

1

Deep Learning on Graphs An Introduction

1.1 Introduction

We start the chapter by answering a few questions about the book. First, we discuss why we should pay attention to deep learning on graphs. In particular, why do we represent real-world data as graphs, why do we want to bridge deep learning with graphs, and what are the challenges for deep learning on graphs? Second, we introduce the content that will be covered by this book, specifically, which topics we will discuss and how we organize these topics. Third, we provide guidance about who should read this book; in particular, who our target audience is and how to read this book for readers with different backgrounds and purposes for reading. To help better understand deep learning on graphs, we briefly review the history under the more general context of feature learning on graphs.

1.2 Why Deep Learning on Graphs?

Because data from real-world applications have very diverse forms, from matrix and tensor to sequence and time series, a natural question that arises is why we attempt to represent data as graphs. There are two main motivations. First, graphs provide a universal representation of data, as shown in Figure 1.1. Data from many systems across various areas can be explicitly denoted as graphs such as social networks, transportation networks, protein–protein interaction networks, knowledge graphs, and brain networks. In addition as indicated by Figure 1.1, numerous other types of data can be transformed into the form of graphs. Second, a huge number of real-world problems can be addressed as a small set of computational tasks on graphs. For example, inferring node attributes, detecting anomalous nodes (e.g., spammers or terrorists),

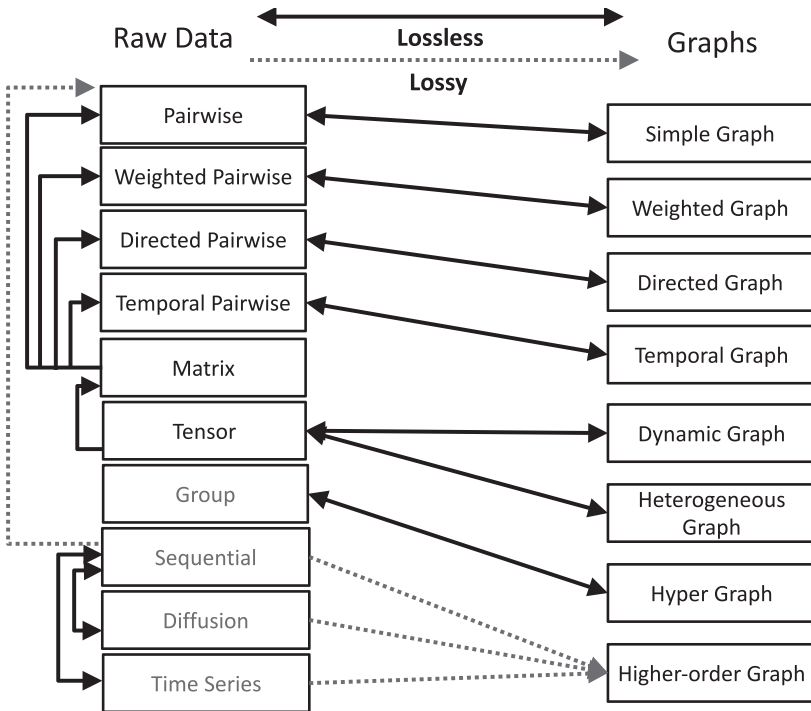


Figure 1.1 Representing real-world data as graphs. The figure is reproduced with permission from J. Xu (2017). Solid lines are utilized to denote lossless representations and dotted lines are used to indicate lossy representations. Note that we replace “network” in the original figure with “graph.”

identifying genes relevant to diseases, and suggesting medications for patients can be summarized as problems of node classification, and recommendations, polypharmacy side effect prediction, drug–target interaction identification, and knowledge graph completion are essentially problems of link prediction.

Nodes on graphs are inherently connected, which suggests that nodes are not independent and identically distributed. Thus, traditional machine learning techniques cannot be directly applied for the computational tasks on graphs. There are two main directions to develop solutions. As shown in Figure 1.2, we will use node classification as an illustrative example to discuss these two directions. One direction is to build a new mechanism specific to graphs. The classification problem designed for graphs is known as collective classification (Sen et al., 2008) as demonstrated in Figure 1.2a. Different from traditional classification, collective classification for a node considers not only

the mapping between its features and its label but also the mapping for its neighborhood. The other direction is to flat a graph by constructing a set of features to denote its nodes where traditional classification techniques can be applied as illustrated in Figure 1.2b. This direction can take advantage of traditional machine learning techniques; thus it has become increasingly popular and dominant. The key to the success of this direction is how to construct a set of features for nodes (or how to construct node representations). Deep learning has been proven to be powerful in representation learning, which has greatly advanced various domains such as computer vision, speech recognition, and natural language processing. Therefore, bridging deep learning with graphs presents unprecedented opportunities. However, deep learning on graphs also faces immense challenges. First, traditional deep learning has been designed for regular data such as images and sequences where the size is fixed, whereas graphs can have varied sizes and nodes in a graph are unordered and can have distinct neighborhoods. Second, the structural information for regular data is simple, whereas that for graphs is complicated, especially given that there are various types of complex graphs, as shown in Figure 1.1, and nodes and edges can associate with rich side information; thus traditional deep learning is not sufficient to capture such rich information. Embracing the unprecedented opportunities as well as the immense challenges, a new research field has been cultivated – deep learning on graphs.

1.3 What Content Is Covered?

The high-level organization of the book is illustrated in Figure 1.3. The book consists of four parts to best accommodate readers with diverse backgrounds and purposes for reading. Part I introduces basic concepts, Part II discusses the most established methods, Part III presents the most representative applications, and Part IV describes advances of methods and applications that are believed to be important and promising for future research. In each chapter, we first motivate the content that will be covered, then present the content with necessary examples or technical details, and finally provide more relevant content as further reading. Next, we briefly elaborate on each chapter.

- Part I: Foundations. These chapters focus on the basics of graphs and deep learning that will lay foundations for deep learning on graphs. In Chapter 2, we introduce the key concepts and properties of graphs, graph Fourier transform, and graph signal processing and formally define various types of complex graphs and computational tasks on graphs. In Chapter 3, we

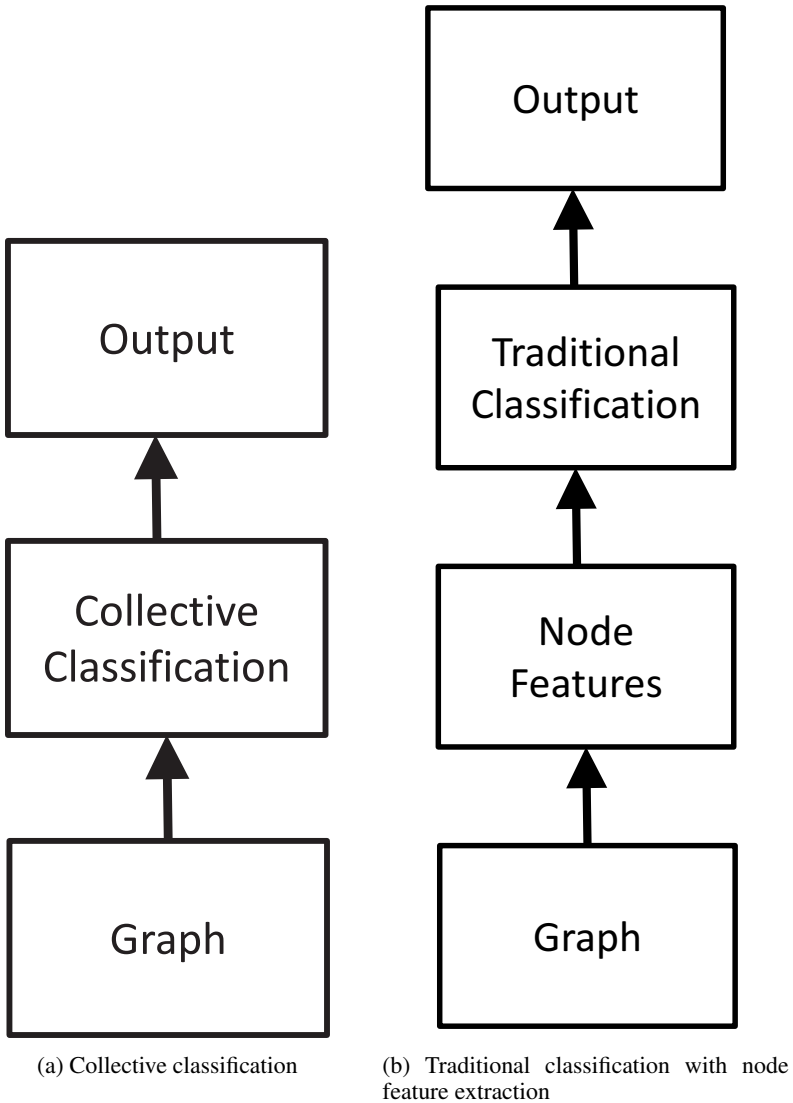


Figure 1.2 Two major directions to develop solutions for node classification on graphs.

discuss a variety of basic neural network models, key approaches to training deep models, and practical techniques to prevent overfitting during training.

- Part II: Methods. These chapters cover the most established methods of deep learning on graphs from the basic to advanced settings. In Chapter 4,

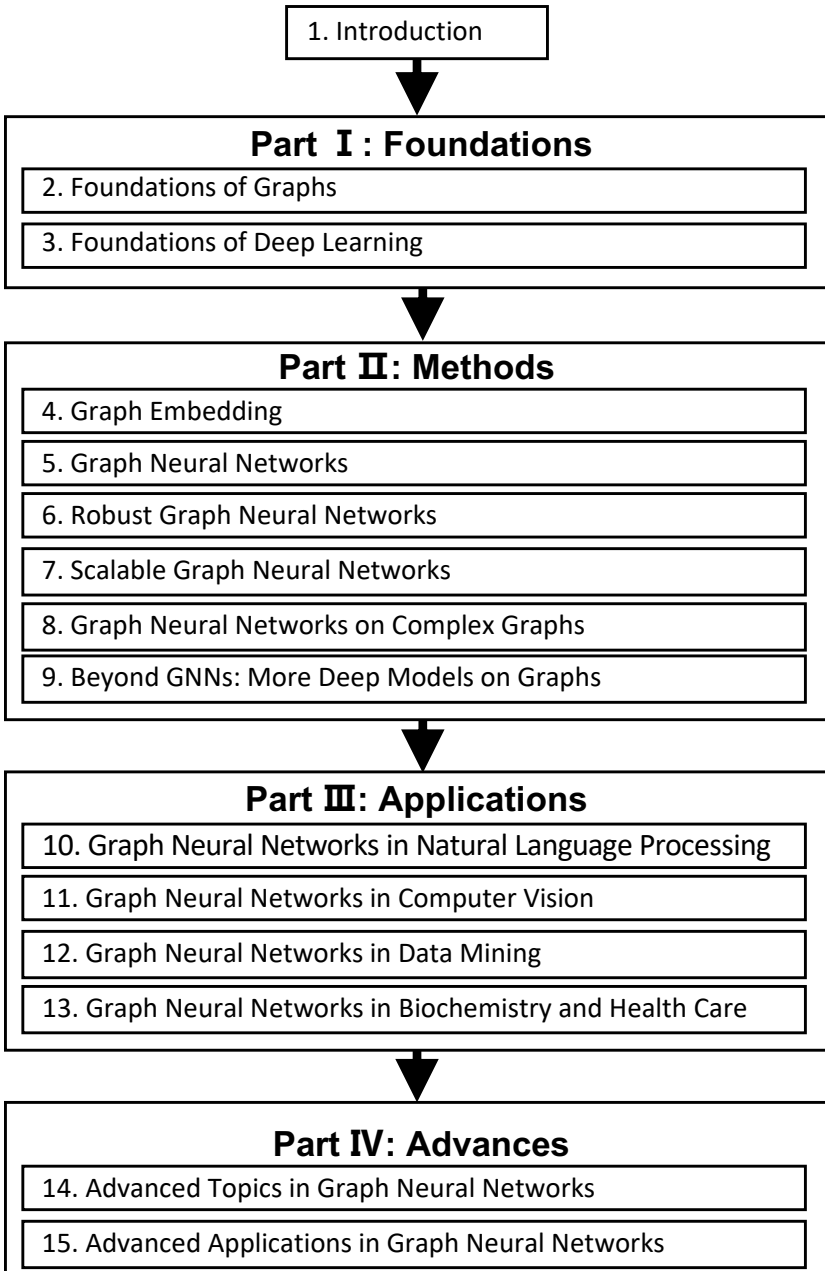


Figure 1.3 The high-level organization of the book.

we introduce a general graph embedding framework from the information preservation perspective, provide technical details on representative algorithms to preserve numerous types of information on graphs, and present embedding approaches specifically designed for complex graphs. A typical graph neural network (GNN) model consists of two important operations, that is, the graph filter operation and the graph pooling operation. In Chapter 5, we review the state-of-the-art graph filter and pooling operations and discuss how to learn the parameters of GNNs for a given downstream task. As generalizations of traditional deep models to graphs, GNNs inherit the drawback of traditional deep models and are vulnerable to adversarial examples. In Chapter 6, we focus on concepts and definitions of graph adversarial attacks and detail representative adversarial attack and defense techniques. GNNs perform the recursive expansion of neighborhoods across layers. The expansion of the neighborhood for a single node can rapidly involve a large portion of the graph or even the whole graph. Hence, scalability is a pressing issue for GNNs. We provide detailed representative techniques to scale GNNs in Chapter 7. In Chapter 8, we discuss GNN models that have been designed for more complicated graphs. To enable deep learning techniques to advance more graph tasks under wider settings, we introduce numerous deep graph models beyond GNNs in Chapter 9.

- Part III: Applications. Graphs provide a universal representation for real-world data; thus methods of deep learning on graphs have been applied to various fields. These chapters present the most representative applications of GNNs, including natural language processing in Chapter 10, computer vision in Chapter 11, data mining in Chapter 12, and biochemistry and health care in Chapter 13.
- Part IV: Advances. These chapters focus on recent advances in both methods and applications. In Chapter 14, we introduce advanced GNNs in terms of expressiveness, depth, fairness, interpretability, and self-supervised learning. We discuss more areas to which GNNs have been applied, including combinatorial optimization, physics, program representation, reinforcement learning, and computer networks in Chapter 15.

1.4 Who Should Read This Book?

This book is easily accessible to readers with a computer science background. Basic knowledge of calculus, linear algebra, probability, and statistics can aid in understanding its technical details. This book has a wide range of target audiences. One target audience is senior undergraduate and graduate students

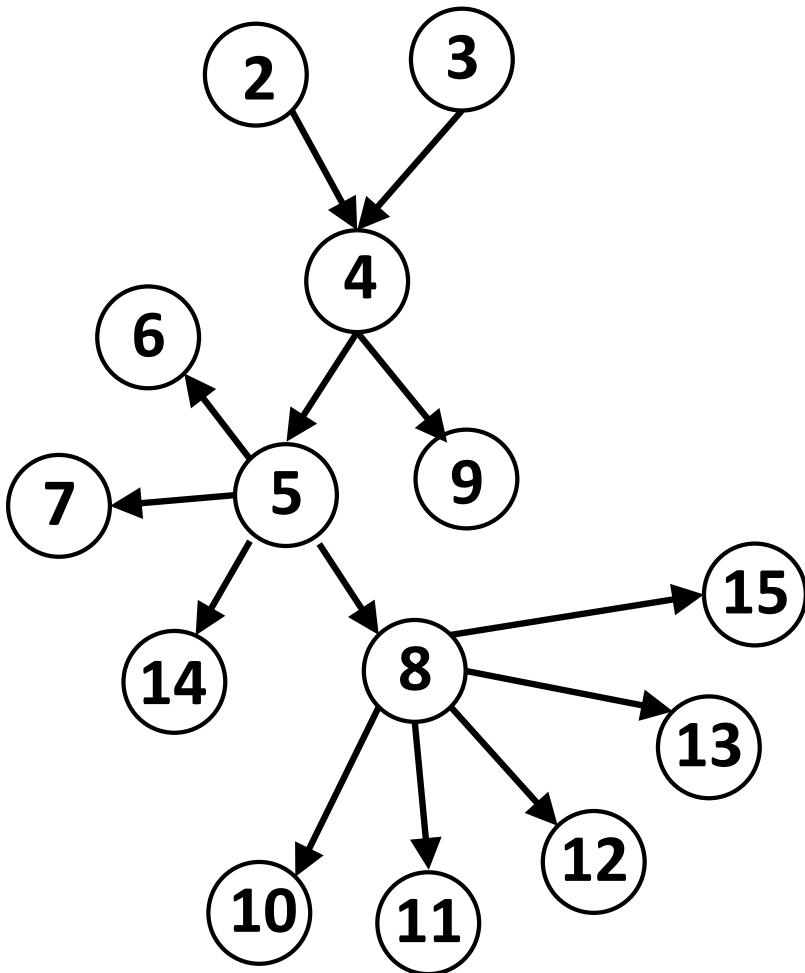


Figure 1.4 The guidance on how to read this book. Note that the number in the circle indicates the corresponding chapter as shown in Figure 1.3.

who are interested in data mining, machine learning, and social network analysis. This book can be independently used for a graduate seminar course on deep learning on graphs. It also can be utilized as a part of a course. For example, Parts II and IV can be considered as advanced topics in courses on data mining, machine learning, and social network analysis; and Part III can be used as advanced methods in solving traditional tasks in computer vision, natural language processing, and health care. Practitioners and project

managers who want to learn the basics and tangible examples of deep learning on graphs and adopt GNNs into their products and platforms are also our target audience. Graph neural networks have been applied to benefit numerous disciplines beyond computer science. Thus, another target audience is researchers who do not have a computer science background but want to apply GNNs to advance their disciplines.

Readers with different backgrounds and purposes for reading should go through the book differently. Suggested guidance on how to read this book is illustrated in Figure 1.4. If readers aim to understand GNN methods for simple graphs (Chapter 5), knowledge about foundations of graphs and deep learning and graph embedding is necessary (Chapters 2–4). If readers want to apply GNNs to advance health care (Chapter 13), they should first read prerequisite materials on foundations of graphs and deep learning, graph embedding, and GNN methods for simple and complex graphs (Chapters 2–5, and 8). For Part III, we do assume that the readers have the necessary background in the corresponding fields of application. In addition, readers should feel free to skip some chapters if they are equipped with the corresponding background. For example, if readers have knowledge regarding foundations of graphs and deep learning, they can skip Chapters 2 and 3 and only read Chapters 4 and 5 to understand GNNs for simple graphs.

1.5 Feature Learning on Graphs: A Brief History

As aforementioned, to take advantage of traditional machine learning for computational tasks on graphs, finding vector node representations is desired. As shown in Figure 1.5, there are two main ways to achieve this goal: feature engineering and feature learning. Feature engineering relies on hand-designed features such as node degree statistics, whereas feature learning is automatic learning of node features. On the one hand, we often do not have prior knowledge of what features are important especially for a given downstream task; thus features from feature engineering could be suboptimal for the downstream task and the process requires an immense amount of human effort. On the other hand, feature learning is automatic learning of features and the process can be guided by the downstream task; consequently, the learned features are likely to be suitable for the downstream task and often obtain better performance than those obtained via feature engineering. Furthermore, the process requires minimal human intervention and can be easily adapted to new tasks. Thus, feature learning on graphs has been extensively studied and various types of

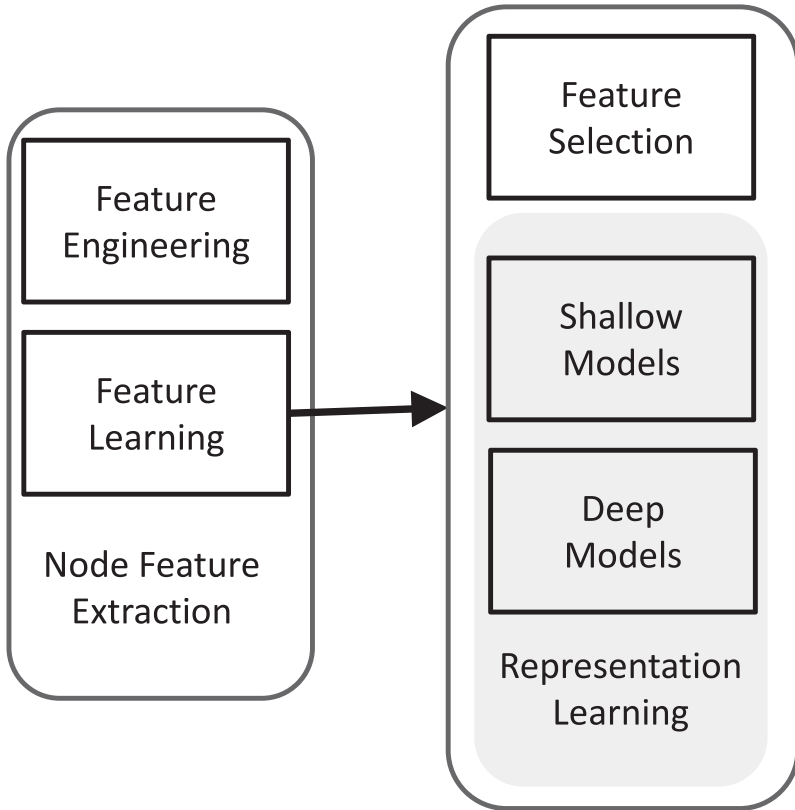


Figure 1.5 Node feature extraction.

feature learning techniques have been proposed to meet different requirements and scenarios. We roughly divide these techniques into feature selection on graphs, which aims to remove irrelevant and redundant node features, and representation learning on graphs, which targets the generation of a set of new node features. In this section, we briefly review these two groups of techniques that provide a general and historical context for readers to understand deep learning on graphs.

1.5.1 Feature Selection on Graphs

Real-world data are often high-dimensional and there exist noisy, irrelevant, and redundant features (or dimensions), particularly when considering

a given task. Feature selection aims to automatically select a small subset of features that have minimal redundancy but maximal relevance to the target, such as the class labels under the supervised setting. In many applications, the original features are crucial for knowledge extraction and model interpretation. For example, in genetic analysis for cancer studies in addition to differentiating cancerous tissues, it is of great importance to identify the genes (i.e., original features) that induce cancerogenesis. In these demanding applications, feature selection is particularly preferred because it maintains the original features and their semantics usually provide critical insights to the learning problem. Traditional feature selection assumes that data instances are independent and identically distributed (i.i.d.). However, data samples in many applications are embedded in graphs that are inherently not i.i.d. This has motivated the research area of feature selection on graphs. Given a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$, where \mathcal{V} is the node set and \mathcal{E} is the edge set, we assume that each node is originally associated with a set of d features $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$. Feature selection on graphs aims to select K features from \mathcal{F} to denote each node where $K \ll d$. The problem was first investigated under the supervised setting (Tang and Liu, 2012a; Gu and Han, 2011). A linear classifier is employed to map from the selected features to the class labels and a graph regularization term is introduced to capture structural information for feature selection. In particular, the term aims to ensure that connected nodes with the selected features can be mapped into similar labels. The problem was further studied under the unsupervised setting (Wei et al., 2015, 2016; Tang and Liu, 2012b). In Tang and Liu (2012b), pseudo labels are extracted from the structural information that serves as the supervision to guide the feature selection process. In Wei et al. (2016), both the node content and structural information are assumed to be generated from a set of high-quality features that can be obtained by maximizing the likelihood of the generation process. The problem has also been extended from simple graphs to complex graphs such as dynamic graphs (Li et al., 2016), multidimensional graphs (Tang et al., 2013b), signed graphs (Cheng et al., 2017; Huang et al., 2020), and attributed graphs (Li et al., 2019b).

1.5.2 Representation Learning on Graphs

Different from feature selection on graphs, representation learning on graphs is to learn a set of new node features. It has been extensively studied for decades and has been greatly accelerated by deep learning. In this subsection, we will give it a brief historical review from shallow models to deep models.

At the early stage, representation learning on graphs was studied under the context of spectral clustering (Shi and Malik, 2000; Ng et al., 2002), graph-based dimension reduction (Belkin and Niyogi, 2003; Tenenbaum et al., 2000;

Roweis and Saul, 2000), and matrix factorization (Zhu et al., 2007; Tang et al., 2013a; Koren et al., 2009). In spectral clustering (Shi and Malik, 2000; Ng et al., 2002), data points are considered as nodes of a graph and then clustering is used to partition the graph into communities of nodes. One key step for spectral clustering is spectral embedding. It aims to embed nodes into a low-dimensional space where traditional clustering algorithms such as k -means can be applied to identify clusters. Techniques of graph-based dimension reduction can be directly applied to learn node representations. These approaches typically build an affinity graph using a predefined distance (or similarity) function based on the raw features of data samples and then they aim to learn node representations to preserve structural information of this affinity graph. For example, IsoMap (Tenenbaum et al., 2000) preserves the global geometry via geodesics, locally linear embedding (Roweis and Saul, 2000) and eigenmap (Belkin and Niyogi, 2003) preserve local neighborhoods in the affinity graph. The aforementioned methods often need to perform eigendecomposition on the affinity matrix (or adjacency matrix or Laplacian matrix). Thus, they are often computationally expensive. Matrix is one of the most popular approaches to denote graphs such as adjacency matrices, incidence matrices, and Laplacian matrices. As a result, matrix factorization can be naturally applied to learn node representations. Using the adjacency matrix to denote a graph as an example, it aims to embed nodes into a low-dimensional space where the new node representations can be utilized to reconstruct the adjacency matrix. A document corpus can be denoted as a bipartite graph where documents and words are nodes and an edge exists between a word and a document if the word appears in the document. Latent semantic indexing has employed truncated singular value decomposition to learn representations of documents and words (Deerwester et al., 1990). In recommender systems, interactions between users and items can be captured as a bipartite graph and matrix factorization has been employed to learn representations of users and items for recommendations (Koren et al., 2009). Matrix factorization is also leveraged to learn node representations for node classification (Zhu et al., 2007; Tang et al., 2016a), link prediction (Menon and Elkan, 2011; Tang et al., 2013a), and community detection (Wang et al., 2011). Actually, a family of modern graph embedding algorithms we will introduce later can also be unified as matrix factorization (Qiu et al., 2018b).

Word2vec is a technique to generate word embeddings (Mikolov et al., 2013). It takes a large corpus of text as input and produces a vector representation for each unique word in the corpus. The huge success of Word2vec in various natural language processing tasks has motivated increasing efforts to apply Word2vec, especially the SkipGram model, to learn node representations in the graph domain. DeepWalk (Perozzi et al., 2014) takes the first step to

achieve this goal. Specifically, nodes in a given graph are treated as words of an artificial language and sentences in this language are generated by random walks on the graph. Then, it uses the SkipGram model to learn node representations, which preserves co-occurring nodes in these random walks. A large body of works have been developed in three major directions: (1) Developing advanced methods to preserve co-occurring nodes (Tang et al., 2015; Grover and Leskovec, 2016; Cao et al., 2015); (2) Preserving other types of information such as a node's structural role (Ribeiro et al., 2017), community information (Wang et al., 2017c), and node status (Ma et al., 2017; Lai et al., 2017; Gu et al., 2018); and (3) Designing frameworks for complex graphs such as directed graphs (Ou et al., 2016), heterogeneous graphs (Chang et al., 2015; Dong et al., 2017), bipartite graphs (Gao et al., 2018b), multidimensional graphs (Ma et al., 2018d), signed graphs (Wang et al., 2017b), hyper graphs (Tu et al., 2018), and dynamic graphs (Nguyen et al., 2018; Li et al., 2017a).

Given the power and the success of deep neural networks (DNNs) in representation learning, increasing efforts have been made to generalize DNNs to graphs. These methods are known as graph neural networks and can be roughly divided into spatial approaches and spectral approaches. Spatial approaches explicitly leverage the graph structure, such as spatially close neighbors, and the first spatial approach was introduced by Scarselli et al. in 2005. Spectral approaches utilize the spectral view of graphs by taking advantage of graph Fourier transform and the inverse graph Fourier transform (Bruna et al., 2013). In the era of deep learning, GNNs have been rapidly developed in the following aspects:

- A huge number of new GNN models have been introduced, including spectral approaches (Defferrard et al., 2016; Kipf and Welling, 2016a) and spatial approaches (Atwood and Towsley, 2016; Niepert et al., 2016; Gilmer et al., 2017; Monti et al., 2017; Veličković et al., 2017; Hamilton et al., 2017a).
- For graph-focused tasks such as graph classification, representation of the whole graph is desired. Thus, numerous pooling methods have been introduced to obtain the graph representation from node representations (Li et al., 2015; Ying et al., 2018c; Gao and Ji, 2019; Ma et al., 2019b).
- Traditional DNNs are vulnerable to adversarial attacks. GNNs inherit this drawback. A variety of graph adversarial attacks have been studied (Zügner et al., 2018; Zügner and Günnemann, 2019; Dai et al., 2018; Ma et al., 2020a) and various defense techniques have been developed (Dai et al., 2018; Zhu et al., 2019a; Tang et al., 2019; Jin et al., 2020b).

- As aforementioned, scalability is a pressing issue for GNNs. Many strategies have been studied to allow GNNs scale to large graphs (Chen et al., 2018a,b; Huang et al., 2018).
- GNN models have been designed to handle complex graphs such as heterogeneous graphs (Zhang et al., 2018b; Wang et al., 2019i; Chen et al., 2019b), bipartite graphs (He et al., 2019), multidimensional graphs (Ma et al., 2019c), signed graphs (Derr et al., 2018), hyper graphs (Feng et al., 2019b; Yadati et al., 2019), and dynamic graphs (Pareja et al., 2019).
- Diverse deep architectures have been generalized for graphs, such as autoencoder (Wang et al., 2016; Cao et al., 2016), variational autoencoder (Kipf and Welling, 2016b), recurrent neural networks (Tai et al., 2015; Liang et al., 2016), and generative adversarial networks (Wang et al., 2018a).
- Because graphs are a universal data representation, GNNs have been applied to advance a number of fields, such as natural language processing, computer vision, data mining, and health care.

1.6 Conclusion

In this chapter, we discuss the opportunities and challenges when we bridge deep learning with graphs that have motivated the focus of this book – deep learning on graphs. The book will cover the essential topics of deep learning on graphs, which are organized into four parts to accommodate readers with diverse backgrounds and purposes for reading, including foundations, methods, applications, and advances. This book can benefit a wide range of readers, including senior undergraduate students, graduate students, practitioners, project managers, and researchers from various disciplines. To provide more context for readers, we give a brief historical review on the area of feature learning on graphs.

1.7 Further Reading

In this chapter, we have briefly reviewed the history of feature selection on graphs. If readers want to know more about feature selection, there are quite a few books (Liu and Motoda, 2007, 2012) and comprehensive surveys (Tang et al., 2014a). An open source feature selection repository scikit-feature has been developed that consists of most of the popular feature selection algorithms (Li et al., 2017b). This is the first comprehensive book on the topic of deep learning on graphs. There are books on the general topics of deep

learning (Goodfellow et al., 2016; Aggarwal, 2018), deep learning on speech recognition (Yu and Deng, 2016; Kamath et al., 2019), and deep learning in natural language processing (Deng and Liu, 2018; Kamath et al., 2019).