbo:deathYear (46661)	xmld:Year (3689462)	1355	1494	34	9	1	2	1	1
dbo:Actor (4695)	DTP dbo:imdbid (7646)	rdfs:Literal (1244446)	1180	1528	1	1	1	2	1	1
dbo:Actor (4695)	DTP dbo:deathDate (46726)	xmld:Date (184226)	1057	1200	2	1	1	2	1	1
dbo:Actor (4695)	DTP dbo:activeYearEndYear (140123)	xmld:Year (3689462)	1031	1051	29	10	1	1	1	1
dbo:Actor (4695)	DTP dbo:alias (15373)	rdfs:Literal (1244446)	933	1198	10	1	1	5	1	1

Fig. 2. ABSTAT Browse GUI

Fig. 2 shows the home page of ABSTAT. The menu on the left side can be used to explore semantic profiles. The **Overview** page gives an overview of the uploaded data sets, ontologies and computed profiles. **Summarize** page gives a configuration interface for custom summarizations including data sets and ontologies uploading. **Consolidate** allows to persist and index the computed profiles into the search engine. **Browse** is the GUI for constraint-based pattern exploration. **Search** is the GUI for full-text searching. Patterns, predicates and types that match the keyword will be returned. Search can be processed over the whole set of indexed profiles or on those originated from a specific data sets. Statistics, data set names and pattern symbols will be shown in the results of the query. **Manage** allows to remove data sets, ontologies and profiles. **APIs** lists the available APIs for machine-friendly profile exploration.

Patterns of the semantic profile are sorted by frequency in descendant order. The user can also put constraints on subjects and/or predicates and/or objects. In every text box a simple suggestion menu will recommend types/predicates that occur in the patterns. Then patterns are filtered in order to match the user constraints. Fig. 3 shows the patterns that match the predicate `dbo:knownFor` and the object type `dbo:Film`. For each pattern several statistics are returned. Considering the one in the black box, the frequency of the pattern shows how many times does this pattern occur in the data set. The number of instances shows how many instances have this pattern including those for which the types `Person` and `Film` and the predicate `knownFor` can be inferred. Max (Min, Avg) subj-obj cardinality is the maximal (minimal, average) number of distinct entities of type `Person` linked to a single entity of type `Film` through the predicate `knownFor`. Max (Min, Avg) subj-objs is the maximal (minimal, average) number of distinct entities of type `Film` linked to a single entity of type `Person` through the predicate `knownFor`. Frequency is given also for types and predicates.

Previous experiments suggest that ABSTAT summaries help users in understanding a data set, e.g., by facilitating query formulation, and provide support to the assessment of data quality by finding outliers in the vocabulary usage

	subject type (occurrences)	predicate (occurrences)	object type (occurrences)	frequency	instances	Max subjs-obj	Avg subjs-obj	Min subjs-obj	Max subj-objs	Avg subj-objs	Min subj-objs
filter	subject	dbo:knownFor	dbo:Film								
☆	dbo:Person (611330)	OP dbo:knownFor (41404)	dbo:Film (101906)	1160	1208	9	1	1	12	2	1
	foaf:Person (1179233)	OP dbo:knownFor (41404)	dbo:Film (101906)	1066	1066	5	1	1	12	2	1
☆	dbo:Actor (4669)	OP dbo:knownFor (41404)	dbo:Film (101906)	31	31	1	1	1	4	2	1
☆	dbo:ScreenWriter (700)	OP dbo:knownFor	dbo:Film (101906)	11	11	1	1	1	5	2	1

Fig. 3. Semantic Profile of DBpedia 2014 data set

[5]. In addition, we have recently found that rich profiles as the ones computed in ABSTAT 1.0 support automatic feature selection for semantic recommender systems, outperforming other purely statistical measures like Information Gain [3, 1]. Finally, ABSTAT 1.0 supports vocabulary suggestions, similarly to [4]. In the future, ABSTAT will provide more significant statistics such statistics about class hierarchy depth, classes and properties per entity, etc.

Acknowledgements

This research has been supported in part by EU H2020 projects EW-Shopp - Grant n. 732590, and EuBusinessGraph - Grant n. 732003.

References

- [1] T. di Noia, A. Maurino, C. Magarelli, M. Palmonari, and A. Rula. Using ontology-based data summarization to develop semantics-aware recommender systems. In *The Semantic Web - ESWC 2018 Satellite Events, Heraklion, Crete, Greece, June 3-7, 2018*, 2018.
- [2] M. Palmonari, A. Rula, R. Porrini, A. Maurino, B. Spahiu, and V. Ferme. ABSTAT: linked data summaries with abstraction and statistics. In *The Semantic Web: ESWC 2015 Satellite Events - ESWC 2015 Satellite Events Portorož, Slovenia, May 31 - June 4, 2015, Revised Selected Papers*, pages 128–132, 2015.
- [3] A. Ragone, P. Tomeo, C. Magarelli, T. D. Noia, M. Palmonari, A. Maurino, and E. D. Sciascio. Schema-summarization in linked-data-based feature selection for recommender systems. In *Proceedings of the Symposium on Applied Computing, SAC 2017, Marrakech, Morocco, April 3-7, 2017*, pages 330–335, 2017.
- [4] J. Schaible, T. Gottron, and A. Scherp. Termpicker: enabling the reuse of vocabulary terms by exploiting data from the linked open data cloud. In *International Semantic Web Conference*, pages 101–117. Springer, 2016.
- [5] B. Spahiu, R. Porrini, M. Palmonari, A. Rula, and A. Maurino. ABSTAT: ontology-driven linked data summaries with pattern minimalization. In *The Semantic Web - ESWC 2016 Satellite Events, Heraklion, Crete, Greece, May 29 - June 2, 2016, Revised Selected Papers*, pages 381–395, 2016.