
GRAPH REPRESENTATION LEARNING

Assistant prof. (that time)

↓
WILLIAM L. HAMILTON

McGill University

2020

↖
12 Nobel laureates

~~Hamilton's~~ graduate)

Yoshua Bengio

(ACM Turing Award)

PRE-PUBLICATION DRAFT OF A BOOK TO BE PUBLISHED BY
MORGAN & CLAYPOOL PUBLISHERS.

UNEDITED VERSION RELEASED WITH PERMISSION.
ALL RELEVANT COPYRIGHTS HELD BY THE AUTHOR AND
PUBLISHER EXTEND TO THIS PRE-PUBLICATION DRAFT.

Citation: William L. Hamilton. (2020). Graph Representation Learning.
Synthesis Lectures on Artificial Intelligence and Machine Learning, Vol. 14,
No. 3, Pages 1-159.

Abstract

Graph-structured data is ubiquitous throughout the natural and social sciences, from telecommunication networks to quantum chemistry. Building relational inductive biases into deep learning architectures is crucial if we want systems that can learn, reason, and generalize from this kind of data. Recent years have seen a surge in research on graph representation learning, including techniques for deep graph embeddings, ^① generalizations of convolutional neural networks to graph-structured data, and neural message-passing approaches inspired by belief propagation. ^② These advances in graph representation learning have led to new state-of-the-art results in numerous domains, including chemical synthesis, 3D-vision, recommender systems, question answering, and social network analysis.

The goal of this book is to provide a synthesis and overview of graph representation learning. We begin with a discussion of the goals of graph representation learning, as well as key methodological foundations in graph theory and network analysis. Following this, we introduce and review methods for learning node embeddings, including random-walk based methods and applications to knowledge graphs. We then provide a technical synthesis and introduction to the highly successful graph neural network (GNN) formalism, which has become a dominant and fast-growing paradigm for deep learning with graph data. The book concludes with a synthesis of recent advancements in deep generative models for graphs—a nascent, but quickly growing subset of graph representation learning.

belief (bias)
before learning (graph)

?

Preliminary

① no-free-lunch theorem and machine learning

⇒ inductive bias is important.

(e.g., consider a linear regression.

The assumption that a linear regression is appropriate is an inductive bias.)

② Induction. (what is induction?)

③ Lin & Tegmark '16 Why does deep and cheap learning work so well?

⇒ Assumption ① Low polynomial order ② Locality ③ Symmetry

~~why deep ⇒ A Hierarchical process~~

Contents

Preface	vi
Acknowledgments	vii
1 Introduction	1
1.1 What is a graph?	2
1.1.1 Multi-relational Graphs	2
1.1.2 Feature Information	3
1.2 Machine learning on graphs	4
1.2.1 Node classification	4
1.2.2 Relation prediction	6
1.2.3 Clustering and community detection	7
1.2.4 Graph classification, regression, and clustering	7
2 Background and Traditional Approaches	9
2.1 Graph Statistics and Kernel Methods	9
2.1.1 Node-level statistics and features	10
2.1.2 Graph-level features and graph kernels	13
2.2 Neighborhood Overlap Detection	16
2.2.1 Local overlap measures	16
2.2.2 Global overlap measures	18
2.3 Graph Laplacians and Spectral Methods	21
2.3.1 Graph Laplacians	21
2.3.2 Graph Cuts and Clustering	23
2.3.3 Generalized spectral clustering	26
2.4 Towards Learned Representations	27
I Node Embeddings	28
3 Neighborhood Reconstruction Methods	29
3.1 An Encoder-Decoder Perspective	30
3.1.1 The Encoder	30
3.1.2 The Decoder	31

3.1.3	Optimizing an Encoder-Decoder Model	31
3.1.4	Overview of the Encoder-Decoder Approach	32
3.2	Factorization-based approaches	32
3.3	Random walk embeddings	34
3.3.1	Random walk methods and matrix factorization	36
3.4	Limitations of Shallow Embeddings	36
4	Multi-relational Data and Knowledge Graphs	38
4.1	Reconstructing multi-relational data	39
4.2	Loss functions	40
4.3	Multi-relational decoders	42
4.3.1	Representational abilities	44
II	Graph Neural Networks	46
5	The Graph Neural Network Model	47
5.1	Neural Message Passing	48
5.1.1	Overview of the Message Passing Framework	48
5.1.2	Motivations and Intuitions	50
5.1.3	The Basic GNN	51
5.1.4	Message Passing with Self-loops	52
5.2	Generalized Neighborhood Aggregation	52
5.2.1	Neighborhood Normalization	53
5.2.2	Set Aggregators	54
5.2.3	Neighborhood Attention	56
5.3	Generalized Update Methods	58
5.3.1	Concatenation and Skip-Connections	59
5.3.2	Gated Updates	61
5.3.3	Jumping Knowledge Connections	61
5.4	Edge Features and Multi-relational GNNs	62
5.4.1	Relational Graph Neural Networks	62
5.4.2	Attention and Feature Concatenation	63
5.5	Graph Pooling	64
5.6	Generalized Message Passing	66
6	Graph Neural Networks in Practice	68
6.1	Applications and Loss Functions	68
6.1.1	GNNs for Node Classification	69
6.1.2	GNNs for Graph Classification	70
6.1.3	GNNs for Relation Prediction	70
6.1.4	Pre-training GNNs	71
6.2	Efficiency Concerns and Node Sampling	72
6.2.1	Graph-level Implementations	72
6.2.2	Subsampling and Mini-Batching	72
6.3	Parameter Sharing and Regularization	73

7	Theoretical Motivations	75
7.1	GNNs and Graph Convolutions	75
7.1.1	Convolutions and the Fourier Transform	76
7.1.2	From Time Signals to Graph Signals	77
7.1.3	Spectral Graph Convolutions	81
7.1.4	Convolution-Inspired GNNs	84
7.2	GNNs and Probabilistic Graphical Models	88
7.2.1	Hilbert Space Embeddings of Distributions	88
7.2.2	Graphs as Graphical Models	89
7.2.3	Embedding mean-field inference	90
7.2.4	GNNs and PGMs More Generally	92
7.3	GNNs and Graph Isomorphism	92
7.3.1	Graph Isomorphism	93
7.3.2	Graph Isomorphism and Representational Capacity	93
7.3.3	The Weisfieler-Lehman Algorithm	94
7.3.4	GNNs and the WL Algorithm	95
7.3.5	Beyond the WL Algorithm	97
III	Generative Graph Models	102
8	Traditional Graph Generation Approaches	103
8.1	Overview of Traditional Approaches	103
8.2	Erdős-Rényi Model	104
8.3	Stochastic Block Models	104
8.4	Preferential Attachment	105
8.5	Traditional Applications	107
9	Deep Generative Models	108
9.1	Variational Autoencoder Approaches	109
9.1.1	Node-level Latents	111
9.1.2	Graph-level Latents	113
9.2	Adversarial Approaches	115
9.3	Autoregressive Methods	117
9.3.1	Modeling Edge Dependencies	117
9.3.2	Recurrent Models for Graph Generation	118
9.4	Evaluating Graph Generation	120
9.5	Molecule Generation	121
	Conclusion	123
	Bibliography	125

Preface

→ From 2013

The field of graph representation learning has grown at an incredible—and sometimes unwieldy—pace over the past seven years. I first encountered this area as a graduate student in 2013, during the time when many researchers began investigating deep learning methods for “embedding” graph-structured data. In the years since 2013, the field of graph representation learning has witnessed a truly impressive rise and expansion—from the development of the standard graph neural network paradigm to the nascent work on deep generative models of graph-structured data. The field has transformed from a small subset of researchers working on a relatively niche topic to one of the fastest growing sub-areas of deep learning.

However, as the field has grown, our understanding of the methods and theories underlying graph representation learning has also stretched backwards through time. We can now view the popular “node embedding” methods as well-understood extensions of classic work on dimensionality reduction. We now have an understanding and appreciation for how graph neural networks evolved—somewhat independently—from historically rich lines of work on spectral graph theory, harmonic analysis, variational inference, and the theory of graph isomorphism. This book is my attempt to synthesize and summarize these methodological threads in a practical way. My hope is to introduce the reader to the current practice of the field, while also connecting this practice to broader lines of historical research in machine learning and beyond.

Intended audience This book is intended for a graduate-level researcher in machine learning or an advanced undergraduate student. The discussions of graph-structured data and graph properties are relatively self-contained. However, the book does assume a background in machine learning and a familiarity with modern deep learning methods (e.g., convolutional and recurrent neural networks). Generally, the book assumes a level of machine learning and deep learning knowledge that one would obtain from a textbook such as Goodfellow et al. [2016]’s *Deep Learning Book*.

William L. Hamilton
August 2020

Acknowledgments

Over the past several years, I have had the good fortune to work with many outstanding collaborators on topics related to graph representation learning—many of whom have made seminal contributions to this nascent field. I am deeply indebted to all these collaborators and friends: my colleagues at Stanford, McGill, University of Toronto, and elsewhere; my graduate students at McGill—who taught me more than anyone else the value of pedagogical writing; and my PhD advisors—Dan Jurafsky and Jure Leskovec—who encouraged and seeded this path for my research.

I also owe a great debt of gratitude to the students of my winter 2020 graduate seminar at McGill University. These students were the early “beta testers” of this material, and this book would not exist without their feedback and encouragement. In a similar vein, the exceptionally detailed feedback provided by Petar Veličković, as well as comments by Mariá C. V. Nascimento, Jian Tang, Àlex Ferrer Campo, Seyed Mohammad Sadegh Mahdavi, Yawei Li, Xiaofeng Chen, and Gabriele Corso were invaluable during revisions of the manuscript.

No book is written in a vacuum. This book is the culmination of years of collaborations with many outstanding colleagues—not to mention months of support from my wife and partner, Amy. It is safe to say that this book could not have been written without their support. Though, of course, any errors are mine alone.

William L. Hamilton
August 2020

Chapter 1

Introduction

Graphs are a ubiquitous data structure and a universal language for describing complex systems. In the most general view, a graph is simply a collection of objects (i.e., nodes), along with a set of interactions (i.e., edges) between pairs of these objects. For example, to encode a social network as a graph we might use nodes to represent individuals and use edges to represent that two individuals are friends (Figure 1.1). In the biological domain we could use the nodes in a graph to represent proteins, and use the edges to represent various biological interactions, such as kinetic interactions between proteins.

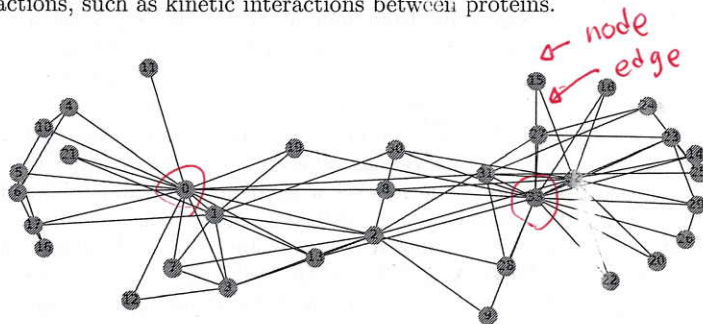


Figure 1.1: The famous *Zachary Karate Club Network* represents the friendship relationships between members of a karate club studied by Wayne W. Zachary from 1970 to 1972. An edge connects two individuals if they socialized outside of the club. During Zachary's study, the club split into two factions—centered around nodes 0 and 33—and Zachary was able to correctly predict which nodes would fall into each faction based on the graph structure [Zachary, 1977].

The power of the graph formalism lies both in its focus on relationships between points (rather than the properties of individual points), as well as in its generality. The same graph formalism can be used to represent social networks, interactions between drugs and proteins, the interactions between atoms

Ask the listeners to list some graphs

in a molecule, or the connections between terminals in a telecommunications network—to name just a few examples.

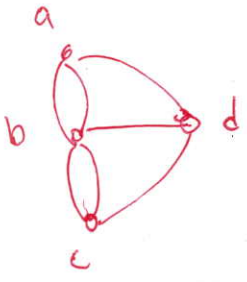
Graphs do more than just provide an elegant theoretical framework, however. They offer a mathematical foundation that we can build upon to analyze, understand, and learn from real-world complex systems. In the last twenty-five years, there has been a dramatic increase in the quantity and quality of graph-structured data that is available to researchers. With the advent of large-scale social networking platforms, massive scientific initiatives to model the interactome, food webs, databases of molecule graph structures, and billions of interconnected web-enabled devices, there is no shortage of meaningful graph data for researchers to analyze. The challenge is unlocking the potential of this data.

1995 ~ windows 95 (the Internet)

!

This book is about how we can use machine learning to tackle this challenge. Of course, machine learning is not the only possible way to analyze graph data.¹ However, given the ever-increasing scale and complexity of the graph datasets that we seek to analyze, it is clear that machine learning will play an important role in advancing our ability to model, analyze, and understand graph data.

Not the only way but may be interesting.



1.1 What is a graph?

Before we discuss machine learning on graphs, it is necessary to give a bit more formal description of what exactly we mean by “graph data”. Formally, a graph $G = (\mathcal{V}, \mathcal{E})$ is defined by a set of nodes \mathcal{V} and a set of edges \mathcal{E} between these nodes. We denote an edge going from node $u \in \mathcal{V}$ to node $v \in \mathcal{V}$ as $(u, v) \in \mathcal{E}$. In many cases we will be concerned only with simple graphs, where there is at most one edge between each pair of nodes, no edges between a node and itself, and where the edges are all undirected, i.e., $(u, v) \in \mathcal{E} \leftrightarrow (v, u) \in \mathcal{E}$.

(no) self-loop

A convenient way to represent graphs is through an adjacency matrix $A \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$. To represent a graph with an adjacency matrix, we order the nodes in the graph so that every node indexes a particular row and column in the adjacency matrix. We can then represent the presence of edges as entries in this matrix: $A[u, v] = 1$ if $(u, v) \in \mathcal{E}$ and $A[u, v] = 0$ otherwise. If the graph contains only undirected edges then A will be a symmetric matrix, but if the graph is directed (i.e., edge direction matters) then A will not necessarily be symmetric. Some graphs can also have weighted edges, where the entries in the adjacency matrix are arbitrary real-values rather than $\{0, 1\}$. For instance, a weighted edge in a protein-protein interaction graph might indicated the strength of the association between two proteins.

Simple graph $\Rightarrow \{0, 1\}$

Write down an adjacency matrix

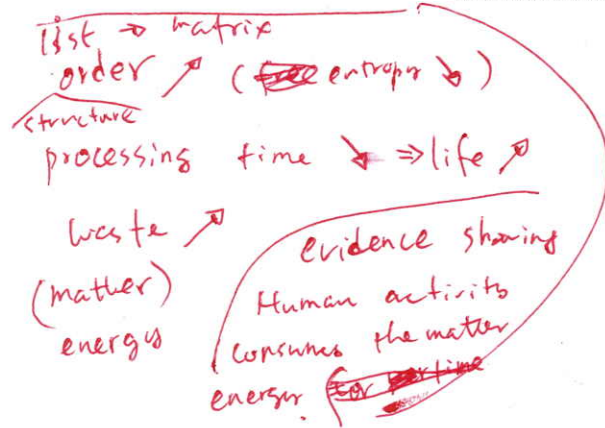
- Inefficient
- but convenient for parallel processing

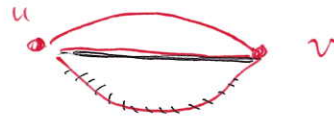
a.f. adjacency list

1.1.1 Multi-relational Graphs

Beyond the distinction between undirected, directed and weighted edges, we will also consider graphs that have different types of edges. For instance, in graphs representing drug-drug interactions, we might want different edges to

¹The field of network analysis independent of machine learning is the subject of entire textbooks and will not be covered in detail here [Newman, 2018].

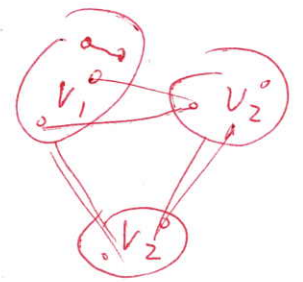




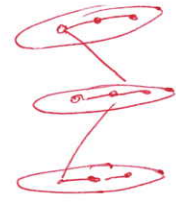
1.1. WHAT IS A GRAPH?

correspond to different side effects that can occur when you take a pair of drugs at the same time. In these cases we can extend the edge notation to include an edge or relation type τ , e.g., $(u, \tau, v) \in \mathcal{E}$, and we can define one adjacency matrix \mathbf{A}_τ per edge type. We call such graphs multi-relational, and the entire graph can be summarized by an adjacency tensor $\mathcal{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{R}| \times |\mathcal{V}|}$, where \mathcal{R} is the set of relations. Two important subsets of multi-relational graphs are often known as heterogeneous and multiplex graphs.

Heterogeneous graphs In heterogeneous graphs, nodes are also imbued with types, meaning that we can partition the set of nodes into disjoint sets $\mathcal{V} = \mathcal{V}_1 \cup \mathcal{V}_2 \cup \dots \cup \mathcal{V}_k$ where $\mathcal{V}_i \cap \mathcal{V}_j = \emptyset, \forall i \neq j$. Edges in heterogeneous graphs generally satisfy constraints according to the node types, most commonly the constraint that certain edges only connect nodes of certain types, i.e., $(u, \tau_i, v) \in \mathcal{E} \rightarrow u \in \mathcal{V}_j, v \in \mathcal{V}_k$. For example, in a heterogeneous biomedical graph, there might be one type of node representing proteins, one type of representing drugs, and one type representing diseases. Edges representing "treatments" would only occur between drug nodes and disease nodes. Similarly, edges representing "polypharmacy side-effects" would only occur between two drug nodes. Multipartite graphs are a well-known special case of heterogeneous graphs, where edges can only connect nodes that have different types, i.e., $(u, \tau_i, v) \in \mathcal{E} \rightarrow u \in \mathcal{V}_j, v \in \mathcal{V}_k \wedge j \neq k$.



\neq k-partite graph \rightarrow inter \rightarrow difference \rightarrow inner \rightarrow



Multiplex graphs In multiplex graphs we assume that the graph can be decomposed in a set of k layers. Every node is assumed to belong to every layer, and each layer corresponds to a unique relation, representing the intra-layer edge type for that layer. We also assume that inter-layer edges types can exist, which connect the same node across layers. Multiplex graphs are best understood via examples. For instance, in a multiplex transportation network, each node might represent a city and each layer might represent a different mode of transportation (e.g., air travel or train travel). Intra-layer edges would then represent cities that are connected by different modes of transportation, while inter-layer edges represent the possibility of switching modes of transportation within a particular city.

1.1.2 Feature Information

Lastly, in many cases we also have attribute or feature information associated with a graph (e.g., a profile picture associated with a user in a social network). Most often these are node-level attributes that we represent using a real-valued matrix $\mathbf{X} \in \mathbb{R}^{|\mathcal{V}| \times m}$, where we assume that the ordering of the nodes is consistent with the ordering in the adjacency matrix. In heterogeneous graphs we generally assume that each different type of node has its own distinct type of attributes. In rare cases we will also consider graphs that have real-valued edge features in addition to discrete edge types, and in some cases we even associate real-valued features with entire graphs.

m can be different

1

2

3

Graph or network? We use the term “graph” in this book, but you will see many other resources use the term “network” to describe the same kind of data. In some places, we will use both terms (e.g., for social or biological networks). So which term is correct? In many ways, this terminological difference is a historical and cultural one: the term “graph” appears to be more prevalent in machine learning community^a, but “network” has historically been popular in the data mining and (unsurprisingly) network science communities. We use both terms in this book, but we also make a distinction between the usage of these terms. We use the term *graph* to describe the abstract data structure that is the focus of this book, but we will also often use the term *network* to describe specific, real-world instantiations of this data structure (e.g., social networks). This terminological distinction is fitting with their current popular usages of these terms. *Network analysis* is generally concerned with the properties of real-world data, whereas *graph theory* is concerned with the theoretical properties of the mathematical graph abstraction.

^aPerhaps in some part due to the terminological clash with “neural networks.”

1.2 Machine learning on graphs

Machine learning is inherently a problem-driven discipline. We seek to build models that can learn from data in order to solve particular tasks, and machine learning models are often categorized according to the type of task they seek to solve: Is it a *supervised* task, where the goal is to predict a target output given an input datapoint? Is it an *unsupervised* task, where the goal is to infer patterns, such as clusters of points, in the data?

Machine learning with graphs is no different, but the usual categories of supervised and unsupervised are not necessarily the most informative or useful when it comes to graphs. In this section we provide a brief overview of the most important and well-studied machine learning tasks on graph data. As we will see, “supervised” problems are popular with graph data, but machine learning problems on graphs often blur the boundaries between the traditional machine learning categories.

1.2.1 Node classification

Suppose we are given a large social network dataset with millions of users, but we know that a significant number of these users are actually bots. Identifying these bots could be important for many reasons: a company might not want to advertise to bots or bots may actually be in violation of the social network’s terms of service. Manually examining every user to determine if they are a bot would be prohibitively expensive, so ideally we would like to have a model that could classify users as a bot (or not) given only a small number of manually



↳ supervised

↳ unsupervised

labeled examples.

This is a classic example of *node classification*, where the goal is to predict the label y_u —which could be a type, category, or attribute—associated with all the nodes $u \in \mathcal{V}$, when we are only given the true labels on a *training set* of nodes $\mathcal{V}_{\text{train}} \subset \mathcal{V}$. Node classification is perhaps the most popular machine learning task on graph data, especially in recent years. Examples of node classification beyond social networks include classifying the function of proteins in the interactome [Hamilton et al., 2017b] and classifying the topic of documents based on hyperlink or citation graphs [Kipf and Welling, 2016a]. Often, we assume that we have label information only for a very small subset of the nodes in a single graph (e.g., classifying bots in a social network from a small set of manually labeled examples). However, there are also instances of node classification that involve many labeled nodes and/or that require generalization across disconnected graphs (e.g., classifying the function of proteins in the interactomes of different species).

At first glance, node classification appears to be a straightforward variation of standard supervised classification, but there are in fact important differences. The most important difference is that the nodes in a graph are not *independent and identically distributed (i.i.d.)*. Usually, when we build supervised machine learning models we assume that each datapoint is statistically *independent* from all the other datapoints; otherwise, we might need to model the dependencies between all our input points. We also assume that the datapoints are *identically distributed*; otherwise, we have no way of guaranteeing that our model will generalize to new datapoints. Node classification completely breaks this i.i.d. assumption. Rather than modeling a set of i.i.d. datapoints, we are instead modeling an interconnected set of nodes.

In fact, the key insight behind many of the most successful node classification approaches is to explicitly leverage the connections between nodes. One particularly popular idea is to exploit *homophily*, which is the tendency for nodes to share attributes with their neighbors in the graph [McPherson et al., 2001]. For example, people tend to form friendships with others who share the same interests or demographics. Based on the notion of homophily we can build machine learning models that try to assign similar labels to neighboring nodes in a graph [Zhou et al., 2004]. Beyond homophily there are also concepts such as *structural equivalence* [Donnat et al., 2018], which is the idea that nodes with similar local neighborhood structures will have similar labels, as well as *heterophily*, which presumes that nodes will be preferentially connected to nodes with different labels.² When we build node classification models we want to exploit these concepts and model the relationships between nodes, rather than simply treating nodes as independent datapoints.

Supervised or semi-supervised? Due to the atypical nature of node classification, researchers often refer to it as *semi-supervised* [Yang et al., 2016]. This terminology is used because when we are training node classi-

²For example, gender is an attribute that exhibits heterophily in many social networks.

applications of
node classification

$$|\mathcal{V}_{\text{train}}| \ll |\mathcal{V}|$$

difference from a
standard supervised
learning (i.i.d. assump.)
graph model?

e.g., social network

math

!!
inductive
bias

②

③

classification models, we usually have access to the full graph, including all the unlabeled (e.g., test) nodes. The only thing we are missing is the labels of test nodes. However, we can still use information about the test nodes (e.g., knowledge of their neighborhood in the graph) to improve our model during training. This is different from the usual supervised setting, in which unlabeled datapoints are completely unobserved during training.

The general term used for models that combine labeled and unlabeled data during training is semi-supervised learning, so it is understandable that this term is often used in reference to node classification tasks. It is important to note, however, that standard formulations of semi-supervised learning still require the i.i.d. assumption, which does not hold for node classification. Machine learning tasks on graphs do not easily fit our standard categories!

difference ②

1.2.2 Relation prediction

Node classification is useful for inferring information about a node based on its relationship with other nodes in the graph. But what about cases where we are missing this relationship information? What if we know only some of protein-protein interactions that are present in a given cell, but we want to make a good guess about the interactions we are missing? Can we use machine learning to infer the edges between nodes in a graph?

This task goes by many names, such as link prediction, graph completion, and relational inference, depending on the specific application domain. We will simply call it *relation prediction* here. Along with node classification, it is one of the more popular machine learning tasks with graph data and has countless real-world applications: recommending content to users in social platforms [Ying et al., 2018a], predicting drug side-effects [Zitnik et al., 2018], or inferring new facts in a relational databases [Bordes et al., 2013]—all of these tasks can be viewed as special cases of relation prediction.

The standard setup for relation prediction is that we are given a set of nodes \mathcal{V} and an incomplete set of edges between these nodes $\mathcal{E}_{\text{train}} \subset \mathcal{E}$. Our goal is to use this partial information to infer the missing edges $\mathcal{E} \setminus \mathcal{E}_{\text{train}}$. The complexity of this task is highly dependent on the type of graph data we are examining. For instance, in simple graphs, such as social networks that only encode “friendship” relations, there are simple heuristics based on how many neighbors two nodes share that can achieve strong performance [Lü and Zhou, 2011]. On the other hand, in more complex multi-relational graph datasets, such as biomedical knowledge graphs that encode hundreds of different biological interactions, relation prediction can require complex reasoning and inference strategies [Nickel et al., 2016]. Like node classification, relation prediction blurs the boundaries of traditional machine learning categories—often being referred to as both supervised and unsupervised—and it requires inductive biases that are specific to the graph domain. In addition, like node classification, there are

many variants of relation prediction, including settings where the predictions are made over a single, fixed graph [Lü and Zhou, 2011], as well as settings where relations must be predicted across multiple disjoint graphs [Teru et al., 2020].

1.2.3 Clustering and community detection

Both node classification and relation prediction require inferring *missing* information about graph data, and in many ways, those two tasks are the graph analogues of supervised learning. Community detection, on the other hand, is the graph analogue of unsupervised clustering.

Suppose we have access to all the citation information in Google Scholar, and we make a *collaboration graph* that connects two researchers if they have co-authored a paper together. If we were to examine this network, would we expect to find a dense “hairball” where everyone is equally likely to collaborate with everyone else? It is more likely that the graph would segregate into different *clusters* of nodes, grouped together by research area, institution, or other demographic factors. In other words, we would expect this network—like many real-world networks—to exhibit a *community* structure, where nodes are much more likely to form edges with nodes that belong to the same community.

This is the general intuition underlying the task of community detection. The challenge of community detection is to infer latent community structures given only the input graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. The many real-world applications of community detection include uncovering functional modules in genetic interaction networks [Agrawal et al., 2018] and uncovering fraudulent groups of users in financial transaction networks [Pandit et al., 2007].

1.2.4 Graph classification, regression, and clustering

The final class of popular machine learning applications on graph data involve classification, regression, or clustering problems over entire graphs. For instance, given a graph representing the structure of a molecule, we might want to build a regression model that could predict that molecule’s toxicity or solubility [Gilmer et al., 2017]. Or, we might want to build a classification model to detect whether a computer program is malicious by analyzing a graph-based representation of its syntax and data flow [Li et al., 2019]. In these *graph classification* or *regression* applications, we seek to learn over graph data, but instead of making predictions over the individual components of a single graph (i.e. the nodes or the edges), we are instead given a dataset of multiple different graphs and our goal is to make independent predictions specific to each graph. In the related task of *graph clustering*, the goal is to learn an unsupervised measure of similarity between pairs of graphs.

Of all the machine learning tasks on graphs, graph regression and classification are perhaps the most straightforward analogues of standard supervised learning. Each graph is an i.i.d. datapoint associated with a label, and the goal is to use a labeled set of training points to learn a mapping from datapoints

← semi-

bias

(詐數)

applications

{ classification
regression

different → lists

(i.e., graphs) to labels. In a similar way graph clustering is the straightforward extension of unsupervised clustering for graph data. The challenge in these graph-level tasks, however, is how to define useful features that take into account the relational structure within each datapoint.

comment: iid assumption

In general, iid is required to get a good result.

(training data and test data are iid from the same distribution)

However, it is usually not satisfied.

⇒ Machine Learning needs inductive bias on the causal relationship of the data. e.g;

因果關係。

Suppose we want to a classifier ~~to~~ to classifier a camel and a cow. If training data (pictures) uses

- camel in the dessert, and
- cow in ~~the field with~~ grassland

We know that ML may learn the difference of dessert and grassland ⇒ ~~bad~~ ~~bad~~ ~~perfor~~ wrong result for a camel in the grassland. ⇒ We need to ML the know

the causal relationship that the result should focus on the camel/cow objects, not the background (correlation).