

# A Collection of Benchmark Datasets for Systematic Evaluations of Machine Learning on the Semantic Web

Petar Ristoski<sup>1</sup>, Gerben Klaas Dirk de Vries<sup>2</sup>, and Heiko Paulheim<sup>1</sup>

<sup>1</sup> University of Mannheim, Germany  
Research Group Data and Web Science  
{petar.ristoski,heiko}@informatik.uni-mannheim.de  
<sup>2</sup> WizeNoze, Amsterdam, Netherlands  
g.k.d.devries@outlook.com

## Abstract.

**Resource type:** Datasets

**Permanent URL:** <http://w3id.org/sw4ml-datasets>

In the recent years, several approaches for machine learning on the Semantic Web have been proposed. However, no extensive comparisons between those approaches have been undertaken, in particular due to a lack of publicly available, acknowledged benchmark datasets. In this paper, we present a collection of 22 benchmark datasets of different sizes. Such a collection of datasets can be used to conduct quantitative performance testing and systematic comparisons of approaches.

**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 repository<sup>3</sup>. 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 matching<sup>4</sup>, the

<sup>3</sup> <http://archive.ics.uci.edu/ml/>

<sup>4</sup> <http://oaei.ontologymatching.org/>

*Berlin SPARQL Benchmark*<sup>5</sup> for triple store performance, the Lehigh University Benchmark (LUBM)<sup>6</sup> for reasoning, or the Question Answering over Linked Data (QALD) dataset<sup>7</sup> 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 evaluations 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*<sup>8</sup> or the *Semantic-Web enabled Recommender Systems Challenge*<sup>9</sup>, 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 commonly used in machine learning experiments, (ii) datasets that were generated from official observations, and (iii) datasets generated from existing RDF datasets. Each of the datasets in the first two categories are initially linked to DBpedia<sup>10</sup>. This has two main reasons, (1) DBpedia being a cross-domain knowledge 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 YAGO<sup>11</sup> and Wikidata<sup>12</sup>. Such links could be exploited for systematic evaluation of the relevance of the data of different LOD dataset in 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:

<sup>5</sup> <http://wifo5-03.informatik.uni-mannheim.de/bizer/berlinsparqlbenchmark/>

<sup>6</sup> <http://swat.cse.lehigh.edu/projects/lubm/>

<sup>7</sup> <http://greentackle.techfak.uni-bielefeld.de/~cunger/qald/>

<sup>8</sup> <http://knowalod2016.informatik.uni-mannheim.de/en/linked-data-mining-challenge/>

<sup>9</sup> <http://challenges.2014.eswc-conferences.org/index.php/RecSys>

<sup>10</sup> <http://dbpedia.org>

<sup>11</sup> <http://yago-knowledge.org/>

<sup>12</sup> <http://www.wikidata.org>

- The *Auto MPG* dataset<sup>13</sup> 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 universities<sup>14</sup>. 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* dataset<sup>15</sup> 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 based on 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 2015<sup>16</sup>. The target is to predict the company’s market value as a classification 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 species<sup>17</sup>.
- The *Facebook Movies* dataset contains a list of movies and the number of Facebook likes for each movie<sup>18</sup>. We first selected 10,000 movies from DBpedia, 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.com<sup>19</sup>, which contains an average rating of all time reviews for a list of movies [12]. The initial

<sup>13</sup> <http://archive.ics.uci.edu/ml/datasets/Auto+MPG>

<sup>14</sup> <http://www.amstat.org/publications/jse/jse.data.archive.htm>

<sup>15</sup> <http://www.amstat.org/publications/jse/v1n1/datasets.lock.html>

<sup>16</sup> <http://www.forbes.com/global2000/list/>

<sup>17</sup> <http://a-z-animals.com/>

<sup>18</sup> We use the Facebook Graph API: <https://developers.facebook.com/docs/graph-api>

<sup>19</sup> <http://www.metacritic.com/browse/movies/score/metacore/all>

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.com<sup>20</sup>, 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 Organization<sup>21</sup>. 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 accidents<sup>22</sup>.
- The *Energy Savings Country* dataset contains a list of countries with the total amount of energy savings of primary energy in 2010<sup>23</sup>, which was downloaded from WorldBank<sup>24</sup>. We use the dataset both for regression and classification.
- Similarly, the *Inflation Country* dataset contains a list of countries with the inflation rate for 2011<sup>25</sup>.
- The *Scientific Journals Country* dataset contains a list of countries with a number of scientific and technical journal articles published in 2011<sup>26</sup>.
- 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 classification 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 FinkiLOD<sup>27</sup>. Furthermore, each drug is linked to DrugBank<sup>28</sup>. 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 contains 178 members of a research group, however the smallest group contains only 4

---

<sup>20</sup> <http://www.metacritic.com/browse/albums/score/metacore/all>

<sup>21</sup> <http://apps.who.int/gho/data/view.main.HIV1510>

<sup>22</sup> <http://apps.who.int/gho/data/view.main.51310>

<sup>23</sup> [http://data.worldbank.org/indicator/10.1\\_ENERGY.SAVINGS](http://data.worldbank.org/indicator/10.1_ENERGY.SAVINGS)

<sup>24</sup> <http://www.worldbank.org/>

<sup>25</sup> <http://data.worldbank.org/indicator/NY.GDP.DEFL.KD.ZG>

<sup>26</sup> <http://data.worldbank.org/indicator/IP.JRN.ARTC.SC>

<sup>27</sup> <http://linkeddata.finki.ukim.mk/>

<sup>28</sup> <http://wifo5-03.informatik.uni-mannheim.de/drugbank/>

- 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 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 random 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 toolkit<sup>29</sup>. It contains information about complex molecules that are potentially carcinogenic, which is given by the `isMutagenic` property.
  - The *BGS* dataset was created by the British Geological Survey and describes geological measurements in Great Britain<sup>30</sup>. 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 online<sup>31</sup>.

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.

## 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 classification 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.

<sup>29</sup> <http://dl-learner.org>

<sup>30</sup> <http://data.bgs.ac.uk/>

<sup>31</sup> <http://w3id.org/sw4ml-datasets>

Table 1: Datasets statistics

Dataset		types										categories										rel out										rel in										rel-vals out										rel-vals in									
Name	Source	Task	I/OD	#links	avg	med	max	min	min	avg	med	max	min	min	avg	med	max	min	min	avg	med	max	min	min	avg	med	max	min	min	avg	med	max	min	min	avg	med	max	min																							
Auto MPG	UCI ML	R	DBpedia YAGO Wikidata	371	29.70	31	46	5	11.20	10	25	2	13.48	13	27	3	5.62	5	25	1	16.50	15	70	0	36.65	23	509	0	3.236	24	60	28.418	0	59.33	21	755	3																								
AAUP	JSE	R/C (c=3)	DBpedia YAGO Wikidata	960	24.40	28	41	0	9.38	9	20	0	12.68	15	28	0	8.20	7	36	0	11.74	11	66	0	62.18	23	2,488	0	2,455.27	110	28,418	1	296.92	20	31,777	0																									
Auto 93	JSE	R	DBpedia YAGO Wikidata	93	28.76	31	43	5	11.13	10	25	3	12.69	12	22	8	4.92	5	7	2	14.35	11	64	4	22.60	18	64	2	4.025	90	46	28,418	4	19.91	19	57	3																								
Zoo	UCI ML	C (c=3)	DBpedia YAGO Wikidata	101	8.61	11	26	0	4.67	3	34	0	8.22	9	15	3	3.54	3	8	1	13.26	11	87	1	146.28	24	3,686	2	26,173.23	28	418	28,418	3	125.82	92	785	0																								
Forbes	Forbes	R/C (c=2)	DBpedia YAGO Wikidata	1,585	14.77	19	62	0	4.87	4	52	0	10.15	11	27	0	2.76	2	27	0	10.44	10	136	0	14.30	4	1,925	0	10,531.37	107	28,418	1	30.14	8	2,881	0																									
Cities	Mercer	R/C (c=3)	DBpedia YAGO Wikidata	212	31.28	35	53	0	6.98	7	26	0	18.08	19	38	0	25.66	25	68	0	16.26	13	131	0	1,474.57	678	19,810	0	8,087.34	3,555	72,320	5	5,298.23	1,599	99,865	1																									
FB Books	Facebook	R/C (c=2)	DBpedia YAGO Wikidata	1,600	19.08	20	42	0	5.15	5	23	0	11.15	11	20	0	1.64	2	7	0	7.04	7	60	0	2.80	2	42	0	4,735.50	8	28,418	1	7.47	4	165	0																									
FB Movies	Facebook	R/C (c=2)	DBpedia YAGO Wikidata	1,600	24.90	27	55	0	12.50	11	60	0	12.43	13	21	0	1.46	1	12	0	11.65	12	51	0	4.96	2	110	0	4,682.42	43	28,418	1	20.75	12	230	0																									
Metacritic Albums	Metacritic	R/C (c=2)	DBpedia YAGO Wikidata	1,600	17.92	19	36	0	4.27	4	26	0	10.85	12	17	2	2.63	3	7	0	8.92	9	63	0	5.28	3	50	0	2,749.90	10	28,418	1	0.99	7	54	1																									



Table 3: Datasets statistics

Dataset			types				rel out				rel in				rel-vals out				rel-vals in			
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	C (c=11)	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
MUTAG	C (c=2)	340	1.0	1	1	1	9.8	10	14	5	\	\	\	\	65.8	56	465	4	\	\	\	\
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

**Acknowledgements** 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

- Bloehdorn, S., Sure, Y.: Kernel Methods for Mining Instance Data in Ontologies. *The Semantic Web* pp. 58–71 (2007)
- de Boer, V., Wielemaker, J., van Gent, J., Hildebrand, M., Isaac, A., van Ossenburggen, J., Schreiber, G.: Supporting linked data production for cultural heritage institutes: The amsterdam museum case study. In: *The Semantic Web: Research and Applications*, pp. 733–747. Springer (2012)
- Demšar, J.: Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research* 7, 1–30 (2006)
- 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 ACM SIGKDD Workshop on Outlier Detection and Description*. pp. 16–21. ACM (2013)
- Färber, I., Günemann, 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)
- Jovanovic, M., Bogojeska, A., Trajanov, D., Kocarev, L.: Inferring cuisine-drug interactions using the linked data approach. *Scientific reports* 5 (2015)
- Paulheim, H.: Generating possible interpretations for statistics from linked open data. In: *9th Extended Semantic Web Conference (ESWC)* (2012)
- Rettinger, A., Lösch, U., Tresp, V., d’Amato, C., Fanizzi, N.: Mining the semantic web. *Data Mining and Knowledge Discovery* pp. 613–662 (2012)
- Ristoski, P., Bizer, C., Paulheim, H.: Mining the web of linked data with rapid-miner. *Web Semantics: Science, Services and Agents on the WWW* (2015)
- Ristoski, P., Paulheim, H.: Analyzing statistics with background knowledge from linked open data. In: *Workshop on Semantic Statistics* (2013)
- Ristoski, P., Paulheim, H.: Semantic web in data mining and knowledge discovery: A comprehensive survey. *Web Semantics* 36, 1–22 (2016)
- Ristoski, P., Paulheim, H., Svátek, V., Zeman, V.: The linked data mining challenge 2015. In: *KNOW@ LOD* (2015)
- Ristoski, P., Paulheim, H., Svátek, V., Zeman, V.: The linked data mining challenge 2016. In: *KNOW@LOD* (2016)
- Schmachtenberg, M., Bizer, C., Paulheim, H.: Adoption of the linked data best practices in different topical domains. In: *The Semantic Web–ISWC* (2014)
- Silla, Jr., C.N., Freitas, A.A.: A survey of hierarchical classification across different application domains. *Data Min. Knowl. Discov.* pp. 31–72 (2011)
- Tresp, V., Bundschuh, M., Rettinger, A., Huang, Y.: Towards machine learning on the semantic web. In: *Uncertainty Reasoning for the Semantic Web I* (2008)
- de Vries, G.K.D.: A fast approximation of the Weisfeiler-Lehman graph kernel for RDF data. In: *ECML/PKDD* (1). pp. 606–621 (2013)