Benchmark Datasets for the Semantic Web: A Collection for Systematic Machine Learning Evaluations

¹RACHITA RAUL, Gandhi Institute of Excellent Technocrats, Bhubaneswar, India ²PRIYANKA MOHANTY, Indus College of Engineering, Bhubaneswar, Odisha, India

Abstract. Several methods for machine learning on the Semantic Web have been put out in recent years. However, because to a dearth of publicly accessible, recognised benchmark datasets, no in-depth comparisons of those methodologies have been made. We present a collection of 22 benchmark datasets in this work, ranging in size from small to large. Such a set of datasets can be utilised to carry out systematic comparisons of methods and quantitative performance assessment.

Keywords: Linked Open Data Machine learning Datasets Benchmarking

1 Introduction

In the recent years, applying machine learning to Semantic Web data has drawn a lot of attention. Many approaches have been proposed for different tasks at hand, ranging from reformulating machine learning problems on the Semantic Web as traditional, propositional machine learning tasks to developing entirely novel algorithms. However, systematic comparative evaluations of different approaches are scarce; approaches are rather evaluated on a handful of often project-specific datasets, and compared to a baseline and/or one or two other systems.

In contrast, evaluations in the machine learning area are often more rigorous. Approaches are usually compared using a larger number of standard datasets, most often from the UCI repository¹. With a larger set of datasets used in the evaluation, statements about statistical significance are possible as well [3].

At the same time, collections of benchmark datasets have become quite well accepted in other areas of Semantic Web research. Notable examples include the Ontology Alignment Evaluation Initiative (OAEI) for ontology matching², the *Berlin SPARQL Benchmark* ³ for triple store performance, the Lehigh University Benchmark (LUBM)⁴ for reasoning, or the Question Answering over Linked Data (QALD) dataset⁵ for natural language query systems.

In this paper, we introduce a collection of datasets for benchmarking machine learning approaches for the Semantic Web. Those datasets are either existing RDF datasets, or external classification or regression problems, for which the instances have been enriched with links to the Linked Open Data cloud [14]. Furthermore, by varying the number of instances for a dataset, scalability eval- uations are also made possible.

2 Related Work

Recent surveys on the use of Semantic Web for machine learning organize the proposed approaches in several categories, i.e., approaches that use Semantic Web data for machine learning [16], approaches that perform machine learning on the Semantic Web [11], and approaches that use machine learning techniques to create and improve Semantic Web data [8,16]. Furthermore, there are some challenges, like the *Linked Data Mining Challenge*⁶ or the *Semantic-Web enabled Recommender Systems Challenge*⁷, which usually focus on only a few datasets and a very specific problem setting.

3 Datasets

Our dataset collection has three categories: (i) existing datasets that are com- monly used in machine learning experiments, (ii) datasets that were gener- ated from official observations, and (iii) datasets generated from existing RDF datasets. Each of the datasets in the first two categories are initially linked to DBpedia⁸. This has two main reasons, (1) DBpedia being a cross-domain knowl- edge base usable in datasets from very different topical domains, and (2) tools like DBpedia Lookup and DBpedia Spotlight making it easy to link external datasets to DBpedia. However, DBpedia can be seen as an entry point to the Web of Linked Data, with many datasets linking to and from DBpedia. In fact, we use the RapidMiner Linked Open Data extension [9], to retrieve external links for each entity to YAGO⁹ and Wikidata¹⁰. Such links could be exploited for systematic evaluation of the relevance of the data of different LOD datasetin different learning tasks.

In the dataset collection, there are four datasets that are commonly used for machine learning. For these datasets, we first enrich the instances with links to LOD datasets, and reuse the already defined target variable to perform machine learning experiments:

- The Auto MPG dataset¹¹ captures different characteristics of cars, and the target is to predict the fuel consumption (MPG) as a regression task.

- The AAUP (American Association of University Professors) dataset contains

a list of universities, including eight target variables describing the salary of different staff at the universities¹². We use the average salary as a target variable both for regression and classification, discretizing

the target variable into "high", "medium" and "low", using equal frequency binning.

- The *Auto* 93 dataset¹³ captures different characteristics of cars, and the target

is to predict the price of the vehicles as a regression task.

- The *Zoo* dataset captures different characteristics of animals, and the target is to predict the type of the animals as a classification task.

For those datasets, cars, universities, and animals are linked to DBpedia basedon their name.

The second category of datasets contains a list of datasets where the target variable is an observation from different real-world domains, as captured by official sources. Again, the instances were enriched with links to LOD datasets. There are thirteen datasets in this category:

- The *Forbes* dataset contains a list of companies including several features of the companies, which was generated from the Forbes list of leading companies 2015¹⁴. The target is to predict the company's market value as a classifica-tion and regression task. To use it for the task of classification we discretize the target variable into "high", "medium", and "low", using equal frequency binning.

- The *Cities* dataset contains a list of cities and their quality of living, as captured by Mercer [7]. We use the dataset both for regression and classification.

- The *Endangered Species* dataset classifies animals into endangered species¹⁵.

- The *Facebook Movies* dataset contains a list of movies and the number of Facebook likes for each movie¹⁶. We first selected 10, 000 movies from DBpe- dia, which were then linked to the corresponding Facebook page, based on the movie's name and the director. The final dataset contains 1, 600 movies, which was created by first ordering the list of movies based on the number of Facebook likes, and then selecting the top 800 movies and the bottom 800 movies. We use the dataset for regression and classification.

- Similarly, the *Facebook Books* dataset contains a list of books and the number of Facebook likes. Each book was linked to the corresponding Facebook page using the book's title and the book's author. Again, we selected the top 800 books and the bottom 800 books, based on the number of Facebook likes.

- The *Metacritic Movies* dataset is retrieved from Metacritic.com¹⁷, which con-

tains an average rating of all time reviews for a list of movies [12]. The initial dataset contained around 10, 000 movies, from which we selected 1, 000 movies from the top of the list, and 1, 000 movies from the bottom of the list. We use the dataset both for regression and classification.

- Similarly, the *Metacritic Albums* dataset is retrieved from Metacritic.com¹⁸,

which contains an average rating of all time reviews for a list of albums [13].

- The *HIV Deaths Country* dataset contains a list of countries with the number of deaths caused by HIV, as captured by the World Health Organization¹⁹. We use the dataset both for regression and classification.

- Similarly, the *Traffic Accidents Deaths Country* dataset contains a list of countries with the number of deaths caused by traffic accidents²⁰.

- The *Energy Savings Country* dataset contains a list of countries with the total amount of energy savings of primary energy in 2010²¹, which was downloaded from WorldBank²². We use the dataset both for regression and classification.

- Similarly, the *Inflation Country* dataset contains a list of countries with the inflation rate for 2011²³.

- The *Scientific Journals Country* dataset contains a list of countries with a number of scientific and technical journal articles published in 2011^{24} .

- The *Unemployment French Region* dataset contains a list of regions in France with the unemployment rate, used in the SemStats 2013 challenge [10].

Again, for those datasets, the instances (cities, countries, etc.) are linked to DBpedia. For datasets which are used for classification and regression, the regression target was discretized using equal frequency binning, usually into a *high* and a *low* class.

The third, and final, category contains datasets that were generated from existing RDF datasets, where the value of a certain property is used as a classi- fication target. There are five datasets in this category:

- The *Drug-Food Interaction* dataset contains a list of drug-recipe pairs and their interaction, i.e., "negative" and "neutral" [6]. The dataset was retrieved from FinkiLOD²⁵. Furthermore, each drug is linked to DrugBank²⁶. We drew a stratified random sample of 2, 000 instances from the complete dataset. When generating the features, we ignore the foodInteraction property in DrugBank, since it highly correlates with the target variable.

- The *AIFB* dataset describes the AIFB research institute in terms of its staff, research group, and publications. In [1] the dataset was first used to predict the affiliation (i.e., research group) for people in the dataset. The dataset con- tains 178 members of a research group, however the smallest group contains only 4 people, which is removed from the dataset, leaving 4 classes. Also, we remove the employs relation, which is the inverse of the *affiliation* relation.

- The AM dataset contains information about artifacts in the Amsterdam Museum [2]. Each artifact in the

ISSN: 2278-4632 Vol-10 Issue-09 No.03 September 2020

dataset is linked to other artifacts and details about its production, material, and content. It also has an artifact category, which serves as a prediction target. We have drawn a stratified ran- dom sample of 1,000 instances from the complete dataset. We also removed the material relation, since it highly correlates with the artifact category.

- The MUTAG dataset is distributed as an example dataset for the DL-Learner toolkit²⁷. It contains information about complex molecules that are poten- tially carcinogenic, which is given by the isMutagenic property.

- The *BGS* dataset was created by the British Geological Survey and describes geological measurements in Great Britain²⁸. It was used in [17] to predict the lithogenesis property of named rock units. The dataset contains 146 named rock units with a lithogenesis, from which we use the two largest classes.

An overview of the datasets is given in Tables 1, 2, and 3. For each dataset, we depict the number of instances, the machine learning tasks in which the dataset is used (*C* stands for classification and *R* stands for regression), the source of the dataset, and the LOD datasets to which the dataset is linked. For each dataset, we depict basic statistics of the properties of the LOD datasets, i.e., average, median, maximum and minimum number of *types, categories, outgoing relations* (rel out), *incoming relations* (rel in), outgoing relations including values (rel-vals out) and incoming relations including values (rel-vals in). The datasets, as well as a detailed description, a link quality evaluation, and licensing information, can be found online²⁹.

From the given statistics, we can infer the following observations: (i) DBpedia contains significantly less *owl:sameAs* links to YAGO, compared to Wikidata;

(ii) DBpedia provides the highest number of types and categories on average per entity; (iii) Wikidata contains the highest number of outgoing and incoming relations for most of the datasets; (iv) YAGO contains the highest number of outgoing and incoming relations values for most of the datasets.

Dataset						type	s		catego	ries				rel c	out			rel i	n		rel-vals	out			rel-vals in			
Name	Source	Task	LOD	#link s	avg	med	max	min	avg	med	max	min	avg	med	max	min	avg	med	max	min	avg	med	max	min	avg	med	max	min
Auto MPG	UCI ML	R	DBpedia YAGO Wikidata	371 331 371	29.70 13.99 1.05	31 16 1	46 21 3	5 0 0	11.20 9.26 0.29	10 9 0	25 23 3	2 0 0	13.48 8.76 20.20	13 9 18	27 18 61	3 0 9	5.62 16.96 5.32	5 2 5	25 138 31	1 1 1	16.50 77.08 13.92	15 70 12	70 278 54	0 0 4	36.65 3,236.24 59.33	23 60 21	509 28,418 755	0 0 3
AAUP	JSE	R/C (c=3)	DBpedia YAGO Wikidata	960 889 959	24.40 10.49 2.13	28 11 2	41 17 5	0 0 0	9.38 3.31 0.88	9 3 1	20 11 2	0 0 0	12.68 11.37 30.71	15 12 29	28 15 83	0 0 0	8.20 13.61 8.51	7 3 7	36 138 44	0 1 0	11.74 85.83 22.38	11 68 21	66 446 97	0 0 0	62.18 2,455.27 296.92	23 110 20	2,488 28,418 31,777	0 1 0
Auto 93	JSE	R	DBpedia YAGO Wikidata	93 80 93	28.76 13.80 1.00	31 16 1	43 19 2	5 0 0	11.13 9.09 0.12	10 10 0	25 18 1	3 0 0	12.69 8.37 17.31	12 10 17	22 11 26	8 0 9	4.92 21.09 3.56	5 2 3	7 138 8	2 2 1	14.35 59.33 11.23	11 59 11	64 129 25	4 0 4	22.60 4,025.90 19.91	18 46 19	64 28,418 57	2 4 3
Zoo	UCI ML	C (c=3)	DBpedia YAGO Wikidata	101 8 101	8.61 0.74 1.00	11 0 1	26 13 2	0 0 0	4.67 0.15 0.67	3 0 1	34 6 2	0 0 0	8.22 0.63 29.69	9 1 35	15 8 57	3 0 3	3.54 127.23 8.28	3 138 7	8 138 27	1 2 0	13.26 5.39 18.20	11 1 21	87 156 45	1 0 1	146.28 26,173.23 125.82	24 28,418 92	3,686 28,418 785	2 3 0
Forbes	Forbes	R/C (c=2)	DBpedia YAGO Wikidata	1,585 1,003 1,189	14.77 7.28 0.82	19 10 1	62 33 4	0 0 0	4.87 2.35 0.22	4 2 0	52 42 3	0 0 0	10.15 7.57 16.59	11 11 16	27 21 137	0 0 0	2.76 52.07 5.00	2 2 5	27 138 52	0 1 0	10.44 34.42 12.69	10 27 10	136 510 207	0 0 0	14.30 10,531.37 30.14	4 107 8	1,925 28,418 2,881	0 1 0
Cities	Mercer	R/C (c=3)	DBpedia YAGO Wikidata	212 187 212	31.28 16.66 2.11	35 19 2	53 30 9	0 0 1	6.98 4.46 3.40	7 4 4	26 15 6	0 0 0	18.08 13.75 69.08	19 15 67	38 32 153	0 0 6	25.66 23.56 39.99	25 9 37	68 138 108	0 2 1	16.26 222.54 105.29	13 214 89	131 681 390	0 0 2	1,474.57 8,087.34 5,298.23	678 3,555 1,599	19,810 72,320 99,865	0 5 1
FB Books	Facebook	R/C (c=2)	DBpedia YAGO Wikidata	1,600 1,334 1,578	19.08 8.37 1.00	20 10 1	42 24 3	0 0 0	5.15 2.03 0.01	5 2 0	23 15 1	0 0 0	11.15 8.41 21.19	11 10 22	20 13 55	0 0 0	1.64 25.32 3.15	2 3 3	7 138 17	0 1 0	7.04 24.37 16.41	7 22 16	60 149 69	0 0 0	2.80 4,735.50 7.47	2 8 4	42 28,418 165	0 1 0
FB Movies	Facebook	R/C (c=2)	DBpedia YAGO Wikidata	1,600 1,339 1,585	24.90 12.08 1.01	27 14 1	55 32 4	0 0 0	12.50 6.51 0.04	11 6 0	60 27 1	0 0 0	12.43 8.39 48.75	13 10 48	21 17 107	0 0 0	1.46 26.89 2.22	1 6 1	12 138 22	0 1 0	11.65 55.01 56.37	12 47 53	51 280 372	0 0 0	4.96 4,682.42 20.75	2 43 12	110 28,418 230	0100
Metacritic Albums	Metacritic	R/C (c=2	DBpedia YAGO Wikidata	1,600 1,444 1,576	17.92 7.22 0.99	19 8 1	36 19 2	0 0 0	4.27 3.22 0.00	4 3 0	26 20 1	0	10.85 8.05 17.64	12 9 18	17 10 45	2 0 0	2.63 16.02 4.00	3 3 5	7 138 9	0 1 1	8.92 40.27 11.73	9 32 12	63 361 49	0 0 0	5.28 2,749.90 8.77	3 10 7	50 28,418 54	0

Table 1. Datasets statistics

									11	aDIG	: 4.	D	ataset	s si	aus	sucs												
	Data	set				typ	es		catego	ries				rel o	out			rel i	n		rel-vals	out				rel-va	s in	
Name	Source	Task	LOD	#links	avg	med	max	min	avg	med	max	min	avg	med	max	min av	/g	med	max	min	avg	med	max	min	avg	med	max	min
Metacritic		R/C	DBpedia	2,000	24.38	27	45	0	11.87	11	42	(12.54	14	19	3 1.3	35	1	7	0	11.42	12	30	0	3.56	2	31	0
Movies	Metacritic	(c=2)	YAGO	1,588	11.79	14	19	0	6.43	6	28	(8.34	10	11	0 28	3.22	6	138	1	48.84	43	216	0	4,960.84	37	28,418	1
)	Wikidata	1,981	0.98	1	1	0	0.03	0	1	(47.86	49	99	01.9	98	1	13	0	52.70	53	237	0	15.77	11	117	0
HIV		R/C	DBpedia	114	35.69	37	52	0	12.61	13	23	(23.59	24	28	3 34	1.26	31	89	6	27.75	25	162	10	4,828.36	1,065	70,426	24
Deaths	WHO	(c=2)	YAGO	108	13.90	15	24	0	9.28	9	18	(28.41	31	35	015	5.18	9	138	5	302.34	244	1,267	0	12,464.42	4,879	112,032	550
Country)	Wikidata	114	4.12	4	8	1	4.83	5	6	(120.87	119	173	7 55	5.68	51	148	2	229.46	210	595	2	45,671.15	4,971	669,273	66
Trafic	WHO	R/C	DBpedia	146	36.40	38	53	0	13.12	13	23	(23.40	24	28	1 37	7.87	36	94	8	27.44	24	162	0	7,528.18	1,587	218,957	77
Accidents		(c=2)	YAGO	139	14.29	15	27	0	9.62	10	16		28.44	51	35	014	1.61	9	138	5	345.03	290	2,104	0	17,882.47	6,126	423,559	693
Country)	W1K1data	146	4.42	4	10	1	4.94	5	6		124.31	121	191	/ 61	.68	55	148	2	242.38	213	/13	2	85,5/5.10	/,369	1,557,157	66
Energy	WorldBan	R/C	DBpedia	162	36.07	38	53	0	13.12	13	23	(23.46	24	28	1 36	5.64	33	94	8	26.72	23	162	0	6,876.80	1,440	218,957	//
Savingss	k	(c=2)	YAGO	152	14.09	15	27	0	9.52	10	16	(27.82	31	35	016	5.40	9	138	5	329.28	279	2,104	0	16,969.96	5,821	423,559	151
Country)	Wikidata	162	4.41	4	10	1	4.92	5	6		123.36	119	191	7 60	0.02	55	148	2	238.69	210	/13	2	//,485.01	5,810	1,557,157	66
Inflation	WorldBan	R/C	DBpedia	160	36.00	38	53	0	13.11	13	23	(23.46	24	28	136	5.74	33	94	8	26.85	24	162	0	6,947.59	1,440	218,957	77
Country	k	(c=2)	YAGO	150	14.09	15	27	0	9.44	10	16	(27.80	31	35	016	0.48	9	138	5	331.16	279	2,104	0	17,114.88	5,821	423,559	693
			Wikidata	160	4.39	4	10	1	4.88	5	6	(123.23	119	191	7 60).12	55	148	2	237.94	210	/13	2	/8,453.16	5,810	1,557,157	66
Scientific		R/C	DBpedia	160	36.00	38	53	0	13.11	13	23	(23.46	24	28	1 36	5.74	33	94	8	26.85	24	162	0	6,947.59	1,440	218,957	77
Journals	WorldBan	(c=2)	YAGO	150	14.09	15	27	0	9.44	10	16	(27.80	31	35	016	5.48	9	138	5	331.16	279	2,104	0	17,114.88	5,821	423,559	693
Country	k)	Wikidata	160	4.39	4	10	1	4.88	5	6	(123.23	119	191	7 60	0.12	55	148	2	237.94	210	713	2	78,453.16	5,810	1,557,157	66

. . . .

ISSN: 2278-4632 Vol-10 Issue-09 No.03 September 2020

Unemployment French Region	SemStats	R/C (c=2)	DBpedia YAGO Wikidata	26 26 26	16.38 8.92 1.35	21 8 1	32 14 3	0 3.7 8 2.7 1 2.5	73 77 58	3 2 3	15 8 4	0 7.81 1 12.42 1 86.23	9 12 84	10 14 119	3 12 74	14.19 3.73 34.00	13 4 33	24 6 51	7 7.19 3 81.12 21 83.12	7 60 79	19 299 157	1 28 58	975.88 1,793.19 332.69	969 1,424 193	2,292 4,527 1,464	37 88 137
Endangered Species	a-z-animals	R/C (c=2)	DBpedia YAGO Wikidata	301 65 301	11.84 2.48 1.05	12 0 1	33 16 4	0 6.3 0 0.7 0 0.4	32 76 14	5 0 0	34 12 6	0 10.77 0 1.78 0 34.32	11 0 37	25 9 137	0 0 3	2.96 108.62 6.94	3 138 6	55 138 78	1 12.65 1 9.53 0 22.04	11 0 22	87 136 400	0 0 1	566.25 22,286.36 21,909.94	15 28,418 70	114,742 28,418 6,460,930	1 1 0
Drug-Food Interaction	FinkiLO D	C (c=2)	DBpedia YAGO Wikidata DrugBan k FinkiLO D	1,989 588 1,908 2,000 2,000	8.83 4.46 1.96 2.00 1.00	4 0 2 2 1	38 31 3 2 1	0 5.4 0 0.6 0 0.0 2 \ 1	46 58 01	5 0 0	18 6 1	0 12.65 0 2.15 0 45.92 61.68 3.00	14 0 47 64 3	15 8 79 71 3	0 0 41 3	1.40 99.69 2.78 1.70 0.00	1 138 2 2 0	5 138 17 2 0	0 3.63 1 7.28 1 34.99 0 41.96 0 1.00	3 0 27 41 1	1 6 15 13	12 0 51 0 59 0 32 14 1 1	34.71 20,427.08 32.25 62.49 0.00	24 28,418 26 50 0	158 28,418 487 211 0	0 1 4 0 0

		ty pes				rel out				rel	in		re	l-va	ls		rel- ls in					
		-	-												out		va					
Name	Task	#links	avg	med	max	min	avg	med	max	min	avg	med	max	min	avg	med	max	min	avg	med	max	min
AIFB	C (c=4)	176	1.4	1	2	1	7.1	7	9	5	2.0	2	5	0	18.2	7	219	2	19.8	9	246	0
AM	С	1,000	1.0	1	1	1	19.8	20	29	- 9	0.6	1	3	0	21.9	20	283	7	3.2	1	273	0
	(c=11)																					
MUTA	C (c=2)	340	1.0	1	1	1	9.8	10	14	5	/	/	/	/	65.8	56	465	4	1	1	/	/
G																						1
BGS	C(c=2)	146	1.0	1	1	1	29.7	31	36	21	1.4	2	4	0	25.2	24	54	15	2.7	2	12	0

4 Conclusion and Outlook

In this paper, we have introduced a collection of 22 benchmark datasets for machine learning on the Semantic Web. So far, we have concentrated on classi- fication and regression tasks. There are methods to derive clustering and outlier detection benchmarks from classification and regression datasets [4,5], so that extending the dataset collection for such unsupervised tasks is possible as well. Furthermore, as many datasets on the Semantic Web use extensive hierarchies in the form of ontologies, building benchmark datasets for tasks like *hierarchical multi-label classification* [15] would also be an interesting extension.

Acknowledgments. The work presented in this paper has been partly funded by the German Research Foundation (DFG) under grant number PA 2373/1-1 (Mine@LOD), and the Dutch national program COMMIT.

References

1. Bloehdorn, S., Sure, Y.: Kernel methods for mining instance data in ontologies. In: Aberer, K., et al. (eds.) ASWC/ISWC -2007. LNCS, vol. 4825, pp. 58–71. Springer, Heidelberg (2007). doi:10.1007/978-3-540-76298-0 5

2. Boer, V., Wielemaker, J., Gent, J., Hildebrand, M., Isaac, A., Ossenbruggen, J., Schreiber, G.: Supporting linked data production for cultural heritage institutes: the Amsterdam museum case study. In: Simperl, E., Cimiano, P., Polleres, A., Corcho, O., Presutti, V. (eds.) ESWC 2012. LNCS, vol. 7295, pp. 733–747. Springer, Heidelberg (2012). doi:10.1007/978-3-642-30284-8 56

3. Demšar, J.: Statistical comparisons of classifiers over multiple data sets. J. Mach. Learn. Res. 7, 1–30 (2006)

4. Emmott, A.F., Das, S., Dietterich, T., Fern, A., Wong, W.K.: Systematic construction of anomaly detection benchmarks from real data. In: Proceedings of the ACMSIGKDD Workshop on Outlier Detection and Description, pp. 16–21. ACM (2013)

5. Farber, I., Günnemann, S., Kriegel, H.P., Kröger, P., Müller, E., Schubert, E., Seidl, T., Zimek, A.: On using class-labels in evaluation of clusterings. In: MultiClust: Workshop on Discovering, Summarizing and Using Multiple Clusterings (2010)

6. Jovanovik, M., Bogojeska, A., Trajanov, D., Kocarev, L.: Inferring cuisine-drug interactions using the linked data approach. Scientific reports 5 (2015)

7. Paulheim, H.: Generating possible interpretations for statistics from linked open data. In: Simperl, E., Cimiano, P., Polleres, A., Corcho, O., Presutti, V. (eds.) ESWC 2012. LNCS, vol. 7295, pp. 560–574. Springer, Heidelberg (2012). doi:10. 1007/978-3-642-30284-8 44

8. Rettinger, A., Lösch, U., Tresp, V., d'Amato, C., Fanizzi, N.: Mining the semantic web. Data Min. Knowl. Disc. **24**(3), 613–662 (2012)

9. Ristoski, P., Bizer, C., Paulheim, H.: Mining the web of linked data with rapid-miner. Web Semant. Sci. Serv. Agents WWW **35**, 142–151 (2015)

10. Ristoski, P., Paulheim, H.: Analyzing statistics with background knowledge fromlinked open data. In: Workshop on Semantic Statistics (2013)

11. Ristoski, P., Paulheim, H.: Semantic web in data mining and knowledge discovery: a comprehensive survey. Web Semant. **36**, 1–22 (2016)

12. Ristoski, P., Paulheim, H., Svátek, V., Zeman, V.: The linked data mining challenge 2015. In: KNOW@ LOD (2015)

13. Ristoski, P., Paulheim, H., Svátek, V., Zeman, V.: The linked data mining challenge 2016. In: KNOW@LOD (2016)

14. Schmachtenberg, M., Bizer, C., Paulheim, H.: Adoption of the linked data best practices in different topical domains. In: Mika, P., et al. (eds.) ISWC 2014. LNCS, vol. 8796, pp. 245–260. Springer, Heidelberg (2014). doi:10.1007/978-3-319-11964-9 16

15. Silla Jr., C.N., Freitas, A.A.: A survey of hierarchical classification across different application domains.

ISSN: 2278-4632 Vol-10 Issue-09 No.03 September 2020

Data Min. Knowl. Disc. 22(1), 31–72 (2011)

16. Tresp, V., Bundschus, M., Rettinger, A., Huang, Y.: Towards machine learning on the semantic web. In: da Costa, P.C.G., et al. (eds.) URSW 2005-2007. LNCS, vol. 5327, pp. 282–314. Springer, Heidelberg (2008)