

Graph Machine Learning

Take graph data to the next level by applying machine learning techniques and algorithms

Claudio Stamile

Aldo Marzullo

Enrico Deusebio



BIRMINGHAM—MUMBAI

Graph Machine Learning

Copyright © 2021 Packt Publishing

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the authors, nor Packt Publishing or its dealers and distributors, will be held liable for any damages caused or alleged to have been caused directly or indirectly by this book.

Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

Group Product Manager: Kunal Parikh
Publishing Product Manager: Devika Battike
Senior Editor: Roshan Kumar
Content Development Editor: Sean Lobo
Technical Editor: Sonam Pandey
Copy Editor: Safis Editing
Project Coordinator: Aparna Ravikumar Nair
Proofreader: Safis Editing
Indexer: Vinayak Purushotham
Production Designer: Joshua Misquitta

First published: May 2021

Production reference: 1270521

Published by Packt Publishing Ltd.
Livery Place
35 Livery Street
Birmingham
B3 2PB, UK.

ISBN 978-1-80020-449-2

www.packt.com

Alla memoria di mio Zio, Franchino Avolio. Alle ruote delle bici troppo sgonfie, all'infanzia che mi ha regalato.

In memory of my uncle, Franchino Avolio. To the wheels of bikes that are too flat, to the childhood he gave me.

– Claudio Stamile

To my family, my roots.

– Aldo Marzullo

To Lili, for always reminding me with your 'learning' process how wonderful the human brain and life are.

– Enrico Deusebio

Contributors

About the authors

Claudio Stamile received an M.Sc. degree in computer science from the University of Calabria (Cosenza, Italy) in September 2013 and, in September 2017, he received his joint Ph.D. from KU Leuven (Leuven, Belgium) and Université Claude Bernard Lyon 1 (Lyon, France). During his career, he has developed a solid background in artificial intelligence, graph theory, and machine learning, with a focus on the biomedical field. He is currently a senior data scientist in CGnal, a consulting firm fully committed to helping its top-tier clients implement data-driven strategies and build AI-powered solutions to promote efficiency and support new business models.

Aldo Marzullo received an M.Sc. degree in computer science from the University of Calabria (Cosenza, Italy) in September 2016. During his studies, he developed a solid background in several areas, including algorithm design, graph theory, and machine learning. In January 2020, he received his joint Ph.D. from the University of Calabria and Université Claude Bernard Lyon 1 (Lyon, France), with a thesis entitled *Deep Learning and Graph Theory for Brain Connectivity Analysis in Multiple Sclerosis*. He is currently a postdoctoral researcher at the University of Calabria and collaborates with several international institutions.

Enrico Deusebio is currently the chief operating officer at CGnal, a consulting firm that helps its top-tier clients implement data-driven strategies and build AI-powered solutions. He has been working with data and large-scale simulations using high-performance facilities and large-scale computing centers for over 10 years, both in an academic and industrial context. He has collaborated and worked with top-tier universities, such as the University of Cambridge, the University of Turin, and the Royal Institute of Technology (KTH) in Stockholm, where he obtained a Ph.D. in 2014. He also holds B.Sc. and M.Sc. degrees in aerospace engineering from Politecnico di Torino.

About the reviewers

Kacper Kubara is a technical co-founder of Artemo and a data engineer at Annual Insight, and is currently pursuing a postgraduate degree in AI at the University of Amsterdam. Despite the focus of his research being graph representation learning, he is also interested in the tools and methods that help to bridge the gap between the AI industry and academia.

Tural Gulmammadov has been leading a group of data scientists and machine learning engineers at Oracle to tackle applied machine learning problems from various industries. He is dedicated to and motivated by the applications of graph theory and discrete mathematics in machine learning over distributed computational environments. He is a cognitive science, statistics, and psychology enthusiast, as well as a chess player, painter, seasonal horse rider, and paddler.

Table of Contents

Preface

Section 1 – Introduction to Graph Machine Learning

1

Getting Started with Graphs

Technical requirements	4	Segregation metrics	30
Introduction to graphs with networkx	5	Centrality metrics	32
Types of graphs	9	Resilience metrics	35
Graph representations	14	Benchmarks and repositories	36
Plotting graphs	18	Examples of simple graphs	36
networkx	18	Generative graph models	38
Gephi	21	Benchmarks	40
Graph properties	26	Dealing with large graphs	47
Integration metrics	27	Summary	48

2

Graph Machine Learning

Technical requirements	52	The generalized graph embedding problem	57
Understanding machine learning on graphs	52	The taxonomy of graph embedding machine learning algorithms	64
Basic principles of machine learning	53		
The benefit of machine learning on graphs	55	The categorization of embedding	

algorithms	65	Summary	68
------------	----	---------	----

Section 2 – Machine Learning on Graphs

3

Unsupervised Graph Learning

Technical requirements	72	Our first autoencoder	95
The unsupervised graph embedding roadmap	72	Denoising autoencoders	100
Shallow embedding methods	74	Graph autoencoders	102
Matrix factorization	74	Graph neural networks	104
Skip-gram	81	Variants of GNNs	106
Autoencoders	92	Spectral graph convolution	107
TensorFlow and Keras – a powerful combination	93	Spatial graph convolution	110
		Graph convolution in practice	111
		Summary	114

4

Supervised Graph Learning

Technical requirements	116	Manifold regularization and semi-supervised embedding	132
The supervised graph embedding roadmap	116	Neural Graph Learning	134
Feature-based methods	117	Planetoid	144
Shallow embedding methods	121	Graph CNNs	145
Label propagation algorithm	121	Graph classification using GCNs	145
Label spreading algorithm	127	Node classification using GraphSAGE	148
Graph regularization methods	131	Summary	150

5

Problems with Machine Learning on Graphs

Technical requirements	152	Embedding-based methods	158
Predicting missing links in a graph	153	Detecting meaningful structures such as communities	
Similarity-based methods	154		

	163	Detecting graph similarities and graph matching	169
Embedding-based community detection	164	Graph embedding-based methods	171
Spectral methods and matrix factorization	165	Graph kernel-based methods	171
Probability models	166	GNN-based methods	171
Cost function minimization	167	Applications	172
		Summary	173

Section 3 – Advanced Applications of Graph Machine Learning

6

Social Network Graphs

Technical requirements	178	Embedding for supervised and unsupervised tasks	188
Overview of the dataset	178	Task preparation	189
Dataset download	179	node2vec-based link prediction	190
Loading the dataset using networkx	180	GraphSAGE-based link prediction	191
Network topology and community detection	181	Hand-crafted features for link prediction	197
Topology overview	182	Summary of results	199
Node centrality	183	Summary	200
Community detection	186		

7

Text Analytics and Natural Language Processing Using Graphs

Technical requirements	202	Knowledge graphs	210
Providing a quick overview of a dataset	203	Bipartite document/entity graphs	212
Understanding the main concepts and tools used in NLP	204	Building a document topic classifier	233
Creating graphs from a corpus of documents	209	Shallow learning methods	234
		Graph neural networks	238
		Summary	249

8

Graph Analysis for Credit Card Transactions

Technical requirements	252	Embedding for supervised and unsupervised fraud detection	270
Overview of the dataset	252	Supervised approach to fraudulent transaction identification	271
Loading the dataset and graph building using networkx	254	Unsupervised approach to fraudulent transaction identification	274
Network topology and community detection	260	Summary	277
Network topology	260		
Community detection	264		

9

Building a Data-Driven Graph-Powered Application

Technical requirements	280	Graph processing engines	285
Overview of Lambda architectures	280	Graph querying layer	288
Lambda architectures for graph-powered applications	283	Selecting between Neo4j and GraphX	293
		Summary	293

10

Novel Trends on Graphs

Technical requirements	296	Graph machine learning and neuroscience	302
Learning about data augmentation for graphs	296	Graph theory and chemistry and biology	304
Sampling strategies	297	Graph machine learning and computer vision	304
Exploring data augmentation techniques	298	Recommendation systems	305
Learning about topological data analysis	299	Summary	305
Topological machine learning	300	Why subscribe?	307
Applying graph theory in new domains	302		

Other Books You May Enjoy

Index

Preface

Graph Machine Learning provides a new set of tools for processing network data and leveraging the power of the relationship between entities that can be used for predictive, modeling, and analytics tasks.

You will start with a brief introduction to graph theory and Graph Machine Learning, learning to understand their potential. As you proceed, you will become well versed with the main machine learning models for graph representation learning: their purpose, how they work, and how they can be implemented in a wide range of supervised and unsupervised learning applications. You'll then build a complete machine learning pipeline, including data processing, model training, and prediction, in order to exploit the full potential of graph data. Moving on, you will cover real-world scenarios, such as extracting data from social networks, text analytics, and natural language processing using graphs and financial transaction systems on graphs. Finally, you will learn how to build and scale out data-driven applications for graph analytics to store, query, and process network information, before progressing to explore the latest trends on graphs.

By the end of this machine learning book, you will have learned the essential concepts of graph theory and all the algorithms and techniques used to build successful machine learning applications.

Who this book is for

This book is for data analysts, graph developers, graph analysts, and graph professionals who want to leverage the information embedded in the connections and relations between data points, unravel hidden structures, and exploit topological information to boost their analysis and models' performance. The book will also be useful for data scientists and machine learning developers who want to build machine learning-driven graph databases. A beginner-level understanding of graph databases and graph data is required. An intermediate-level working knowledge of Python programming and machine learning is also expected to make the most out of this book.

What this book covers

Chapter 1, Getting Started with Graphs, introduces the basic concepts of graph theory using the NetworkX Python library.

Chapter 2, Graph Machine Learning, introduces the main concepts of graph machine learning and graph embedding techniques.

Chapter 3, Unsupervised Graph Learning, covers recent unsupervised graph embedding methods.

Chapter 4, Supervised Graph Learning, covers recent supervised graph embedding methods.

Chapter 5, Problems with Machine Learning on Graphs, introduces the most common machine learning tasks on graphs.

Chapter 6, Social Network Analysis, shows an application of machine learning algorithms on social network data.

Chapter 7, Text Analytics and Natural Language Processing Using Graphs, shows the application of machine learning algorithms to natural language processing tasks.

Chapter 8, Graph Analysis for Credit Card Transactions, shows the application of machine learning algorithms to credit card fraud detection.

Chapter 9, Building a Data-Driven Graph-Powered Application, introduces some technologies and techniques that are useful for dealing with large graphs.

Chapter 10, Novel Trends on Graphs, introduces some novel trends (algorithms and applications) in graph machine learning.

To get the most out of this book

A Jupyter or a Google Colab notebook is sufficient to cover all the examples. For some chapters, Neo4j and Gephi are also required.

Software/Hardware covered in the book	OS Requirements
Python	Windows, macOS X, and Linux (any)
Neo4j	Windows, macOS X, and Linux (any)
Gephi	Windows, macOS X, and Linux (any)
Google Colab or Jupyter Notebook	Windows, macOS X, and Linux (any)

If you are using the digital version of this book, we advise you to type the code yourself or access the code via the GitHub repository (link available in the next section). Doing so will help you avoid any potential errors related to the copying and pasting of code.

Download the example code files

You can download the example code files for this book from GitHub at <https://github.com/PacktPublishing/Graph-Machine-Learning>. In case there's an update to the code, it will be updated on the existing GitHub repository.

We also have other code bundles from our rich catalog of books and videos available at <https://github.com/PacktPublishing/>. Check them out!

Download the color images

We also provide a PDF file that has color images of the screenshots/diagrams used in this book. You can download it here: https://static.packt-cdn.com/downloads/9781800204492_ColorImages.pdf.

Conventions used

There are a number of text conventions used throughout this book.

Code in text: Indicates code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, user input, and Twitter handles. Here is an example: "Mount the downloaded WebStorm-10*.dmg disk image file as another disk in your system."

A block of code is set as follows:

```
html, body, #map {  
    height: 100%;  
    margin: 0;  
    padding: 0  
}
```

When we wish to draw your attention to a particular part of a code block, the relevant lines or items are set in bold:

```
Jupyter==1.0.0
networkx==2.5
matplotlib==3.2.2
node2vec==0.3.3
karateclub==1.0.19
scipy==1.6.2
```

Any command-line input or output is written as follows:

```
$ mkdir css
$ cd css
```

Bold: Indicates a new term, an important word, or words that you see on screen. For example, words in menus or dialog boxes appear in the text like this. Here is an example: "Select **System info** from the **Administration** panel."

Tips or important notes
Appear like this.

Get in touch

Feedback from our readers is always welcome.

General feedback: If you have questions about any aspect of this book, mention the book title in the subject of your message and email us at customercare@packtpub.com.

Errata: Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you have found a mistake in this book, we would be grateful if you would report this to us. Please visit www.packtpub.com/support/errata, selecting your book, clicking on the Errata Submission Form link, and entering the details.

Piracy: If you come across any illegal copies of our works in any form on the internet, we would be grateful if you would provide us with the location address or website name. Please contact us at copyright@packt.com with a link to the material.

If you are interested in becoming an author: If there is a topic that you have expertise in, and you are interested in either writing or contributing to a book, please visit authors.packtpub.com.

Reviews

Please leave a review. Once you have read and used this book, why not leave a review on the site that you purchased it from? Potential readers can then see and use your unbiased opinion to make purchase decisions, we at Packt can understand what you think about our products, and our authors can see your feedback on their book. Thank you!

For more information about Packt, please visit [packt . com](http://packt.com).

