

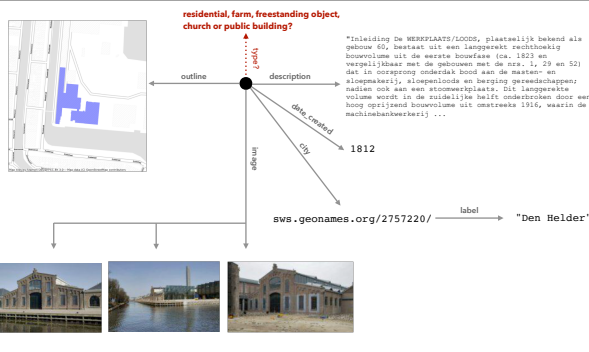
# kgbench

a collection of knowledge graph datasets for evaluating relational and multimodal machine learning

Peter Bloem Xander Wilcke Lucas van Berkel Victor de Boer



Knowledge representation & reasoning  
User centric data science



Imagine if you were told about a building, and shown the outline of its floorplan. Could you guess what its function was?

You might guess that it's not a church (or at least not an old one), and that it's not a single residential house. If I gave you more information, perhaps a description in natural language, you might be able to discern more. I could also show you a picture, and perhaps give you some relational data, including when the building was built and where it's located.

By integrating all of these source of information: visual, natural language, geographical, you can **learn** a lot more than from just a single source.

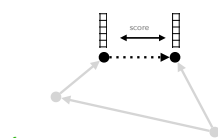
most knowledge is **multimodal** and **relational**

That is the motivation behind this research. Most knowledge that we have, looks like this: various modalities, strung together with relational connections. Knowledge graphs are a very natural representation for this kind of data, as the previous slide shows, so if we want to do machine learning in this kind of setting---learning from all available modalities---and the relational connections between them, we could do worse than to investigate machine learning on knowledge graphs.

## what's available?

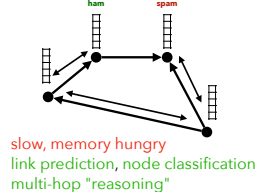
### embedding methods

DistMult, ComplEx, SimplE, etc.



### message passing methods

R-GCN, RGAT, etc.



So let's look at what kind of methods are available so far. Probably the most popular approach are the *embedding models*. These learn a single vector representation for each node in the graph, and crucially, to predict whether a triple is true, look only at the embeddings for the subject and object of that triple. Any information beyond this one-hop neighborhood needs to be included implicitly through the embeddings.

An alternative approach is **message passing**. This extends the embedding vectors representing the nodes by pooling information from deeper in the graph. Such models are slower and more memory hungry than the simple embedding models, but they can use information from deeper in the graph more explicitly.

They can also be used to solve **node classification**

problems, for which simple embedding models are not suitable. We believe that for integrating information in a rich multimodal knowledge, message passing models are the best candidates and that node classification is the simplest testing ground for evaluation.

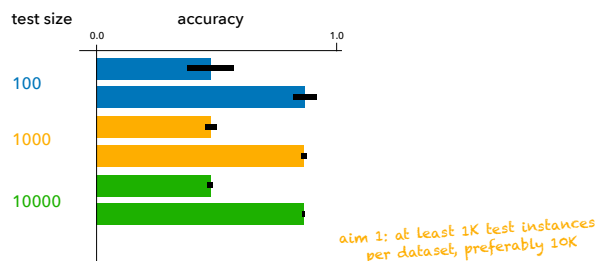
## MP for relational and multimodal data what's stopping us?

Dataset	AIFB	MUTAG	BGS*	AM*
Entities	8 285	23 644	87 688	246 728
Relations	45	23	70	122
Edges	29 043	74 227	230 698	875 946
Labeled	176	340	146	1 000
Classes	4	2	2	11

So, if message passing is such a great framework, what's stopping us? Right, now: a lack of good data. Here are four of the most popular datasets.

The main problem is not the quality of the data, but the number of labeled nodes available as test data. This is the entirety of the labeled data, so if we split off a test set of 20% that can go as low as 30 instances.

### 1. small test sets



The reason that this is a problem is that how well we can compare models depends on the size of the test set. We may compute classification accuracy and be confident that our model is far better than our baseline, but with only 100 instances in our test data, we can see that the error bars are actually very big. If our baseline is 50% accurate and our model is 60% accurate, 100 instances is not enough to then claim that our model has a better performance. Or, if our baseline scores 90% and we score 95%, we may still get overlapping error bars.

At 1000 instances in our data, things start to look a little better, but really for a long-lasting benchmark, we would want something like 10000 instances in our test set (and our validation set as well).

This is our first aim: to create datasets with at least 1000 test instances, but preferably an order of magnitude more.

### 2. large volume per instance

Dataset	AIFB	MUTAG	BGS*	AM*
Entities	8 285	23 644	87 688	246 728
Relations	45	23	70	122
Edges	29 043	74 227	230 698	875 946
Labeled	176	340	146	1 000
Classes	4	2	2	11

aim 2: a medium-sized, information-rich knowledge graph

Now, to achieve this, we can't just take this kind of dataset and scale it up 50 times to get a test set of 10K instances. If we did that, just training a single message passing model would become a great engineering challenge. We do eventually want to train on such data, of course, but for easy model development, we need small datasets, so that we can iterate fast, and allow people with limited hardware resources to perform research.

So, our second aim is to achieve a modest total size of graph, that is still rich in information. To achieve this, we make a trade-off. We remove a lot of information that we judge to be unlikely to be helpful for the classification task in hand. This makes the data less representative of knowledge graphs in the wild, but better suited for fast iteration and model development.

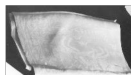
$$\begin{array}{r} \hline \text{AM}^* \\ \hline 246\,728 \\ 122 \\ 875\,946 \\ 1\,000 \\ 11 \\ \hline \end{array}$$

reproduction

\_:N5952ff9311ae4d14bf2f92cabcbc78e0

reproductionIdentifierURL

"..\..\dat\collectie\images\pictura2009\S\_A\_14011\_002\_000.jpg"

$$\begin{array}{r} \hline \text{AM}^* \\ \hline 246\,728 \\ 122 \\ 875\,946 \\ 1\,000 \\ 11 \\ \hline \end{array}$$


## Dutch Monument Graph (dmgfu11, dmgt777k)

integrates information from:

- the Dutch Cultural Heritage Agency  
[www.cultureelerfgoed.nl](http://www.cultureelerfgoed.nl)
- the Dutch Cadastre, Land Registry and Mapping Agency  
[www.kadaster.nl](http://www.kadaster.nl)
- Statistics Netherlands  
[www.cba.nl](http://www.cba.nl)
- Geonames  
[www.geonames.org](http://www.geonames.org)



\*Inleiding de WERKPLAATS/LOCUS, plaatselijk bekend als gebouw 60, bestaat uit een laaggevoelt rechthoekig bouwvolume uit de eerste bouwfase (ca. 1823 en vergelijkbaar met de gebouwen met de nrs. 1, 29 en 52) dat in oorsprong onderdak bood aan de meesters en alogmakers, alogmakers en berging gereedschappen; nadien ook aan een stoomwagplaats. Dit laagprofiel volume wordt in de zuidelijke helft onderbroken door een hoog

classes:  
**residential**  
**farm**  
**freestanding object**  
**church**  
**public building**  
defensive  
archaeological

windmill  
road, waterway  
castle manor  
charity  
catering industry  
part of church  
part of building

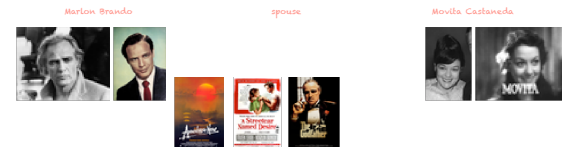


## movie genres (mdgenre)

integrates information from:

- Wikidata (movies and actors)
- IMDb (portraits and posters)

<<http://www.wikidata.org/entity/Q34012>> <<http://www.wikidata.org/prop/direct/P26>> <<http://www.wikidata.org/entity/Q269037>>



classes:  
drama  
action  
literature-based  
fiction  
comedy-drama  
romance  
fantasy  
musical  
documentary  
genre movie  
comedy

## Actor gender (mdgender\*)

- Sanity-check task. Not an official part of the benchmark.
- Potential for study of bias in classification.

<<http://www.wikidata.org/entity/Q34012>> <<http://www.wikidata.org/prop/direct/P21>> <<http://www.wikidata.org/wiki/Q269037>>



classes:  
female  
male  
cisgender  
transgender person  
transgender female  
transgender male  
non-binary  
cisgender male  
gender fluid

The Dutch Monument graph is the data we saw in the first slide. It contains a large amount of multimodal data about monumental buildings in the data.

The dmgt777k version is limited to only the top 5 classes. This reduces the available test and validation data, but it provides a smaller number of classes, and a more uniform class distribution, which is beneficial for training.

For most of these tasks, we have little direct evidence of exactly how much the multimodal data contributes to the problem. As a potential sanity check, we include a gender classification task. The benefit of this task is that we have a reasonable confidence that the gender or sex can be guessed from the portrait image with some accuracy, so this task can be used to make sure that information present in the multimodal literals is propagating through the graph.

However, we emphasize that we do not consider this a suitable benchmark task. The task of classifying sex or gender from appearance or other information is very sensitive, and since benchmarks have a guiding effect for what the community focuses on, we do not believe that this is a suitable general purpose task to train our models on. While wikidata offers sex and gender categories beyond male and female, their frequencies in this data are far too low to be cast in a simple classification framework.

We expect there may be some value to this task in studying gender bias, but for general purpose model evaluation this task should only be used as a sanity check.

## CS publications and authors (db1p)

classes : nr of citations (above or below median)

regression:	classes:
nr of citations	1 citation
	>1 citation

Integrates data from:

- ❖ DBLP (paper names, authors, metadata)
- ❖ Wikidata (extra author metadata)
- ❖ OpenCitations (citation counts)



This dataset is low on multimodal information, but it does provide a well-balanced binary classification objective and a very large amount of test and validation data.

for your convenience

Available in .hdt, .nt.gz. Pre-transformed to CSV:

- ❖ Easy to load into numpy, torch, sk-learn, etc.
- ❖ Python data-loader
- ❖ Baselines and example code

## data loader

```
import kgbench as kg

data = kg.load('amplus') # Load with numpy arrays, and train/validation split

data = kg.load('amplus', torch=True) # Load with pytorch arrays

data = kg.load('amplus', final=True) # Load with numpy arrays and train/test split
```

[illegible]

## baselines and example code

**SO...**

**kgbench.info**

Five knowledge graph node classification tasks with

- ❖ large test sets,
- ❖ manageable size,
- ❖ rich, multimodal information

and

- ❖ a friendly dataloader,
- ❖ baselines/reference implementations.

---

```
pip install git+https://github.com/pbloem/kgbench.git
```

---