

An implementation of RDF-Graph Kernels for Spark

Dennis Kubitza ^{*}

Maximilian Radomsky [†]

February 17, 2018

[Project References on Github](#)

Lab Report ¹

Abstract

Many implemented machine learning Algorithms strongly depend on the specific structure of the observed and unobserved Data, which forces users to fit their Observations in a particular predefined setting or reimplement the algorithms to fit their requirements. For dynamic Data models like [Resource Description Frameworks](#), that operate on schema-free Structures, one class of Algorithms is naturally well suited to compromise both approaches: Kernel Based Algorithms. We follow the examinations of [Lösch et al. \(2012\)](#) and implement their proposed Graph-Kernels for the usage in [Apache Spark](#), especially for further usage in the [Semantic Analytics Stack \(SANSa\)](#). Our implementation combines different approaches from Graph Combinatorics, Data-Mining and Big Data Analysis to ensure scalability in storage and computational performance.

^{*}**Email:** denn_kubi@freenet.de, Rudolf-Breitscheid-Str.1 40595 Düsseldorf, Germany

[†]**Email:** maxradomskyy@gmail.com,

¹as Part of the Examination of Modul 4223, Master of Computer Science, University of Bonn

Contents

1	Introduction	3
2	Theory and Approach	3
2.1	Walk Kernel	4
2.2	Path Kernel	4
2.3	Full Subtree Kernel	4
2.4	Partial Subtree Kernel	4

1 Introduction

While each and every Machine Learning Task is defined by its Input Set, its Set of valid models and the expected behaviour of the learning Agent, some algorithms exist that solve problems under such general assumptions that almost any Data-dependent Problem can be reduced to fit their requirements. Kernel-based Machine learning methods don't require any specific structure for the Data, so long as a scalar valued function exists, suitable for summarizing a Observation or Subobservation as a single value. We call such a function a kernel Functions. Especially for schema-free data like RDFs or Labeled Property Graphs such algorithms are highly valuable, as they neither enforce to refactor the data or rewrite existing algorithms and paradigms. As Kernels only need to fulfill very basic properties, Kernel Based Machine learning is very flexible in the definition of a learning task. Especially for Knowledge Graphs like RDFs we can implement different local and global Models in a rather easy way. In the context of growing Databases the usage of Kernels for different Tasks also offers new possibilities as the calculation of Latent Feature Models, which are commonly used to analyse global relations, can get rather costly, especially as Machine learning is needed to train them. In this Lab we will try to give efficient implementations for the computing some basic graph Kernels, described by [Lösch et al. \(2012\)](#), backed by the Spark environment. This report is structured as follows: In following section [Theory and Approach](#) we will define the Kernels in the theory, and give the idea behind the described Graph Kernels [Lösch et al. \(2012\)](#). As we target the implementation of all 4 Kernels, we will state for each of them the problems we expect to occur and our solution proposals. In the section ?? we will then describe the concrete implementation of the final version of the package and the used Structures. The final part thematizing the the implementation will be ?? where we describe our test procedure and the time advancements from our final version to prior/or different solution attempts of us. In the End we will provide a final statement concerning the workflow we set up and the major difficulties we had to cope with during the implementation.

2 Theory and Approach

A definition of Kernels, that is suitably general and applicable to all Machine Learning Algorithms, but still mathematically precise can be found in ([Shawe-Taylor and Cristianini, 2004](#)). In the Context of RDF-Graphs it is possible to reformulate this definition, as [Lösch et al. \(2012\)](#) did, bridging machine learning to feature-representation models for structured Data.

Definition 1 *Let \mathcal{X} be a Subgraph of an RDF-Graph and let $\phi : \mathcal{X} \Rightarrow \mathbb{R}^k$ be a feature representation. A kernel Function is given by*

$$k(x, y) = \langle \phi(x), \phi(y) \rangle_H$$

, where $\langle \cdot, \cdot \rangle_H$ extends \mathbb{R}^k to a Hilbert space.

In the case of RDF-Data [Lösch et al. \(2012\)](#) listed four different Kernels that are in particular suitable, as they may scale to both local and global features and compute Kernels independently from the type of reference or literals given. This is due to the used feature representation, where the existence of certain characteristic subgraphs are identified with Indicator-Variables in the Feature Representation of two Subgraphs of Interest. In the following Paragraphs we will explain them and possible problems / computation approaches, we want to consider in our implementations.

2.1 Walk Kernel

The path kernel corresponds to a weighted sum of the cardinality of walks up to a length l , or more formally:

$$(\kappa_{l,\lambda})(G_1, G_2) = \sum_{i=1}^l \lambda^i |\{p | p \in \text{walks}_i(G_1 \cap G_2)\}|$$

[Lösch et al. \(2012\)](#). The heavy. For further details we refer to ??.

2.2 Path Kernel

The Path Kernel uses the same principle as the Walk Kernel, but counts the number of paths.

$$(\kappa_{l,\lambda})(G_1, G_2) = \sum_{i=1}^l \lambda^i |\{p | p \in \text{paths}_i(G_1 \cap G_2)\}|$$

[Lösch et al. \(2012\)](#). As it is now not possible to use the path independence, as before, we also need to alter the the approach. Most favorable without the need of accessing previous calculations steps, e.g to check if we already visited a single node.

2.3 Full Subtree Kernel

Already [Lösch et al. \(2012\)](#) mentioned that the computation of a Intersection Graph might get costly. They therefore proposed to limit the Calculations of Kernels, not on arbitrary Subgraphs, but only on certain Subgraphs wich can be identified with a certain central entitiy. This enables a replacement of the Intersection Graph with other suitable structures. One of them is the so called Intersection Tree:

Definition 2 XXX

We can obtain the Full Subtree Kernel by XXX to the Intersection Tree of two entities:

Definition 3 XXX

The main target of computational power is now the Intersection Graph. While ? propose an easy Algorithm to generate them, they did not consider the nessecarity of parrallelization. In ?? we provide a parralized Version of the Algorithm together with interesting parts of its computation. The final enumeration of Full Subtrees can easily be paralized, as each connected subgraph of a Tree is neccesarily again a Tree.

2.4 Partial Subtree Kernel

Like the Full Subtree Kernel, the Partial Subtree Kernel is defined by XXX on the Intersection Tree of two Entities:

Definition 4 XXX

Psychic Costs and Human Capital Investment I also shed new light o

References

- Lösch, U., Bloehdorn, S., and Rettinger, A. (2012). Graph kernels for rdf data. In Simperl, E., Cimiano, P., Polleres, A., Corcho, O., and Presutti, V., editors, *The Semantic Web: Research and Applications*, pages 134–148. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Shawe-Taylor, J. and Cristianini, N. (2004). *Kernel methods for pattern analysis*. Cambridge university press.