

# A Fast Approximation of the Weisfeiler-Lehman Graph Kernel for RDF Data

Referee report

Emilio Cecchini

emilio.cecchini@stud.unifi.it

Lorenzo Palloni

lorenzo.palloni@stud.unifi.it

June 2, 2019

## 1 Summary

The goal of this paper is to introduce a faster version of the Weisfeiler-Lehman graph kernel algorithm when applied to Resource Description Framework (RDF) data.

The *Resource Description Framework* (RDF) is the foundation for knowledge representation on the semantic web. A resource is described by a set of *triples* which are of the form *subject-predicate-object*. The entire collection of triples form a graph where the subjects and the objects are the nodes and the predicates are the edges.

The *Weisfeiler-Lehman test* is an algorithm that is used to compute graph isomorphism. The test proceeds in iterations where the key idea is to augment the node labels by the sorted set of node labels of neighbouring nodes, and compress these augmented labels into new, short labels. These steps are then repeated until the node label sets of the two graphs differ, or the number of iterations reaches the prefixed maximum.

The *Weisfeiler-Lehman kernel* is the state-of-the-art for graph kernels. It computes the number of subtrees shared between two graphs by using the Weisfeiler-Lehman test of graph isomorphism.

This paper introduces an approximation of the Weisfeiler-Lehman kernel, which first extracts a set of subgraphs from the entire RDF graph and then the kernels are computed. For each instance a subgraph up to a certain depth is extracted from the RDF dataset and this subgraph is added to a total graph that the extraction algorithm is building. Thus, vertices and edges are only added if they have not been added to the graph already. For each node and edge, together with their labels, their extraction depth is stored. The relabeling process is the same of the standard Weisfeiler-Lehman test with the extension of the labels on the edges. Finally the kernel is computed by counting the number of common labels at each depth.

## 2 Evaluation

In this paper there is no formal theorem or proof. The author states that this kernel yields an approximation of the standard Weisfeiler-Lehman graph kernel, but he never actually gives any formal proof of the accuracy of that approximation. The comparison with the standard Weisfeiler-Lehman graph kernel can be found only in the experiments section.

In the first experiment, where a classification on the SWRC ontology [2] is performed, the author specifies that the *affiliation* relation and its inverse (the *employs* relation) were removed from the dataset for training purposes. We instead discovered that there are two other relationships that must be removed because they link the instances to their corresponding class, these relationships are *member* and *head*. The fact that these two predicates were not removed from the training dataset led to a higher accuracy than the real one.

This paper proposes a new method on computing graph kernels, but it is limited only to RDF data. This method exploits the fact that usually, in the RDF graphs, the extracted sugraphs share many nodes and edges. This fact limits the number of scenarios in which the method is applicable with good results.

The algorithm described in this paper is an approximation of the Weisfeiler-Lehman graph kernel proposed in [3]. The approximation algorithm is very similar to the standard one described in [3], the only difference is that the label expansion process is also extended to the edges and the concept of *depth* is introduced in order to have bigger graphs without storing duplicated nodes or edges. These two simple modifications seem to lead to a faster version, but there is not much innovation in this new proposed method.

The proposed kernel method is a tool to perform machine learning algorithms on RDF data. There is a small section in the paper where the author introduces the *Resource Description Framework*, but it is never clearly explained what it means to perform a classification on that kind of data.

During the extraction process of the subgraphs of the instances, the algorithm keeps track of the extraction depth to which each node and edge were extracted. In the paper there is confusion about the order of the index of the depth. In the pseudocode of the algorithm the depth is counted backward, that is the root has index equal to the maximum extraction depth while the leaves of the tree have depth equal to zero. While in the explanation of the algorithm the author describes the process with the indexes inverted.

The datasets used in the experiments are still available online. There is a GitHub repository that contains the source code of the experiments but it is quite old and we were not able to compile and to run it.

## References

- [1] Vries Gerben Klaas Dirk, A Fast Approximation of the Weisfeiler-Lehman Graph Kernel for RDF Data, 2013

- [2] Sure, Y., Bloehdorn, S., Haase, P., Hartmann, J., Oberle, D., The swrc ontology - semantic web for research communities. Volume 3803 of LNCS., Covilha, Portugal, Springer (Dezember 2005) 218 – 231
- [3] Shervashidze, N., Schweitzer, P., van Leeuwen, E.J., Mehlhorn, K., Borgwardt, K.M. Weisfeiler-lehman graph kernels, 2011