

POLITECNICO DI TORINO

DEPARTMENT OF CONTROL AND COMPUTER ENGINEERING

Master of Science in Computer Engineering

Master Degree Thesis

Deep Learning on Polito Knowledge Graph

Leveraging Relational GCN for link prediction between nodes of
a newly built publications graph



Supervisors

Prof. Antonio Vetrò
Prof. Juan Carlos De Martin

Candidate

Giovanni GARIFO

ACADEMIC YEAR 2018-2019

To Monia

To my Grandfather

Abstract

Summary here, one page

Acknowledgements

Acknowledgements here, half page

Contents

List of Tables	9
List of Figures	10
1 Introduction	11
1.1 Motivation	11
1.2 Thesis structure	12
1.2.1 Chapter 2	12
1.2.2 Chapter 3	12
1.2.3 Chapter 4	12
2 Background	13
2.1 Semantic Web	13
2.1.1 From a Web of contents to a Web of data	13
2.1.2 The Semantic Web building blocks	14
2.1.3 Knowledge Bases as knowledge repositories	16
2.2 Learning on Graphs	18
2.2.1 Representation learning	18
2.2.2 Deep Learning on graphs	20
3 State of the art	25
4 Approach and Methodology	27
5 Development and Implementation	29
6 Evaluation	31
7 Conclusions	33

List of Tables

List of Figures

2.1	An example of ontology defined using OWL and RDF Schema.	15
2.2	An extract of the Polito Knowledge Graph.	17
2.3	Word vectors allows to perform vector operations, the results obtained reflect the fact that Word2Vec is capable of embed the meaning of such words.	19
2.4	A digital image can be thought of as a graph.	20
2.5	First layer of a GCN updating a node feature vector by embedding the features of adjacent nodes.	22

Chapter 1

Introduction

1.1 Motivation

Graphs are used to empower some of the most complex IT services available today, an example among all being the Google Search Engine ¹. They can be used to represent almost any kind of information, and they are particularly capable of representing the structure of complex systems and describe the relationships between their elements.

Over the last decade, much effort has been made to try to leverage the power of graphs to represent human knowledge and to build search tools capable of querying and understanding the semantic relations within them. RDF graphs are a particular class of graphs that can be used to build knowledge repositories. Given a domain and an ontology, they allow to build a structured representation of the knowledge in such domain.

Modern machine learning techniques can be used to mine latent information from such graphs. One of the main challenges in this field is how to learn meaningful representations of entities and relations that embed the underlying knowledge. Such representations can then be used to evaluate new links inside the graph or to classify unseen nodes. Deep learning techniques have proved to be first class citizens when dealing with representation learning tasks, being able to learn latent representations without any prior knowledge other than the graph structure, so as not to require any feature engineering.

¹https://en.wikipedia.org/wiki/Knowledge_Graph

1.2 Thesis structure

1.2.1 Chapter 2

1.2.2 Chapter 3

1.2.3 Chapter 4

Chapter 2

Background

2.1 Semantic Web

2.1.1 From a Web of contents to a Web of data

The *World Wide Web* has been developed as a tool to easily access documents and to navigate through them by following hyperlinks. This simple description already resembles the structure of a graph: we can think of documents as nodes, and of hyperlinks as edges. The unstoppable growth of the "web graph" led to the emergence of new tools to extricate in such complexity. Search engines have been developed to easily navigate such a giant graph, by scoring search results based on trivial statistics, such as the number of times a document has been linked, as in the case of the PageRank ¹ algorithm developed by Google.

The Web rapidly became one of the most innovative technology ever built, allowing to retrieve information quickly and easily as never before. The next evolutionary step has been to think about a Web not only exploitable by human beings but also by machines. In order to build such a comprehensive system, where information can be not only machine-readable, but machine-understandable, the World Wide Web had to move from a web of content, to a web of data.

The World Wide Web Consortium (W3C) introduced the Semantic Web

¹<https://wikipedia.org/wiki/PageRank>

as an extension to the prior standard of the WWW. Its primary goal has been to define a framework to describe and query semantic information contained in the documents available on the web, so as to allow machines to understand the semantic information contained in web pages. In the vision of Tim Berners-Lee, the father of WWW, this will bring to the transition from a World Wide Web to a Giant Global Graph ², where a web page contains metadata that provides to a machine the needed information to understand the concepts and meanings expressed in it.

2.1.2 The Semantic Web building blocks

The three key components of the Semantic Web standards are:

1. OWL: the Web Ontology Language
2. RDF: the Resource Description Framework
3. SPARQL: The SPARQL Protocol and RDF Query Language

OWL is a language used to define ontologies. In this context, an ontology is defined as a collection of concepts, relations and constraints between these concepts that describes an area of interest or a domain. OWL allows to classify things in terms of their meaning by describing their belonging to classes and subclasses defined by the ontology: if a thing is defined as member of a class, this means that it shares the same semantic meaning as all the other members of such class. The result of such classification is a taxonomy that defines a hierarchy of how things are semantically inter-related in the domain under analysis. The instances of OWL classes are called individuals, and can be related with other individuals or classes by means of properties. Each individual can be characterized with additional information using literals, that represent data values like strings, dates or integers.

RDF is a XML-based framework that defines a standard model for the description, modelling and interchange of resources on the Web.

The first component of the framework is the "RDF Model and Syntax", which defines a data model that describes how the RDF resources should

²https://en.wikipedia.org/wiki/Giant_Global_Graph

be represented. The basic model consist of only three object types: resource, property, and statement. A resource is uniquely identified by an Uniform Resource Identifier (URI). A property can be both a resource attribute or a relation between resources. A statement describes a resource property, and is defined as a triple between a subject (the resource), a predicate (the property) and an object (a literal or another resource).

The second component of the framework is the "RDF Schema" (RDFS), that defines a basic vocabulary for describing RDF resources and the relationships between them. Many vocabularies have been built on top of RDFS, such as the Friend of a Friend (FOAF) vocabulary, for describing social networks, or the one maintained by the Dublin Core Metadata Initiative, that defines common terms used in the definition of metadata for digital resources.

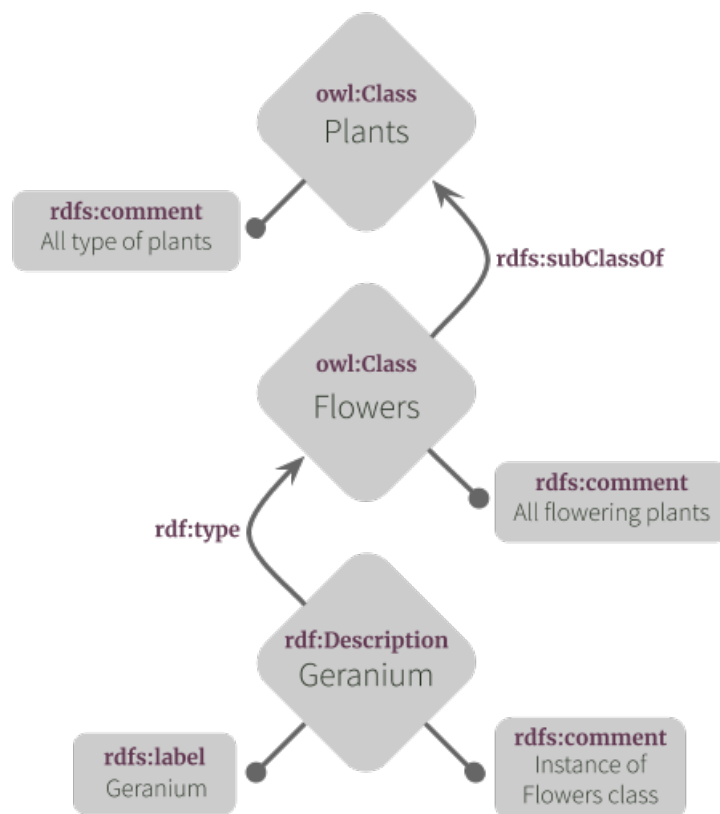


Figure 2.1. An example of ontology defined using OWL and RDF Schema.

SPARQL is a query language for triplestores, a class of Database Management Systems (DBMS) specialized in storing RDF databases. Such DBMS often expose endpoints that can be used to query the database and obtain results. Given the complexity of the data stored, the query language has been designed to be as simple as possible, in example by allowing the use of variables, whose definition is preceded by a question mark.

The syntax of SPARQL is heavily derived from SQL, with some minor adaptations to be more suited for querying graphs data. The following is an example of query which selects all the labels (human-readable description of a resource) of all the entities that matches the given resource type.

```
PREFIX plants:<http://example.org/plants/>

SELECT ?name
WHERE {
    ?subject rdf:type plants:flowers .
    ?subject rdfs:label ?name .
}
```

2.1.3 Knowledge Bases as knowledge repositories

Even if the raise of the Semantic Web has suffered a slowdown in its growth due to the complexity of its vision, many new projects have been born from its enabling technologies. Efforts have been put by profit and non-profit organizations in trying to build complex knowledge repositories starting from the knowledge already present in the Web. An example among all is the DBpedia³ project, which developed a structured knowledge base from the unstructured data available on Wikipedia. Another example is the so called *Knowledge Graph*⁴ made by Google, which is used to enhance its search engine and virtual assistant capabilities, allowing to retrieve punctual information about everything that has been classified in its ontology and described in its knowledge base.

³<https://wiki.dbpedia.org/>

⁴<https://blog.google/products/search/introducing-knowledge-graph-things-not/>

From an implementation perspective, knowledge bases can be created to describe a specific domain by defining an ontology and a vocabulary for such domain using OWL and RDF Schema, and then by describing the concepts of such domain using the RDF Model and Syntax. The RDF document obtained can then be stored in a triplestore and queried using SPARQL. The biggest effort when building knowledge bases is to have a correct understanding and prior knowledge of the domain of interest, to avoid the risk of mischaracterizing and misrepresenting concepts.

If all the requirements and cautions are met, a well formed knowledge base may prove to be a critical resource for an organization. It allows not only to build new services upon it, but also to improve the existing knowledge inside the organization by performing reasoning upon the available knowledge, thus to discover implicit facts that can be derived from existing relationships. Another field of applications is the development of Expert Systems, AI software that emulates the behaviour of a human decision-making process by navigating the knowledge base and taking decisions like in a rule-based system.

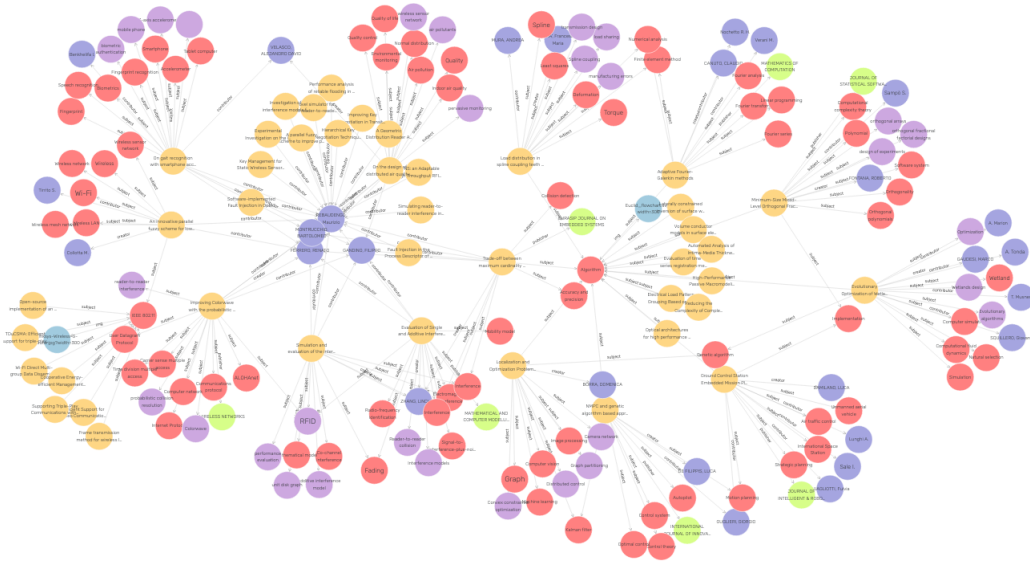


Figure 2.2. An extract of the Polito Knowledge Graph.

Today's knowledge bases are commonly composed by tens of thousands nodes and by hundreds of thousands of edges, such giant data structures pose many challenges. Not only storing and querying giant graphs requires

the adoption of specialized DBMS that are capable of efficiently store and query the RDF input representation, but also doing analysis and gathering statistics from such giant graphs requires the adoption of highly efficient algorithms in order to retrieve the desired output in an acceptable time.

The availability of such a complex and informative data structure leads to the opening of interesting scenarios, especially when thinking about the latent information that can be extracted from it. In fact, a knowledge base is a structured representation of the human knowledge in a specific field, thus its comprehensiveness is restricted by the human understanding.

2.2 Learning on Graphs

2.2.1 Representation learning

Machine learning (ML) algorithms are used to learn models from the available data, with the final goal to obtain a set of parameters that are fine-tuned to identify seen characteristics in the data used for training. The models obtained can be later used to recognize unseen inputs by leveraging the knowledge embedded in such parameters. ML algorithms require the input data to be available in a machine-understandable vector representation. An important task in the ML field is the learning of such representations, task known as representation learning.

Natural Language Processing (NLP) is one of the research branches that in the past years has made a great use of machine learning algorithms both for language recognition and for embedding words meaning into words vectors. One of the most successful algorithms when dealing with representation learning of words is Word2Vec⁵, where the model obtained is trained to learn a vector representation for each word in a vocabulary. In Word2Vec, the concept of meaning of a word is related to the context in which such word is frequently used, so two words are recognized as similar if they're used in similar contexts, thus in the vector space of the learnt representations words that have similar meaning have higher cosine similarity with respect to dissimilar ones. For instance, the cosine similarity between the word vectors of "Man" and "King" is roughly the same as the one between the words "Woman" and "Queen", since such

⁵<https://arxiv.org/abs/1301.3781>

words are used in similar contexes. This has open up new scenarios for language recognition and processing, since it allowed to perform vector operations on such words which brought interesting results, as can be seen in figure 2.3.

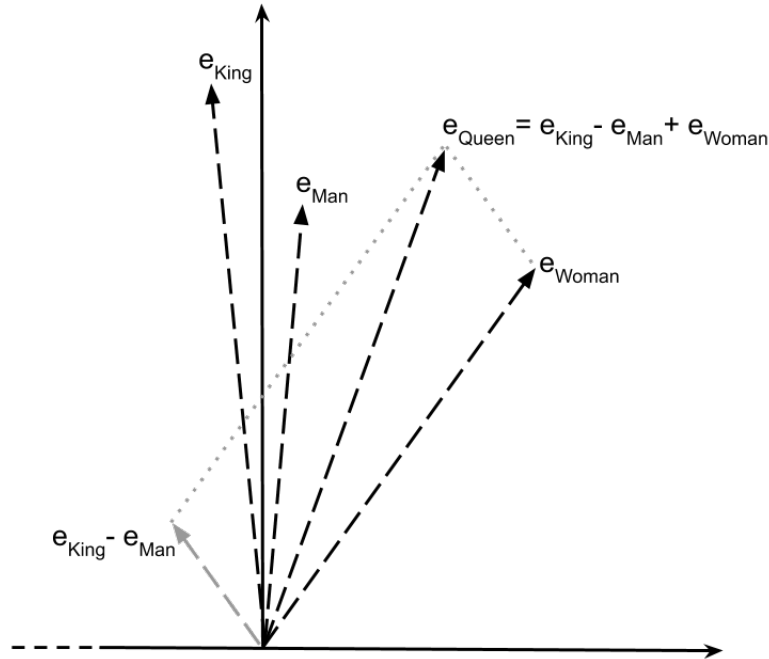


Figure 2.3. Word vectors allows to perform vector operations, the results obtained reflect the fact that Word2Vec is capable of embed the meaning of such words.

The idea that words can be characterized by the context in which they're used can be generalized and used into other fields of research, such as the field of representation learning on graphs.

Graphs are composed by nodes and edges, and are used to describe complex systems, such as social networks or the interactions in a molecular biology system. To apply machine learning algorithms to such data structures, in order to analyze the available data and predict new facts, vector representations of such nodes and edges are needed. Such vector representations are often referred to as "embeddings".

Early approaches required these representations to be learned from feature vectors that were handcrafted, task that required not only a relevant amount of effort, but also a deep understanding of the domain of interest.

This has long been one of the main obstacles when dealing with representation learning tasks, since who has knowledge of the domain and who has to engineer the features were unlikely the same individual.

2.2.2 Deep Learning on graphs

In the latest years a big shift towards deep architectures has been made in machine learning, mainly thanks to the development of highly parallelized architectures that are able to efficiently compute at the hardware level vector and matrix multiplications, operations that are at the basis of any machine learning task. Deep Learning (DL) algorithms are able to extract relevant features from raw data by applying simple mathematical operations, such as convolution, to the input data. An example of one of the most successful applications of DL is in image recognition, where matrix representations of images are convolved with self-trained filters that are able to extract the relevant features needed to recognize patterns present in the input images.

Deep learning techniques have proven to function well also in the field of representation learning for graph data, this should not give rise to surprise, given that as can be seen in figure 2.4, a digital image is composed by pixels which can be thought of as nodes in a graph, where each pixel is connected by an edge to its first neighbours. This suggests that the techniques used when dealing with images can be adapted, with some major changes, to the field of representation learning on graphs, but also in other fields of research, such as learning on manifolds.

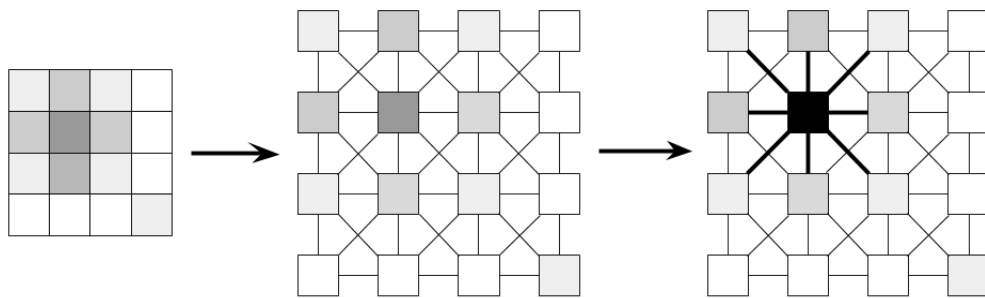


Figure 2.4. A digital image can be thought of as a graph.

The problems when working with graph data is that commonly graphs

are built to describe complex systems, such as the knowledge of a domain or field for knowledge graphs, and thus are composed of a fairly high amount of nodes and edges. The matrices used to store the graph structure can thus explode in dimensionality, so as to become impractical as input data. Moreover, graphs aren't regular structures with given shape and size, like a matrix of pixels for images, but they live in an irregular domain which led to highly irregular structures. The first problem can be solved by randomly sampling the graph at each training epoch, the immediate drawback being that more than one epoch is required to train over all nodes. The second problem can instead be solved by adapting known algorithms to work on highly irregular domains. One of the possible approaches, which has proven to work well, is the one based on convolutions.

Graph Convolutional Networks (GCNs) are a class of semi-supervised deep learning algorithms for graphs which are based on the the same convolution and backpropagation operations as the so famous Convolutional Neural Networks (CNNs) used for feature learning on images. The main difference between CNNs and GCNs is in the architecture of the neural network, and so in how the convolution is performed, instead the backpropagation phase is the same as the one used to update the parameters of CNNs, with the loss function being task-specific. In a CNN the input matrix of each network layer, which is the pixel matrix of the input image for the first, is convolved with a convolutional filter, whose parameters are then updated during the backpropagation phase.

GCNs works differently, but similarly, by convolving at the l -th layer of the network the feature vector of each node with the feature vectors of its l -nearest neighbors, this is done by applying the following transformation:

$$H^{l+1} = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} H^l W^l) \quad (2.1)$$

Where H^l is the output of the previous layer or, for the input layer, the nodes feature matrix where each row is commonly initialized as a one-hot encoded feature vector. \tilde{A} is the adjacency matrix of the graph summed with the identity matrix to add self loops of nodes, \tilde{D} is the node degree matrix of \tilde{A} and is used to normalize it, W^l is the weight matrix of the layer, that is shared among all nodes for each layer, just like the convolutional filter in a CNN, and $\sigma()$ is a non linear activation function, like *ReLU*.

Analyzing the forward rule for a single node make it more clear how the embedding of a node is updated through the convolution:

$$h_i^{(l+1)} = \sigma\left(\sum_{j \in \eta_i} \frac{1}{c_{ij}} h_j^{(l)} W^{(l)}\right) \quad (2.2)$$

Setting aside the normalization costant c_{ij} which is obtained from the multiplication between the adjancency and degree matrices, at each layer the updated feature vector of the node i is obtained by summing over all the nodes in its neighborhood the result of the multiplication between the neighbors feature vectors and the weight matrix of the layer. A consequences of applying this rule to all nodes is that at the l -th layer the feature vectors of nodes that are at a l hop distance from the node i will be embedded in its feature vector, because the feature vectors of such nodes were embedded, at the previous layer, by the immediate neighbors of node i , which are now passing the information of their surrounding to it. So the amount of layers of the network is a parameter that controls how much information for fareset nodes has to be collected from each node embedding.

This means that each node is characterized by it's context, just like it happens in Word2Vec, but in a non-Euclidean domain. So for example, in a social graph where each person is characterized by its friends, interests, places visited and so on, two people will have similar embeddings if they are connected to similar nodes.

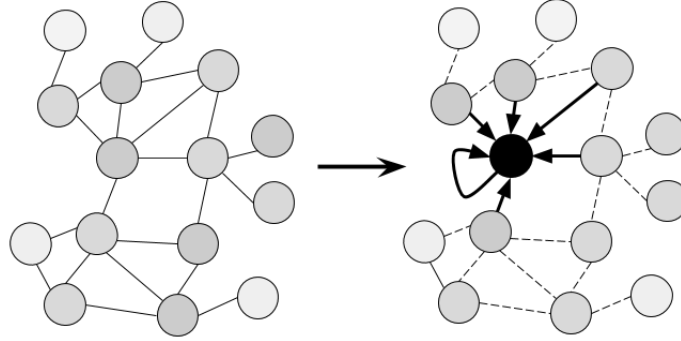


Figure 2.5. First layer of a GCN updating a node feature vector by embedding the features of adjacent nodes.

The node embeddings obtained by applying a GCN or one of its variants can then be used to perform some learning task on the graph, like

classification of unseen nodes or link prediction of non-existent edges, the latter being one of the most interesting task because it empowers most of the recommendation systems available in the industry.

Chapter 3

State of the art

Chapter 4

Approach and Methodology

Chapter 5

Development and Implementation

Chapter 6

Evaluation

Chapter 7

Conclusions

Bibliography

- [1] G. Galilei, *Nuovi studii sugli astri medicei*, Manuzio, Venetia, 1612.
- [2] E. Torricelli, in “La pressione barometrica”, *Strumenti Moderni*, Il Porcellino, Firenze, 1606.
- [3] E. Torricelli e A. Vasari, in “Delle misure”, *Atti Nuovo Cimento*, vol. III, n. 2 (feb. 1607), p. 27–31.
- [4] Duane J.T., *Learning Curve Approach To Reliability Monitoring*, IEEE Transactions on Aerospace, Vol. 2, pp. 563-566, 1964