

Learning Hierarchical Review Graph Representations for Recommendation

Yong Liu, Susen Yang, Yinan Zhang, Chunyan Miao, Zaiqing Nie, and Juyong Zhang

Abstract—The user review data have been demonstrated to be effective in solving different recommendation problems. Previous review-based recommendation methods usually employ sophisticated compositional models, such as Recurrent Neural Networks (RNN) and Convolutional Neural Networks (CNN), to learn semantic representations from the review data for recommendation. However, these methods mainly capture the local dependency between neighbouring words in a word window, and they treat each review equally. Therefore, they may not be effective in capturing the global dependency between words, and tend to be easily biased by noise review information. In this paper, we propose a novel review-based recommendation model, named Review Graph Neural Network (RGNN). Specifically, RGNN builds a specific review graph for each individual user/item, which provides a global view about the user/item properties to help weaken the biases caused by noise review information. A type-aware graph attention mechanism is developed to learn semantic embeddings of words. Moreover, a personalized graph pooling operator is proposed to learn hierarchical representations of the review graph to form the semantic representation for each user/item. We compared RGNN with state-of-the-art review-based recommendation approaches on two real-world datasets. The experimental results indicate that RGNN consistently outperforms baseline methods, in terms of Mean Square Error (MSE).

Index Terms—Review-based Recommendation, Hierarchical Graph Representation Learning, Graph Neural Networks.

1 INTRODUCTION

WITH the explosive growth of online information and contents, recommendation systems are playing an increasingly important role in various scenarios, *e.g.*, E-commerce websites and online social media platforms. In the last decades, various techniques [1], [2] have been proposed to solve different recommendation problems, for example business rating prediction [3], Point-of-Interests recommendation [4], and diversified item recommendation [5]. Despite many research efforts have been devoted to developing recommendation systems, existing methods may still suffer from the data sparsity problem. To remedy this problem, various types of side information have been incorporated into the recommendation systems, *e.g.*, the item knowledge graph [6] and the users' reviews on items [7]. In this work, we focus on how to effectively exploit the user review data which usually describe users' preferences and also include their interests on different aspects of items.

Conventional review-based recommendation methods incorporate topic modeling techniques, *e.g.*, Latent Dirich-

let Allocation (LDA) [8] and word embedding model [9], with matrix factorization to jointly learn semantic latent factors for users and items [10], [11], [12], [13], [14]. Recently, the rapid development of deep learning motivates the applications of different neural compositional models, *e.g.*, Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN), to incorporate the review data for recommendation [15], [16], [17], [18], [19]. In addition, there also exists some recent work performing sentiment analysis on the review data to learn aspect-based representations for both users and items [20], [21]. These review-based recommendation methods mainly focus on the rating prediction problem. Following these work, we also study the rating prediction task in this paper.

Despite the successes of existing review-based recommendation methods, they may still suffer from some of the following limitations. Firstly, the methods that only consider document-level co-occurrences of words [10], [11] can only capture the global-level semantic information. They usually ignore the word orders and the local context of words. Secondly, the word embedding and neural compositional models (*e.g.*, CNN and RNN) are effective in capturing the context information from neighbouring words in a word window. However, these methods treat each review equally. As they only rely on the local context of a word, they may be easily biased by the noise review information that are not relevant to the user interests or item properties. Thirdly, most existing methods use the same model to learn the representations of all users and items, without considering the specific characteristics of an individual user/item.

To overcome above limitations, we propose a novel review-based recommendation method, namely Review Graph Neural Network (RGNN). The contributions made in this paper are summarized as follows.

- Yong Liu is currently with Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly (LILY) and Alibaba-NTU Singapore Joint Research Institute, Nanyang Technological University, Singapore 639798. Email: stephenliu@ntu.edu.sg.
- Yinan Zhang and Chunyan Miao are currently with Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly (LILY), Alibaba-NTU Singapore Joint Research Institute, and School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798. Email: yinan002@e.ntu.edu.sg, ascymiao@ntu.edu.sg.
- Susen Yang and Juyong Zhang are currently with School of Mathematical Sciences, University of Science and Technology of China, Hefei, Anhui, China 230052. Email: susen@mail.ustc.edu.cn, juyong@ustc.edu.cn.
- Zaiqing Nie is currently with Institute for AI Industry Research, Tsinghua University, Beijing, China. Email: zaiqing@air.tsinghua.edu.cn. This work was done when he was working in Alibaba Group, Beijing, China.
- Corresponding author: Juyong Zhang

- We propose to build a directed review graph to describe the reviews of a user/item, where nodes represent the words and edges describe the co-occurrences and order relationships between words. The review graph aggregates the information of all the reviews associated with the user/item to provide a global view about the user/item properties. Then, the reviews with similar contents tend to have more dense connections with each other. Therefore, the review graph can help weaken the biases caused by the noise reviews.
- We develop the Type-aware Graph Attention (TGAT) mechanism to learn the word embeddings, by aggregating the context information from neighboring words and considering the word orders.
- Moreover, we also develop the Personalized Graph Pooling (PGP) operator to generate the hierarchical representations of the user/item review graph, considering the user/item-specific properties. PGP is able to select different informative words for different users/items, based on the intuition that the same word may have different semantic informativeness with respect to different users/items.
- To demonstrate the effectiveness of RGNN, we perform extensive experiments on two real datasets, compared with state-of-the-art review-based recommendation methods.

The remaining parts of this work are organized as follows. In Section 2, we review the most relevant existing work. Then, in Section 3, we introduce the details of the proposed recommendation model. Next, Section 4 reports the experimental results. Finally, Section 5 concludes this work and discusses the potential future work directions.

2 RELATED WORK

This section reviews existing work about review-based recommendation, graph neural networks, and graph pooling.

2.1 Review-based Recommendation Methods

Traditional review-based recommendation methods employ topic modeling techniques on the review data to learn latent feature distributions of users and items. For example, McAuley and Leskovec [10] propose the HFT (*i.e.*, Hidden Factors as Topics) model that learns the latent factors of users and items from review texts by a LDA model. Ling *et al.* [12] propose the RMR (*i.e.*, Ratings Meet Reviews) model that combines the topic model with a mixture of Gaussian models to improve the recommendation accuracy. In [11], Bao *et al.* employ non-negative matrix factorization to derive topics of each review, and a transform function to align the topic distributions with corresponding latent user/item factors. In [13], Tan *et al.* propose a rating-boosted method that combines textual reviews with users' sentiment orientations to learn more accurate topic distributions.

The quick development of deep learning motivates the application of deep learning techniques to extract semantic information from the review data and build the end-to-end recommendation models. For instance, Wang *et al.* [15] propose the pioneer work that employs Stacked Denoising AutoEncoders (SDAE) to learn deep representations from review contents, and probabilistic matrix factorization to

model users' rating behaviors. Zheng *et al.* [16] propose the Deep Cooperative Neural Networks (DeepCoNN) that employs CNN to extract semantic information from the user reviews, and factorization machine to capture the interactions between the user and item features. Recently, the attention mechanism is applied to further improve the recommendation accuracy. In [17], Chen *et al.* extract the review representation using CNN, and scored each review through an attention mechanism to learn informative user and item representations. In [22], Wu *et al.* utilize CNN with attention mechanism to highlight relevant semantic information by jointly considering the reviews written by a user and the reviews written to an item. Moreover, Lu *et al.* [23] propose to jointly learn the user and item information from ratings and reviews by optimizing matrix factorization and an attention-based GRU (*i.e.*, Gated Recurrent Unit) network. Liu *et al.* [18] propose the dual attention mutual learning between ratings and reviews for recommendation. They use both local and mutual attention of the CNN to jointly learn review representations for improving the model interpretability. Wu *et al.* [24] utilize three attention networks to learn representations of sentence, review, and user/item in turn, and a graph attention network to model user-item interactions. To incorporate the user/item characteristics, Liu *et al.* [19] apply a personalized attention model to learn personalized representations of users and items from the review data. In addition, there are some aspect-based recommendation models [20], [21] that apply aspect-based representation learning to automatically detect the most important parts of a review, and jointly learn the aspect, user, and item representations.

2.2 Graph Neural Networks

The graph neural networks (GNNs) extend deep learning techniques to process the graph data. Two main categories of GNN methods are recurrent graph neural networks and graph convolutional networks [25].

The pioneer GNN methods are about recurrent graph neural networks, which learn the representation of a node by propagating its neighborhood information in an iterative manner until convergence. For example, in [26], the node's states are updated by exchanging its neighborhood information recurrently before satisfying the convergence criterion. Gallicchio and Micheli [27] improve the training efficiency of [26] by generalizing the echo state network to graph domains. Moreover, Li *et al.* [28] propose the Gate Graph Neural Networks (GGNN) which employs a gated recurrent unit as the recurrent function. This method no longer needs parameter constraints to ensure convergence.

Graph Convolutional Networks (GCNs) generalize the convolution operation from grid data to graph data. Existing GCN methods usually fall into two groups, *i.e.*, spectral-based methods and spatial-based methods. The spectral-based methods define convolution from the perspective of graph signal processing. For example, Bruna *et al.* [29] employs the eigen decomposition of the graph Laplacian and defines the convolution in the Fourier domain. To reduce computational complexity, Defferrard *et al.* [30] approximate the convolutional filters by Chebyshev polynomials of the diagonal matrix of eigenvalues. Then, Kipf and Welling [31]

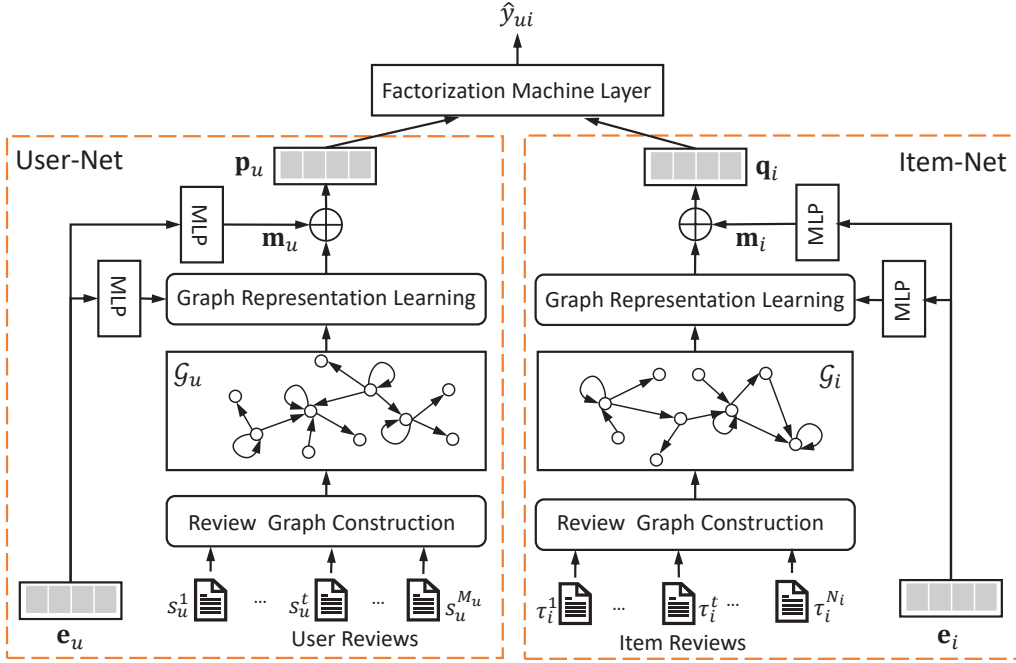


Fig. 1: The overall framework of the proposed RGNN model.

develop a simpler convolutional architecture by the first-order approximation of the CheNet [30]. The spatial-based methods directly operate on the graph and propagate node information along edges. For instance, Micheli [32] proposes to directly sum up a node’s neighborhood information and applies residual connection to memorize information over each layer. Atwood and Towsley [33] treat the graph convolution as a diffusion process, and they assume there is a certain transition probability from a node to its neighbor. As the number of neighbors of nodes varies, Hamilton *et al.* [34] propose the GraphSage model which samples a fixed number of neighbors for each node and utilizes three different aggregators to aggregate the neighbors’ feature information. In [35], Veličković *et al.* propose the Graph Attention Networks (GAT) that assumes the contributions of neighbor nodes to the central node are different. The multi-head attention mechanism is performed to calculate the importance of different neighbors. Moreover, Zhang *et al.* [36] employ a self-attention mechanism to compute an additional attention score for each attention head. In the recent work [37], Xu *et al.* identify the graph structures that cannot be distinguished by popular GNN methods, and propose the Graph Isomorphism Network (GIN) which uses a learnable parameter to adjust the weight of the central node.

2.3 Graph Pooling

The pooling operation usually helps provide better performance in graph representation tasks. Earlier pooling methods are mainly based on the deterministic graph clustering algorithms [38], [39], which are less flexible. The recent graph pooling methods can be divided into two groups: spectral and non-spectral methods. The spectral pooling methods [40], [41] employ laplace matrix and eigen-decomposition to capture the graph topology. However,

these methods are not scalable to large graph, due to high computation complexity. Thus, a lot of research attentions have been attracted to develop the non-spectral pooling methods, which can be further divided into global and hierarchical pooling methods. The global pooling methods [42], [43] pool the representations of all nodes at one step. For example, Vinyals *et al.* [42] propose the Set2Set model that aggregates nodes information by combining Long Short Term Memory (LSTM) with a content-based attention mechanism. Zhang *et al.* [43] develop the SortPool model that first sorts nodes according to their features and then concatenates the embedding of sorted nodes to summarize the graph. Differing from the global pooling methods, the hierarchical pooling methods aim to learn hierarchical graph representations. For example, Ying *et al.* [44] propose the first end-to-end differentiable pooling operator, named DiffPool, which learns a dense soft clustering assignment matrix of nodes by GNN. The gPool [45], [46] and SAG-Pool [47] pooling operators utilize a node-scoring procedure to select a set of top ranked nodes to form an induced graph from original graph. Moreover, Diehl [48] proposes a hard pooling method called EdgePool, which is developed based on edge contraction. In [49], Ranjan *et al.* utilize a self-attention mechanism to capture the importance of each node. This method learns a sparse soft cluster assignment for nodes to pool the graph. In [50], Zhang *et al.* propose a non-parametric pooling operation using a structure learning mechanism with sparse attention to learn a refined graph structure that preserves the key substructure in original graph.

3 THE PROPOSED RECOMMENDATION MODEL

This work focuses on exploiting users’ review data for recommendation. For each user u , the set of her reviews is

denoted by $\mathcal{S}_u = \{s_u^1, s_u^2, \dots, s_u^{M_u}\}$, where s_u^t denotes the t -th review of user u , and M_u denotes the number of reviews generated by the user u . Similarly, let $\mathcal{T}_i = \{\tau_i^1, \tau_i^2, \dots, \tau_i^{N_i}\}$ denote the reviews of each item i , where N_i denotes the number of reviews generated on the item i , and τ_i^t denotes the t -th review of item i . Note that each review may consist of multiple text sentences. In addition, we denote the observed rating given by the user u on the item i by y_{ui} . For a word x , we denote its word embedding by $\mathbf{e}_x \in \mathbb{R}^{1 \times d_0}$, where d_0 denotes the dimensionality of the latent semantic space. Moreover, we use $\mathbf{e}_u \in \mathbb{R}^{1 \times d_1}$ and $\mathbf{e}_i \in \mathbb{R}^{1 \times d_1}$ to denote the embedding of user u and item i , respectively, where d_1 is the dimensionality of latent rating space.

The overall framework of the proposed RGNN model is summarized in Figure 1. As shown in Figure 1, RGNN contains three components, *i.e.*, the User-Net, Item-Net, and a prediction layer. The User-Net and Item-Net are two parallel networks used to learn semantic representations for users and items, respectively. The prediction layer is developed based on Factorization Machine (FM) [51], which aims to predict the rating scores by considering the interactions between learned semantic features. Both User-Net and Item-Net consist of the following three modules:

- The review graph construction module builds the review graph for each user and item, considering the review words as nodes in the graph and capturing the co-occurrence and order relationships between words.
- The graph representation learning module applies a Type-aware Graph Attention Network (TGAT) and Personalized Graph Pooling operator (PGP) to extract the user/item-specific hierarchical representations from the entire review graph.
- Moreover, the feature fusion module integrates the hierarchical representations of the user/item review graph and the user/item embedding to form the semantic representation for the user/item.

As User-Net and Item-Net only differ in their inputs, we only introduce the details of User-Net in the following sections. The same procedure is applied to the Item-Net.

3.1 Review Graph Construction

Various graph based methods have been proposed to process the text data [52], which aim to capture the inherent dependence between words. Here, we utilize the graph-of-words method [53] to construct the review graph with preserving the word orders. For the reviews \mathcal{S}_u of the user u , we firstly apply the phrase pre-processing techniques such as tokenization and text cleaning to select keywords of each review in \mathcal{S}_u . Then, all reviews in \mathcal{S}_u can be encoded into an un-weighted and directed graph, where nodes represent words and edges describe the co-occurrences between words within a sliding window with fixed size ω . Note that we define the word windows considering the positions of selected keywords in the original text.

Intuitively, the order information of words in the text may help capture the sentiment information [54]. For example, the phrases “not really good” and “really not good” convey different levels of negative sentiment information, which is only caused by the word orders. To preserve the orders of words in the review, we define three types of

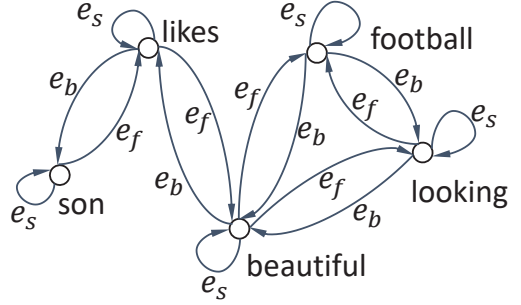


Fig. 2: An example of constructing graph for the review “My son likes this beautiful looking football”. The words “son”, “likes”, “beautiful”, “looking”, and “football” are selected as nodes. Empirically, we set the window size ω to 3 for edge building.

edges in the review graph: **forward type**, **backward type**, and **self-loop type** (respectively denoted by e_f , e_b , and e_s). We denote the set of edge type by $\mathcal{R} = \{e_f, e_b, e_s\}$. In a review $s_u^t \in \mathcal{S}_u$, if the selected keyword (*i.e.*, node) x_2 appears before the selected keyword x_1 and the distance between x_1 and x_2 in the original text of s_u^t is *less than* ω , we construct an edge $E(x_2, x_1)$ from x_2 to x_1 and set the type of edge $E(x_2, x_1)$ as e_f . Moreover, we also construct an edge $E(x_1, x_2)$ from x_1 to x_2 and set the type of $E(x_1, x_2)$ as e_b . In the rare situation¹ where the word x_2 may appear both before and after the word x_1 , we randomly set the type of edge $E(x_1, x_2)$ as e_f or e_b . In the experiments, we empirically set ω to 3. To consider the information of the word itself (*e.g.*, x_1), we add a self-loop edge at each word (*e.g.*, $E(x_1, x_1)$) and set the edge type as e_s . This trick also helps define the attention score in Eq. (1).

Figure 2 shows an example about constructing the review graph for the review “My son likes this beautiful looking football”, where the words “son”, “likes”, “beautiful”, “looking”, and “football” are selected as keywords to build the review graph. In the original review text, the distance between “son” and “likes” is 1 which is less than ω . Thus, we construct an edge from “son” to “likes” and set the edge type as e_f . Moreover, we also construct an edge from “likes” to “son” and set the edge type as e_b . However, the distance between “son” and “beautiful” in original text is 3, which is not smaller than the windows size ($\omega = 3$). Therefore, as shown in Figure 2, there is no connection between the selected keywords “son” and “beautiful”.

For each user u , we denote her review graph by $\mathcal{G}_u = \{\mathcal{X}_u, \mathcal{E}_u\}$, where \mathcal{X}_u denotes the set of nodes (*i.e.*, selected keywords) and \mathcal{E}_u denotes the set of node-edge-node triples (x_h, e_r, x_t) . Here, e_r is type of the edge connecting from node x_t to x_h , which is one of the three edge types defined above. Similarly, we can build the review graph $\mathcal{G}_i = \{\mathcal{X}_i, \mathcal{E}_i\}$ for each item i .

1. We have empirically studied the ratio of edges that have more than one edge type on the experimental datasets. By setting the word window size ω to 3 (*i.e.*, our default setting), there are only 0.76%, 0.86%, and 0.73% of the edges between neighbouring words in the user review graph have more than one edge types, on the Music, Office, and Toys datasets, respectively.

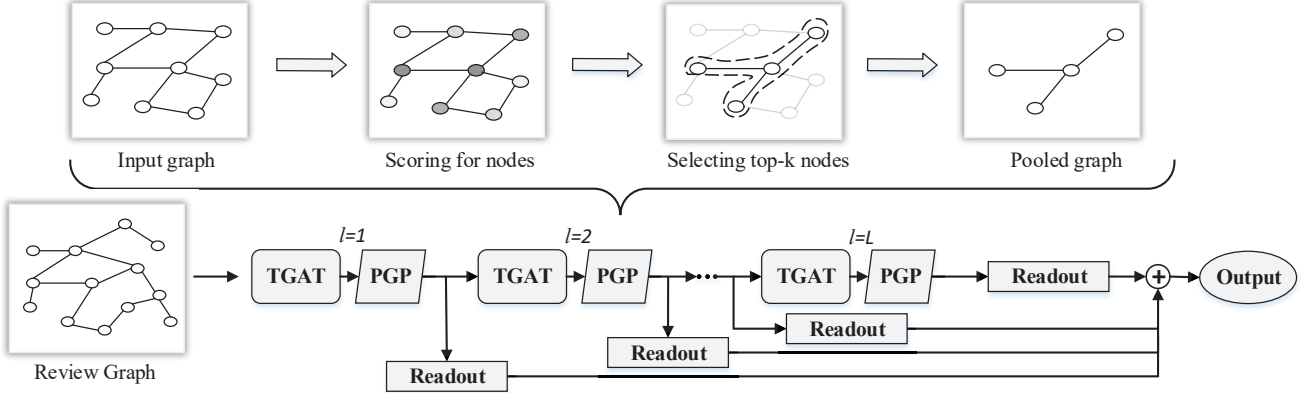


Fig. 3: The procedure of hierarchical graph representation learning.

3.2 Hierarchical Graph Representation Learning

After completing the construction of review graph \mathcal{G}_u , we employ the representation learning procedure as shown in Figure 3 to learn hierarchical representations of \mathcal{G}_u , by repeating the TGAT and PGP operations multiple times. In this section, we introduce the details of the operations in the ℓ -th representation learning layer, where the input graph is $\mathcal{G}_u^\ell = \{\mathcal{X}_u^\ell, \mathcal{E}_u^\ell\}$. Suppose \mathcal{G}_u^ℓ has K_u^ℓ nodes. Then, we also use an adjacency matrix $\mathbf{A}_u^\ell \in \mathbb{R}^{K_u^\ell \times K_u^\ell}$ and a feature matrix $\mathbf{X}_u^\ell \in \mathbb{R}^{K_u^\ell \times d_2^\ell}$ to describe \mathcal{G}_u^ℓ , where each row in $\mathbf{X}_u^\ell(x_1, \cdot)$ denotes the input feature of the node x_1 , and d_2^ℓ denotes the dimensionality of the input feature. Note that the features of nodes in the review graph are initialized by the embeddings of corresponding words. Thus, at the first layer, d_2^ℓ is equal to d_0 . In \mathcal{G}_u^ℓ , if there exists an edge between two nodes x_1 and x_2 , we set $\mathbf{A}_u^\ell[x_1, x_2]$ to 1, and set to 0, otherwise.

3.2.1 Type-aware Graph Attention Network

In the reviews, the semantic meaning of a word can be enriched by its neighboring words, which can be considered as the dependency between a node and its neighboring nodes in the review graph. At the ℓ -th representation learning layer, for a node x_h in the graph \mathcal{G}_u^ℓ , we denote its neighborhood in the graph by $\mathcal{N}_h^\ell = \{x_t | (x_h, e_r, x_t) \in \mathcal{E}_u^\ell\}$. Note that \mathcal{N}_h^ℓ also includes the node x_h , because there exists a self-loop edge at each node in the review graph. A type-aware attention mechanism is developed to effectively aggregate the information of neighboring words, considering the types of connecting edges. Following [35], we define the attention weight $a^\ell(x_h, e_r, x_t)$ of a neighboring node $x_t \in \mathcal{N}_h^\ell$ as,

$$a^\ell(x_h, e_r, x_t) = \frac{\exp[\pi^\ell(x_h, e_r, x_t)]}{\sum_{x_{\bar{t}} \in \mathcal{N}_h^\ell} \exp[\pi^\ell(x_h, e_r, x_{\bar{t}})]}, \quad (1)$$

where $\pi^\ell(x_h, e_r, x_t)$ is implemented by the following relational attention mechanism,

$$\pi^\ell(x_h, e_r, x_t) = \text{LeakyReLU}[(\mathbf{x}_h^\ell \mathbf{W}_1^\ell)(\mathbf{x}_t^\ell \mathbf{W}_2^\ell + \mathbf{e}_r^\ell \mathbf{W}_3^\ell)^\top]. \quad (2)$$

Here, $\mathbf{x}_h^\ell, \mathbf{x}_t^\ell \in \mathbb{R}^{1 \times d_2^\ell}$ are the input embeddings of x_h and x_t at the ℓ -th layer. $\mathbf{e}_r^\ell \in \mathbb{R}^{1 \times d_2^\ell}$ is the embeddings of e_r at the ℓ -th layer. $\mathbf{W}_1^\ell, \mathbf{W}_2^\ell, \mathbf{W}_3^\ell \in \mathbb{R}^{d_2^\ell \times d_2}$ are three different weight matrices used for the linear transformations of the node

and edge embeddings, and d_2 denotes the dimensionality of the hidden space of the graph representation learning layers. The attention weights indicate which neighboring node should contribute more to the central node, when aggregating neighborhood information. Then, we can compute the linear combination of neighbors as follows,

$$\bar{\mathbf{x}}_h^\ell = \sum_{x_t \in \mathcal{N}_h^\ell} a^\ell(x_h, e_r, x_t) \mathbf{x}_t^\ell \mathbf{W}_4^\ell, \quad (3)$$

where $\mathbf{W}_4^\ell \in \mathbb{R}^{d_2^\ell \times d_2}$ denotes the transform matrix. A feature matrix $\bar{\mathbf{X}}_u^\ell \in \mathbb{R}^{K_u^\ell \times d_2}$ is used to denote the output of TGAT. Each row of $\bar{\mathbf{X}}_u^\ell$ corresponds to the feature of a node x in \mathcal{G}_u^ℓ , computed by Eq. (3). Then, $\bar{\mathbf{X}}_u^\ell$ will be used as the input of the PGP operation at the ℓ -th representation learning layer.

Compared with the methods using CNN to extract semantic information, the proposed TGAT operation can not only capture the global contextual word information from all the reviews of a user/item, but also specify varying importance weights to different contextual words. In addition, we explicitly incorporate the edge type e_r into the attention mechanism, which encodes word order information during aggregating neighboring contextual word information.

3.2.2 Personalized Graph Pooling

The graph pooling operation plays an important role in learning graph representations. Various hierarchical graph pooling methods have been proposed to adaptively select a subset of important nodes to form a smaller sub-graph [45], [46], [47], [49], [50]. However, these methods ignore the user/item-specific properties of an individual user/item, thus they are not suitable to learn the review graph representation for recommendation. In practice, different users may emphasize different aspects of items. For example, some users are most interested in the item price, while some other users may be more interested in the item quality. Similarly, the most important properties of different items may also be different. Therefore, it is intuitively that the same word may have different importance to different users/items. Motivated by this intuition, we propose the personalized graph pooling operator to capture the important nodes (*i.e.*, words) for each user/item.

For the user u , we first use a Multilayer Perceptron (MLP) to generate her personalized measurement vector θ_u^l

for measuring the importance of nodes at the ℓ -th representation learning layer,

$$\theta_u^\ell = \text{ReLU}(\mathbf{e}_u \mathbf{W}_5^\ell + \mathbf{b}_5^\ell), \quad (4)$$

where $\mathbf{W}_5^\ell \in \mathbb{R}^{d_1 \times d_2}$ and $\mathbf{b}_5^\ell \in \mathbb{R}^{1 \times d_2}$ are the weight matrix and bias term of MLP at the ℓ -th pooling layer. We adapt graph pooling to the recommendation scenario studied in this work, and define the following importance score β_h^ℓ for each node $x_h \in \mathcal{X}_u^\ell$ as,

$$\beta_h^\ell = |\bar{\mathbf{x}}_h^\ell \theta_u^{\ell \top}| + \frac{1}{|\mathcal{N}_h^\ell|} \sum_{x_t \in \mathcal{N}_h^\ell} |(\bar{\mathbf{x}}_h^\ell - \bar{\mathbf{x}}_t^\ell) \mathbf{w}^{\ell \top}|, \quad (5)$$

where $\mathbf{w}^\ell \in \mathbb{R}^{1 \times d_2}$ is a learnable parameter. In Eq. (5), the first part captures the importance of node x_h to the user u . The second part describes the feature diversity of x_h 's neighborhood at the ℓ -th layer. It is developed based on the criterion that a node x_h can be removed in the pooled graph with negligible information loss, if x_h can be well reconstructed by its neighboring nodes in \mathcal{G}_u^ℓ . Moreover, \mathbf{w}^ℓ can also help balance the contributions of these two parts in the ℓ -th layer.

Then, the score vector $\beta^\ell = [\beta_1^\ell, \beta_2^\ell, \dots, \beta_{K_u^\ell}^\ell]^\top$ is multiplied with the feature matrix $\bar{\mathbf{X}}_u^\ell$ to obtain the following filtered feature matrix,

$$\hat{\mathbf{X}}_u^\ell = \mathbf{D}^\ell \bar{\mathbf{X}}_u^\ell, \quad (6)$$

where $\mathbf{D}^\ell \in \mathbb{R}^{K_u^\ell \times K_u^\ell}$ is a diagonal matrix. The diagonal element $D_{hh}^\ell = \sigma(\beta_h^\ell)$, and $\sigma(\cdot)$ denotes the sigmoid function. Then, we can obtain the adjacency matrix $\mathbf{A}_u^{\ell+1}$ and feature matrix $\mathbf{X}_u^{\ell+1}$ of the pooled graph $\mathcal{G}_u^{\ell+1}$ as follows,

$$\text{idx} = \text{rank}(\beta^\ell, K_u^{\ell+1}), \quad (7)$$

$$\mathbf{A}_u^{\ell+1} = \mathbf{A}_u^\ell(\text{idx}, \text{idx}), \mathbf{X}_u^{\ell+1} = \hat{\mathbf{X}}_u^\ell(\text{idx}, :), \quad (8)$$

where $\text{rank}(\cdot)$ returns the indices of the top- $K_u^{\ell+1}$ values in β^ℓ , $\mathbf{A}_u^\ell(\text{idx}, \text{idx})$ is the row-wise and column-wise indexed adjacency matrix, and $\hat{\mathbf{X}}_u^\ell(\text{idx}, :)$ is the row-wise indexed feature matrix. Then, we have $\mathbf{A}_u^{\ell+1}$ and $\mathbf{X}_u^{\ell+1}$ which are new adjacency matrix and corresponding feature matrix of a smaller graph $\mathcal{G}_u^{\ell+1}$. In this work, we empirically set $K_u^{\ell+1} = \lceil \alpha K_u^\ell \rceil$, where α is the pooling ratio, $\lceil \cdot \rceil$ is the ceiling function.

3.3 Feature Fusion

After repeating the TGAT and PGP operations L times, we can obtain multiple pooled sub-graphs with different size, i.e., $\mathcal{G}_u^2, \mathcal{G}_u^3, \dots, \mathcal{G}_u^{L+1}$. This hierarchical pooling operation extracts the coarser-grained information of review graph layer-by-layer, which can be considered as different level representations of reviews, e.g., word-level, sentence-level, and review-level. To integrate these sub-graph information, we first use the max-pooling and MLP to generate representation of each sub-graph \mathcal{G}_u^ℓ , where $2 \leq \ell \leq L+1$, as follows,

$$\hat{\mathbf{g}}_u^\ell = \max_{t=1}^{d_2} \mathbf{X}_u^\ell(:, t), \quad \mathbf{g}_u^\ell = \text{ReLU}(\hat{\mathbf{g}}_u^\ell \mathbf{W}_6^\ell + \mathbf{b}_6^\ell), \quad (9)$$

where $\mathbf{W}_6^\ell \in \mathbb{R}^{d_2 \times d_1}$ and $\mathbf{b}_6^\ell \in \mathbb{R}^{1 \times d_1}$ are the weight matrix and bias vector. \mathbf{g}_u^ℓ denotes the representation of the sub-graph generated at the ℓ -th layer. Then, we concatenate the

non-linear transform of the user embedding after a MLP layer and the representations of sub-graphs at all layers to generate the semantic representation \mathbf{p}_u of the user u as follows,

$$\begin{aligned} \mathbf{m}_u &= \text{ReLU}(\mathbf{e}_u \mathbf{W}_7 + \mathbf{b}_7), \\ \mathbf{p}_u &= \mathbf{m}_u \oplus \mathbf{g}_u^2 \oplus \mathbf{g}_u^3 \oplus \dots \oplus \mathbf{g}_u^{L+1}, \end{aligned} \quad (10)$$

where \oplus is the concatenating operation, $\mathbf{W}_7 \in \mathbb{R}^{d_1 \times d_1}$ and $\mathbf{b}_7 \in \mathbb{R}^{1 \times d_1}$ are the weight matrix and bias vector. Similarly, we can get the semantic representation of the item i as \mathbf{q}_i .

3.4 Rating Prediction

A FM layer is used to predict the user's preference on an item, considering the higher order interactions between the user and item semantic features. Specifically, we first concatenate the user and item semantic representations as $\mathbf{z} = \mathbf{p}_u \oplus \mathbf{q}_i$, and the prediction \hat{y}_{ui} of the user's preference can be defined as follows,

$$\hat{y}_{ui} = b_0 + b_u + b_i + \mathbf{z} \mathbf{w}^\top + \sum_{i=1}^{d'} \sum_{j=i+1}^{d'} \langle \mathbf{v}_i, \mathbf{v}_j \rangle z_i z_j, \quad (11)$$

where b_0, b_u and b_i are the global bias, user bias, and item bias. $\mathbf{w} \in \mathbb{R}^{1 \times d'}$ is the coefficient vector, and $d' = 2 * (L + 1) * d_1$. $\mathbf{v}_i, \mathbf{v}_j \in \mathbb{R}^{1 \times k}$ are the latent factors associated with i -th and j -th dimension of \mathbf{z} . Empirically, we set $k = d_2$ in this work. z_i is the value of the i -th dimension of \mathbf{z} . In this work, we focus on the rating prediction problem. Thus, the model parameters Θ of RGNN can be learned by solving the following optimization problem,

$$\min_{\Theta} \sum_{(u,i) \in \mathcal{D}} (y_{ui} - \hat{y}_{ui})^2 + \lambda \|\Theta\|_F^2, \quad (12)$$

where λ is the regularization parameter, \mathcal{D} denotes the set of user-item pairs used to update the model parameters. The entire framework can be effectively trained by the end-to-end back propagation algorithm.

3.5 Discussions

In this work, we build a review graph to describe all the reviews of a user/item. If a word appears in multiple review sentences, its neighbors in review graph aggregate all its neighboring words (within a word window) in these sentences. Then, TGAT is proposed to learn the word embedding, considering its global context words among all the relevant reviews. However, the CNN/RNN-based methods separately aggregate the local neighboring words of a word in different review sentences. Moreover, TGAT employs attention mechanism to automatically assign different importance weights to neighboring words. By performing multiple times of PGP operations, the proposed model can capture the long-term dependency between selected words.

Suppose the average number of selected keywords for each user/item is N_w . The time complexity of constructing the review graph is $O(\omega N_w)$. In practice, the review graph construction can be performed offline. The time complexity for RGNN to obtain the hierarchical review representation is $O(\sum_{\ell=1}^L \alpha^\ell N_w d_0^2)$, where L denotes the number of layers, α and d_0 denote the pooling ratio and word embedding

TABLE 1: The statistics of the experimental datasets.

Dataset	# Users	# Items	# Ratings	Density
Music	1,429	900	10,261	0.80%
Office	4,905	2,420	53,228	0.45%
Toys	19,412	11,924	167,597	0.07%
Beauty	22,363	12,101	198,502	0.07%
Games	24,303	10,672	231,577	0.09%
Yelp	126,084	65,786	3,519,533	0.04%

dimension. The total training time complexity of RGNN is $O(LN_w d_0^2)$. Moreover, the time complexity of using CNNs to obtain the review representation is $O(MSN_w d_0)$, where M and S denote the number of convolutional kernels and the convolutional kernel size, respectively. Thus, compared with the CNN based methods, the proposed RGNN model is not complicate or can be even more efficient.

4 EXPERIMENTS

We empirically perform extensive experiments on real datasets to demonstrate the effectiveness of the proposed RGNN method, comparing with state-of-the-art review-based recommendation methods.

4.1 Experimental Setting

4.1.1 Datasets

The experiments are performed on the Amazon review dataset [55] and Yelp dataset², which have been widely used for recommendation research. For Amazon review dataset, we choose the following 5-core review subsets for evaluation: “Musical Instruments”, “Office Products”, “Toys and Games”, “Beauty”, and “Video Games” (respectively denoted by Music, Office, Toys, Beauty, and Games). For Yelp dataset, we only keep the users and items that have at least 10 reviews for experiments. Table 1 summarizes the details of these experimental datasets.

4.1.2 Setup and Metric

For each dataset, we randomly choose 20% of the user-item review pairs (denoted by \mathcal{D}_{test}) for evaluating the model performance in the testing phase, and the remaining 80% of the review pairs (denoted by \mathcal{D}_{train}) are used in the training phase. Then, 10% of the data in \mathcal{D}_{train} are held out as the validation set to choose the model hyper-parameters, and the remaining 90% of the data in \mathcal{D}_{train} are used to update the model parameters. As this work focuses on the rating prediction task, the performances of the recommendation algorithms are typically evaluated by Mean Square Error (MSE), which has been widely used in previous studies [16], [19], [21], [22]. The definition of MSE is as follows,

$$MSE = \frac{1}{|\mathcal{D}_{test}|} \sum_{(u,i) \in \mathcal{D}_{test}} (y_{ui} - \hat{y}_{ui})^2, \quad (13)$$

where \mathcal{D}_{test} denotes the set of all tested user-item rating pairs. From the definition in Eq. (13), we can note that lower MSE values indicate better performances.

2. <https://www.yelp.com/dataset>

4.1.3 Baseline Methods

We compare RGNN with the following recommendation methods:

- **PMF** [56]: This is the probabilistic matrix factorization model, which is a classical collaborative filtering based rating prediction method.
- **SVD++** [57]: This is a classic matrix factorization method that exploits both the user’s explicit preferences on items and the influences of the user’s historical items on the target item.
- **CDL** [15]: This is a hierarchical Bayesian model that employs SDAE for learning features from the content information and collaborative filtering for modeling the rating behaviors.
- **DeepCoNN** [16]: This method contains two parallel networks, which focus on modeling the user behaviors and learning the item properties from the review data.
- **ANR** [21]: This is an aspect-based neural recommendation model that learns aspect-based representations for the user and item by an attention-based module. Moreover, the co-attention mechanism is applied to the user and item importance at the aspect-level.
- **NARRE** [17]: This method uses an attention mechanism to model the importance of reviews and a neural regression model with review-level explanations for rating prediction.
- **CARL** [22]: This is a context-aware representation learning model for rating prediction, which uses convolution operation and attention mechanism for review-based feature learning and factorization machine for modeling high-order feature interactions.
- **NRPA** [19]: This method uses a review encoder to learn the review representation, and a user/item encoder with personalized attention mechanism to learn user/item representations from reviews.
- **RMG** [24]: This method utilizes a multi-view learning framework to incorporate information from the review-content view and user-item interaction view for recommendation.
- **DAML** [18]: This method employs CNN with local and mutual attention mechanism to learn the review features and improve the interpretability of the recommendation model.

4.1.4 Parameter Settings

In this work, we implement the proposed RGNN model by Pytorch [58] and choose Adam [59] as the optimization method to learn the model parameters. The dimensionality of the semantic space d_0 , the rating space d_1 , and the hidden space of the graph representation learning layers d_2 are chosen from [4, 8, 16, 32, 48, 64, 128]. The learning rate γ is varied in [0.0005, 0.001, 0.005], and regularization parameter λ is varied in [0.001, 0.01, 0.05, 0.1, 1.0]. The batch size is chosen in [64, 128, 256]. Moreover, the number of layers L is varied in [1, 2, 3, 4], and the pooling ratio α is set to 0.5.

The hyper-parameters of baseline methods are set following original papers. In PMF and SVD++, the dimension of latent embeddings are chosen from {16, 32, 48, 64, 128}. The learning rate is varied in {0.0001, 0.001, 0.01}, and the batch size is chosen from {64, 128, 256, 512}. For CDL, we set $a=1$, $b=0.01$, $K=50$, and choose λ_u , λ_v , λ_n , λ_w

TABLE 2: The performances of different recommendation algorithms evaluated by MSE. The best results are in **bold** faces and the second best results are underlined. \blacktriangle denotes the average relative improvements achieved by RGNN over baselines on all datasets. * indicates RGNN significantly outperforms all baselines with $p < 0.05$ using student t -test on the dataset.

Dataset	PMF	SVD++	CDL	DeepCoNN	ANR	NARRE	CARL	NRPA	RMG	DAML	RGNN
Music	1.8783	0.7952	1.2987	0.7909	0.7825	0.7688	0.7632	0.7658	0.7560	<u>0.7401</u>	0.7319*
Office	0.9635	0.7213	0.8763	0.7315	0.7237	0.7266	0.7193	0.7343	0.7348	<u>0.7164</u>	0.7125*
Toys	1.6091	0.8276	1.2479	0.8073	0.7974	0.7912	0.8248	0.8100	0.8435	<u>0.7909</u>	0.7786*
Games	1.5260	1.2081	1.6002	1.1234	<u>1.1038</u>	1.1120	1.1308	1.1259	1.1496	1.1086	1.0996*
Beauty	2.7077	1.2129	1.7726	1.2210	1.2021	<u>1.1997</u>	1.2250	1.2034	1.2419	1.2175	1.1885*
Yelp	1.4217	1.2973	1.4042	1.2719	1.2708	<u>1.2675</u>	1.3199	1.2721	1.2981	1.2700	1.2645*
\blacktriangle Improvements	38.97%	4.77%	29.02%	3.16%	2.06%	1.77%	3.43%	2.57%	4.19%	1.14%	-

TABLE 3: Performances of RGNN with respect to different settings of L .

	Music	Office	Toys	Games	Beauty	Yelp
L=1	0.7572	0.7130	0.7830	1.1021	1.1964	1.2671
L=2	0.7319	0.7125	0.7786	1.1025	1.1946	1.2658
L=3	0.7557	0.7145	0.7825	1.0996	1.1903	1.2645
L=4	0.7557	0.7149	0.7802	1.1017	1.1885	1.2684

from $\{0.01, 0.1, 1, 10, 100\}$. The number of hidden factors is selected from $\{50, 100, 200, 300\}$, and dropout ratio is varied in $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. For DeepCoNN, we empirically vary the number of latent factors (*i.e.*, $|\mathbf{x}_u|$ and $|\mathbf{y}_i|$) from 5 to 100 and the number of convolutional kernels from 10 to 400. The window size is set to 3, and the learning rate is chosen from $\{0.0001, 0.001, 0.002, 0.01\}$. In ANR, the dimensionality of latent representations is chosen from $\{64, 128, 256\}$. The window size is set to 3. The number of aspects is selected from $\{2, 4, 6, 8\}$. Moreover, the dimension of aspect embeddings is chosen from $\{5, 10, 15, 25, 50\}$, and the dropout ratio is selected from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. For NARRE, the learning rate is selected from $\{0.0001, 0.001, 0.002, 0.01, 0.05\}$. The dropout ratio is chosen from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. The batch size is selected from $\{50, 100, 150\}$, and the latent factor number is varied in $\{16, 32, 64, 128\}$. The word window size is to 3. For CARL, we select the dimensionality of latent embeddings from $\{15, 50, 100, 200, 300\}$ and the batch size from $\{100, 200\}$. The learning rate is chosen from $\{0.0001, 0.001, 0.01\}$, and the dropout ratio is chosen from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. The word window size is also set to 3. In NRPA, the dimensionality of embeddings is chosen from $\{16, 32, 64, 128\}$. The number of filters in CNN and the dimension of attention vectors are chosen from $\{20, 40, 80, 100\}$. The word window size is set to 3. For RMG, the number of CNN filters is chosen from $\{50, 100, 150, 200\}$, and the dimensionality of embeddings is chosen from $\{50, 100, 150, 200\}$. The batch size is varied in $\{20, 50, 100, 200\}$, and the dropout ratio is selected from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. Empirically, the size of the word window is set to 3. For DAML, the dimensionality of latent feature vectors is chosen from $\{8, 16, 32, 64\}$, and the sliding window size is set to 3. The learning rate is selected from $\{0.00001, 0.0001, 0.001, 0.002, 0.01\}$. The dropout ratio is selected from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$.

For all evaluated methods, we employ early stopping strategy for model training, and the optimal hyperparameters are chosen based on the performances on the validation data.

4.2 Performance Comparison

Table 2 summarizes the results of different recommendation methods. We can make the following observations.

- The review-based recommendation methods significantly outperforms PMF, which demonstrates the importance of review data for improving recommendation accuracy.
- The neural methods using CNNs and attention mechanism (*e.g.*, DeepCoNN and NARRE) achieve better performances than CDL. The potential reason is that CDL employs bag-of-words model to process the review data, thus ignores the word order and local context information.
- The review-based methods based on attention mechanism (*e.g.*, NARRE, NRPA, RMG) outperform DeepCoNN, by learning the importance of words or reviews for better understanding of users' rating behaviours. DAML usually performs best among baseline methods. It utilizes mutual attention of the CNN to study the correlation between the user reviews and item reviews, thus can capture more information about the user preferences and item properties.
- RGNN consistently achieves the best performances on all datasets. The improvements over baselines are statistically significant with $p < 0.05$ using student t -test. In RGNN, the user/item reviews are described by a directed review graph with preserving the word orders, which can capture the global and non-consecutive topology and dependency information between words in reviews. Moreover, the PGP operator extracts hierarchical and coarser-grained representations of reviews while considering the user/item-specific properties, thus can further help improve the recommendation performance.

Moreover, Table 3 summarizes the performances of RGNN with respect to different settings of L . On Music, Office, Toys datasets, the best performances of RGNN are achieved by setting L to 2. On Games, Beauty, and Yelp datasets, RGNN achieves the best performances, when L is set to 3, 4, and 3, respectively. These results demonstrate the effectiveness of the hierarchical graph representation learning.

4.3 Ablation Study

To investigate the importance of each component of RGNN, we firstly study the contributions of the review graph and the factorization machine layer. We consider the following variants of RGNN for experiments:

- **RGNN-RG**: In this variant, we remove both the user and item review graphs. We use the max pooling of the

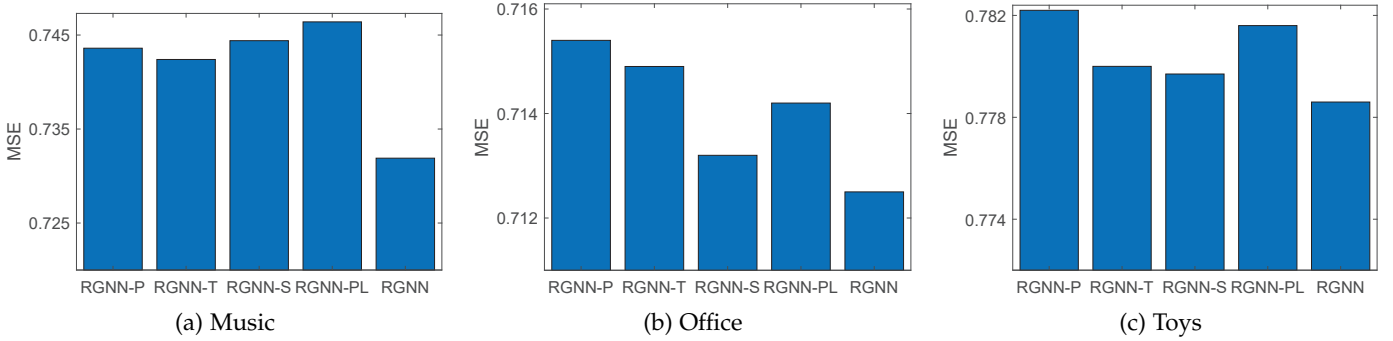


Fig. 4: Performances of RGNN-T, RGNN-P, RGNN-S, RGNN-PL, and RGNN on Music, Office, and Toys datasets.

TABLE 4: Performances of DAML, RGNN-RG, RGNN-FM, and RGNN on the Music, Office, and Toys datasets.

Dataset	DAML	RGNN-RG	RGNN-FM	RGNN
Music	0.7401	0.7705	0.7411	0.7319
Office	0.7164	0.7132	0.7115	0.7125
Toys	0.7909	0.7903	0.7817	0.7786

embeddings of the keywords extracted from S_u as the review representation of the user u , which is denoted by \mathbf{r}_u . Then, in Eq. (10), we concatenate \mathbf{m}_u with \mathbf{r}_u to form the user representation \mathbf{p}_u . Similar operation is performed on the reviews of the item i . Note that FM layer is also used for rating prediction in this variant.

- **RGNN-FM**: For this variant, we remove the FM layer and replace the rating prediction function in Eq. (11) by the product of \mathbf{p}_u and \mathbf{q}_i as $\hat{y}_{ui} = \mathbf{p}_u \mathbf{q}_i^T$.

Table 4 summarizes the performances achieved by DAML, RGNN-RG, RGNN-FM, and RGNN on the Music, Office, and Toys datasets. From Table 4, we can note that RGNN outperforms RGNN-RG by 5.01% and 1.48% on Music and Toys dataset. Compared with the best baseline method DAML that employs FM layer for prediction, RGNN-FM outperforms DAML on the Office and Toys datasets, achieves comparable results with DAML on the Music dataset. These two observations demonstrate the effectiveness of the proposed review graph based representation learning strategy in capturing the users’ preferences on items. Moreover, RGNN outperforms RGNN-FM on Music and Toys datasets. This indicates the FM layer can help improve the rating prediction accuracy. As shown in Table 4, RGNN-FM achieves better results than RGNN-RG on three datasets. This demonstrates that the main improvements achieved by RGNN are caused by the proposed hierarchical review graph representation learning strategy.

Moreover, we perform experiments to analyze the contributions of different components in the hierarchical graph representation learning procedure. Specifically, we evaluate the performances of the following RGNN variants:

- **RGNN-T**: RGNN without considering the word orders (*i.e.*, removing the edge type embedding in Eq. (2)).
- **RGNN-P**: RGNN without considering personalized properties of users/items (*i.e.*, using the same θ^ℓ for all users/items in Eq. (5)).
- **RGNN-S**: RGNN without considering the feature diversity of a node’s neighborhood in estimating the node importance (*i.e.*, removing the second part in Eq. (5)).

- **RGNN-PL**: RGNN without considering the graph pooling (*i.e.*, removing the PGP operators in Figure 3).

The experimental results on Music, Office, and Toys datasets are summarized in Figure 4. As shown in Figure 4, RGNN outperforms RGNN-T on all three datasets. This indicates that the word order is important in extracting semantic information from reviews for recommendation. Moreover, RGNN achieves better performances than RGNN-P, RGNN-S and RGNN-PL in terms of MSE. This result demonstrates that the proposed PGP operator can help improve the recommendation accuracy via capturing more informative review words for a user/item. When estimating the importance of a node (*i.e.*, word), the proposed PGP operator considers the user/item-specific properties, as well as the feature diversity of a node’s neighborhood in the review graph. From Figure 4, we can also note that GNN-PL achieves better MSE than RGNN-P on Office and Toys datasets. RGNN-P adopts the same pooling parameter for all users/items. It tends to select a set of common words to form the review representations for all the users/items. The specific properties of a user/item in the review data may not be exploited. Therefore, RGNN-P may perform poorer than RGNN-PL that employs max pooling on the embeddings of all nodes in the review graph of each individual user/item.

4.4 Impacts of Edge Types

In addition, we also conduct experiments to analyze the impacts of different edge types. Specifically, we study the performances of RGNN with the following settings:

- **RGNN-f**: RGNN without considering the forward edge type e_f .
- **RGNN-b**: RGNN without considering the backward edge type e_b .
- **RGNN-s**: RGNN without considering the self-loop edge type.

Figure 5 summarizes the MSE values achieved by RGNN with different settings of edge types on Music, Office, and Toys datasets. As shown in Figure 5, RGNN-b outperforms RGNN-f on all three datasets. This indicates that the forward type of edges are more important than the backward type of edges in capturing users’ rating behaviors for recommendation. Moreover, RGNN-s achieves the worst MSE values on Music and Toys datasets. This indicates the self-loop type of edges play an important role in the graph representation learning procedure. This observation is consistent with the

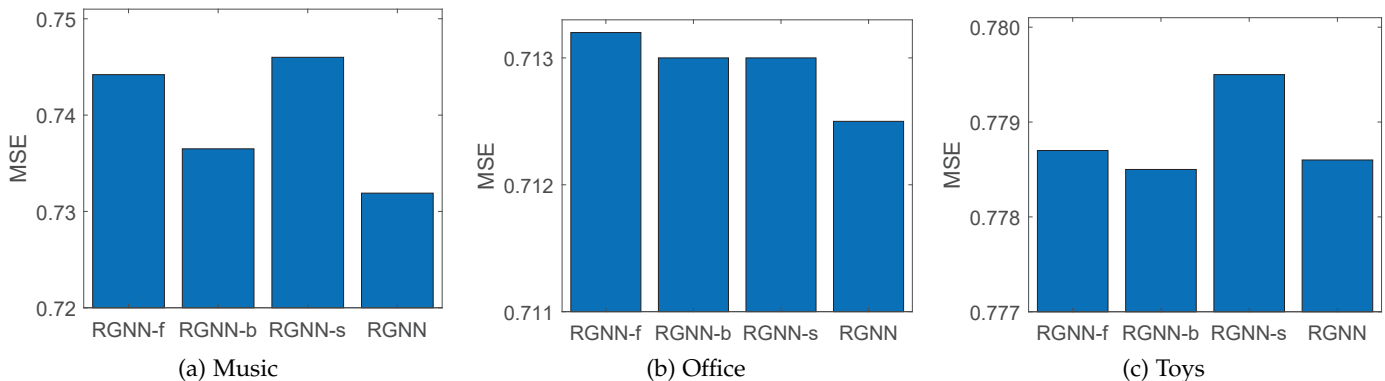


Fig. 5: The performances of RGNN with different settings of edge types on Music, Office, and Toys datasets.

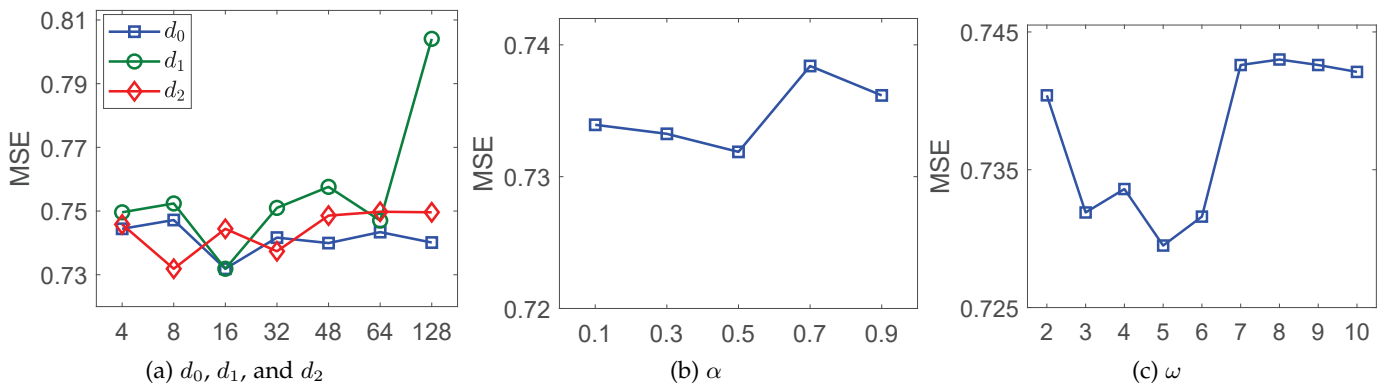


Fig. 6: The performances of RGNN with respect to different settings of d_0 , d_1 , d_2 , α , and ω on Music dataset.

definition of the attention mechanism in Eq. (3). By removing the self-loop type of edges, for each node x_h in the graph at the ℓ -th layer, it is not included in its first-hop neighbors \mathcal{N}_h^ℓ in the graph. According to Eq. (3), the new embedding of x_h (i.e., \bar{x}_h^ℓ) is updated by only considering the input embeddings of other neighboring nodes at the ℓ -th layer, without considering the input embedding of x_h at the ℓ -th layer. In addition, RGNN achieves the best results on Music and Office datasets. This indicates that the combination of all three types of edges usually provides a better description about the review data, thus achieves more accurate rating prediction results.

4.5 Parameter Sensitivity Study

In this section, we perform experiments to evaluate the sensitivity to hyper-parameters for RGNN. The dimensionality of the semantic space d_0 , the dimensionality of the rating space d_1 , and the dimensionality of the hidden space of the graph representation learning layers d_2 , are varied in $\{4, 8, 16, 32, 48, 64, 128\}$. The graph pooling ratio α is adjusted in $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. The size of the word window ω is chosen from $\{2, 3, 4, 5, 6, 7, 8, 9, 10\}$. The results on Music dataset are summarized in Figure 6. We can note that the best MSE can be achieved by setting d_0 , d_1 to 16, and d_2 to 8. Larger dimensionality of the latent semantic space and rating space cause more computation cost, however may not help improve the rating prediction accuracy. For the pooling ratio, we can notice that the best setting for α

on Music dataset is 0.5. Moreover, RGNN usually achieves better performances by setting $\alpha < 0.5$, compared by setting $\alpha > 0.5$. One potential reason is that when $\alpha > 0.5$ more nodes would be kept by the pooling operation, thus the “noise nodes” are more likely to be involved in learning the hierarchical review graph representations. In addition, as shown in Figure 6 (c), RGNN achieves the best MSE value when ω is set to 5 on Music dataset. Better performances can be usually achieved when ω is set in the range between 3 and 6.

4.6 Discussions about Pooled Graphs

The proposed RGNN model employs graph pooling to learn the review representations. However, the pooling operation may produce isolated subgraphs. For example, on music dataset, after two layers of pooling operations, there are about 10 isolated subgraphs (3 of which are isolated nodes), on average. Figure 7 shows an example about the graph connectivity after two layers of pooling operation. As there are different edge types in the review graph, it is very hard for us to build the connections between different subgraphs. Thus, in this work, we do not perform specific operations to connect these subgraphs. At the ℓ -th layer, in the TGAT operator, the embedding of a node is updated considering its neighbouring nodes in the subgraph it belongs to. Then, in the PGP operator, we score the nodes in all subgraphs in this layer for personalized pooling. In the future work, we

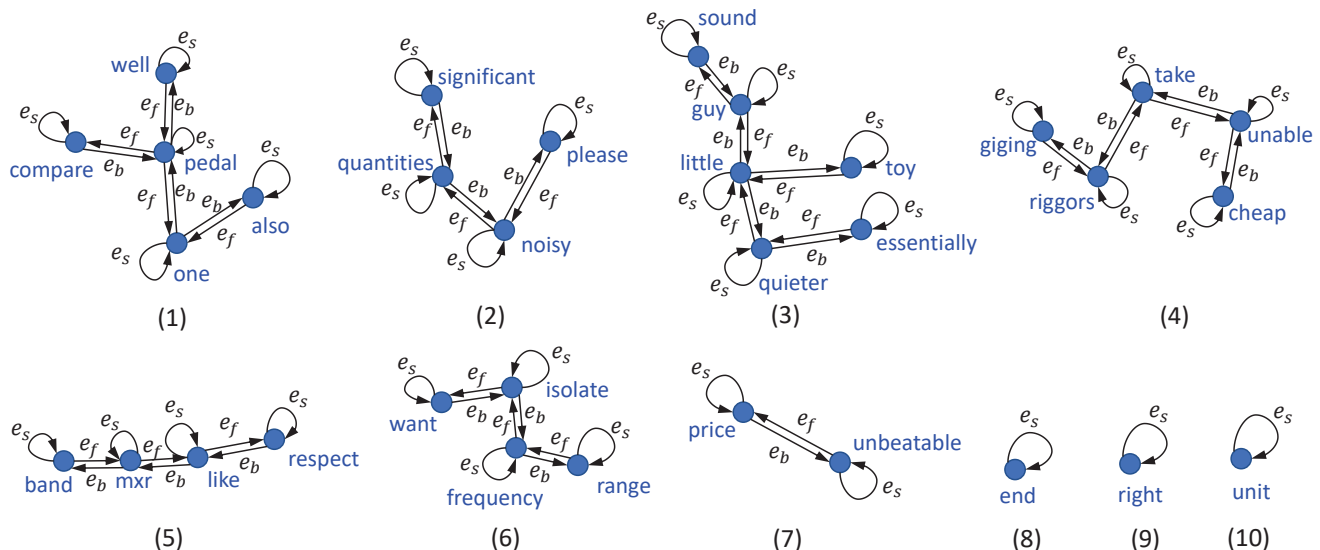


Fig. 7: An example about the graph structure after two layers of graph pooling. In this example, there are 10 isolated subgraphs after graph pooling.

would like to study how to maintain the connectivity of the review graph after the pooling operation.

4.7 Case Study

In addition, we also conduct a case study to further explore whether RGNN can select informative words in the graph pooling layer. We sample a user-item pair in test data and extract the review documents of the user and item. Figure 8 visualizes the heat maps of these review documents. We highlight the words extracted from the first PGP layer of RGNN. The words with dark yellow background have larger importance scores than the words with yellow background (refer to Eq. (5)). As shown in Figure 8, the words “playing”, “sound”, “money”, and “long” are highlighted in user reviews. These words indicate that this user may be more interested in the quality and price of an item. Moreover, in the item reviews, the words “high”, “quality”, “solid”, “sound”, and “cheap” which describe the item properties are selected in the first pooling layer. In the review given by the user on the item, from the words “price”, “quality”, “natural”, and “sound”, we can argue that the user’s high rating on this item is indeed, because this item meets the concerns of the user. Hence, this case again demonstrates that the PGP layer can select informative words from the review graphs, thus can help improve the recommendation accuracy.

5 CONCLUSION AND FUTURE WORK

In this work, we propose a novel review-based recommendation model, namely Review Graph Neural Network (RGNN), which incorporates the advantages of graph attention networks and hierarchical graph pooling in understanding users’ review behaviors for recommendation. In RGNN, the reviews of each user/item are described by a directed graph, which simultaneously considers the co-occurrence of words and the word orders. RGNN employs a type-aware graph attention network and a personalized

graph pooling operator to learn hierarchical semantic representations for users and items from the review data. A factorization machine layer is then developed to predict the user’s rating on an item, considering the interactions between the user and item semantic features. The experimental results on two real-world datasets show that RGNN consistently outperforms state-of-the-art review-based recommendation methods, in terms of MSE. For future work, we would like to study how to maintain the connectivity of the review graph after the graph pooling operation. Moreover, we are also interested in extending RGNN to solve the top- N item recommendation problem.

6 ACKNOWLEDGMENTS

This research is supported, in part, by Alibaba Group through Alibaba Innovative Research (AIR) Program and Alibaba-NTU Singapore Joint Research Institute (JRI), Nanyang Technological University, Singapore. This research is also supported, in part, by the National Research Foundation, Prime Minister’s Office, Singapore under its AI Singapore Programme (AISG Award No: AISG-GC-2019-003) and under its NRF Investigatorship Programme (NRFI Award No. NRF-NRFI05-2019-0002). Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not reflect the views of National Research Foundation, Singapore. This research is also supported, in part, by the National Natural Science Foundation of China (No. 61672481), and the Youth Innovation Promotion Association CAS (No. 2018495).

REFERENCES

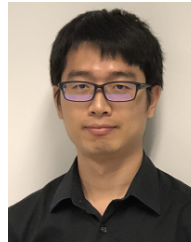
- [1] Y. Shi, M. Larson, and A. Hanjalic, “Collaborative filtering beyond the user-item matrix: A survey of the state of the art and future challenges,” *CSUR*, vol. 47, no. 1, p. 3, 2014.
- [2] S. Zhang, L. Yao, A. Sun, and Y. Tay, “Deep learning based recommender system: A survey and new perspectives,” *CSUR*, vol. 52, no. 1, p. 5, 2019.

Review document of user (ID: 701):	Review document of item (ID: 594):
<p>"I absolutely love these strings. I'm used to playing higher tension strings and these fit the bill perfectly. They have a great full sound and after breaking them in after a couple weeks of playing they give a really good classic P bass sound. Its not for people who want the modern brighter sound but for the classic 50s 60s and 70s sound these are great!", 'Only issue they break easily. But they last a long time and sound great. Well worth the money. Crisp and fresh tone for a long time.' "I love these strings they are a little more mellow than most mando strings so they sit in the mix and let the guitars cut though a little more. I'm hoping they last a while!", I love these cables. I ordered two about a year and change ago and they are my main cables for my bass and guitar rigs. They are still perfect and they get used for many hours each day'</p>	<p>"I've used both the ES-57 and ES-58 for several years now. They are a fantastic mic, great quality and construction...I would probably suggest a higher end if you were studio recording but I set up for live sound in and outdoors and these work great..." "This mic has a very solid construction and looks and feels like a SM57...I think it will do just fine." "Wow, I will never buy another SM-57 again...I am so psyched that I found this microphone, it's perfect. But seriously, the sound quality has never once disappointed me...But everyone else has talked about sound, so I want to talk about build quality.It's solid, has some weight to it, and the plastic at the top feels exactly right." "I run the sound for local bands and I use SM's because they can take abuse and still sound great...A little less bass but that doesn't mean it is a thin sound - it certainly is not. I can get it to sound the same as the SM if I EQ it slightly but I like the natural sound anyway. Also, the signal is a bit hotter but I really like the tone for both vocals and for instruments. Value for money it can't be beat. If you are on a budget, get one....", "I bought this to replace my Blue Snowball because I've reached that point where I require something that sounds a little bit more professional than what I was using...My voice sounds great with this thing...With this, my voice is strong and clear...The materials feel high quality. This feels like a professional microphone...All I can tell you is that this thing sounds freaking amazing... So if you're looking for a vocal/podcast mic, this is a great, cheap option. Are there better options? Sure, but not for this price. Not even close. Recommended wholeheartedly."</p>
<p>Review given by user 701 on item 594:</p>	
<p>This is the equivalent to the Shure SM-57 except for a quarter of the price. The quality is almost the same and in terms of recording they are pretty much the same as well. There allow a more natural sound than the SM57 I believe and are thus a better product!</p>	
<p>Rating given by user 701 on item 594:</p>	
<p>Real rating value: 5.0; Predicted rating value: 5.064</p>	

Fig. 8: Highlighted words according to importance scores in the review documents of a user-item pair. The words with dark yellow background have larger importance scores than the words with yellow background. The words with red background are manually selected to show the user's opinions on the item.

- [3] L. Hu, A. Sun, and Y. Liu, "Your neighbors affect your ratings: on geographical neighborhood influence to rating prediction," in *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, 2014, pp. 345–354.
- [4] X. Liu, Y. Liu, K. Aberer, and C. Miao, "Personalized point-of-interest recommendation by mining users' preference transition," in *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, 2013, pp. 733–738.
- [5] Q. Wu, Y. Liu, C. Miao, B. Zhao, Y. Zhao, and L. Guan, "Pd-gan: adversarial learning for personalized diversity-promoting recommendation," in *IJCAI'19*. AAAI Press, 2019, pp. 3870–3876.
- [6] Z. Sun, J. Yang, J. Zhang, A. Bozzon, L.-K. Huang, and C. Xu, "Recurrent knowledge graph embedding for effective recommendation," in *RecSys'18*. ACM, 2018, pp. 297–305.
- [7] L. Chen, G. Chen, and F. Wang, "Recommender systems based on user reviews: the state of the art," *User Modeling and User-Adapted Interaction*, vol. 25, no. 2, pp. 99–154, 2015.
- [8] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *JMLR*, vol. 3, no. Jan, pp. 993–1022, 2003.
- [9] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [10] J. McAuley and J. Leskovec, "Hidden factors and hidden topics: understanding rating dimensions with review text," in *RecSys'13*, 2013, pp. 165–172.
- [11] Y. Bao, H. Fang, and J. Zhang, "Topicmf: Simultaneously exploiting ratings and reviews for recommendation," in *AAAI'14*, 2014.
- [12] G. Ling, M. R. Lyu, and I. King, "Ratings meet reviews, a combined approach to recommend," in *RecSys'14*, 2014, pp. 105–112.
- [13] Y. Tan, M. Zhang, Y. Liu, and S. Ma, "Rating-boosted latent topics: Understanding users and items with ratings and reviews." in *IJCAI'16*, vol. 16, 2016, pp. 2640–2646.
- [14] W. Zhang, Q. Yuan, J. Han, and J. Wang, "Collaborative multi-level embedding learning from reviews for rating prediction." in *IJCAI'16*, 2016, pp. 2986–2992.
- [15] H. Wang, N. Wang, and D.-Y. Yeung, "Collaborative deep learning for recommender systems," in *KDD'15*. ACM, 2015, pp. 1235–1244.
- [16] L. Zheng, V. Noroozi, and P. S. Yu, "Joint deep modeling of users and items using reviews for recommendation," in *WSDM'17*. ACM, 2017, pp. 425–434.
- [17] C. Chen, M. Zhang, Y. Liu, and S. Ma, "Neural attentional rating regression with review-level explanations," in *WWW'18*, 2018, pp. 1583–1592.
- [18] D. Liu, J. Li, B. Du, J. Chang, and R. Gao, "Daml: Dual attention mutual learning between ratings and reviews for item recommendation," in *KDD'19*. ACM, 2019, pp. 344–352.
- [19] H. Liu, F. Wu, W. Wang, X. Wang, P. Jiao, C. Wu, and X. Xie, "Nrpa: Neural recommendation with personalized attention," *SIGIR'19*, 2019.
- [20] K. Bauman, B. Liu, and A. Tuzhilin, "Aspect based recommendations: Recommending items with the most valuable aspects based on user reviews," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017, pp. 717–725.
- [21] J. Y. Chin, K. Zhao, S. Joty, and G. Cong, "Anr: Aspect-based neural recommender," in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, 2018, pp. 147–156.
- [22] L. Wu, C. Quan, C. Li, Q. Wang, B. Zheng, and X. Luo, "A context-aware user-item representation learning for item recommendation," *TOIS*, vol. 37, no. 2, p. 22, 2019.
- [23] Y. Lu, R. Dong, and B. Smyth, "Coevolutionary recommendation model: Mutual learning between ratings and reviews," in *WWW'18*, 2018, pp. 773–782.
- [24] C. Wu, F. Wu, T. Qi, S. Ge, Y. Huang, and X. Xie, "Reviews meet graphs: Enhancing user and item representations for recom-

- mentation with hierarchical attentive graph neural network," in *EMNLP-IJCNLP 2019*, 2019, pp. 4886–4895.
- [25] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [26] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE Transactions on Neural Networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [27] C. Gallicchio and A. Micheli, "Graph echo state networks," in *The 2010 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2010, pp. 1–8.
- [28] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, "Gated graph sequence neural networks," in *ICLR*, 2015.
- [29] J. Bruna, W. Zaremba, A. Szlam, and Y. Lecun, "Spectral networks and locally connected networks on graphs," in *ICLR*, 2014.
- [30] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NIPS*, 2016, pp. 3844–3852.
- [31] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [32] A. Micheli, "Neural network for graphs: a contextual constructive approach," *IEEE Transactions on Neural Networks*, vol. 20, no. 3, pp. 498–511, 2009.
- [33] J. Atwood and D. Towsley, "Diffusion-convolutional neural networks," in *Advances in Neural Information Processing Systems 29*, 2016, pp. 1993–2001.
- [34] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems 30*, 2017, pp. 1024–1034.
- [35] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [36] J. Zhang, X. Shi, J. Xie, H. Ma, I. King, and D. Y. Yeung, "Gaan: Gated attention networks for learning on large and spatiotemporal graphs," in *UAI*, 2018.
- [37] J. L. Keyulu Xu, Weihua Hu and S. Jegelka, "How powerful are graph neural networks?" in *ICLR*, 2019.
- [38] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Advances in neural information processing systems*, 2016, pp. 3844–3852.
- [39] M. Fey, J. Eric Lenssen, F. Weichert, and H. Müller, "Splinecnn: Fast geometric deep learning with continuous b-spline kernels," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 869–877.
- [40] F. M. Bianchi, D. Grattarola, and C. Alippi, "Mincut pooling in graph neural networks," *arXiv preprint arXiv:1907.00481*, 2019.
- [41] Y. Ma, S. Wang, C. C. Aggarwal, and J. Tang, "Graph convolutional networks with eigenpooling," *arXiv preprint arXiv:1904.13107*, 2019.
- [42] O. Vinyals, S. Bengio, and M. Kudlur, "Order matters: Sequence to sequence for sets," *arXiv preprint arXiv:1511.06391*, 2015.
- [43] M. Zhang, Z. Cui, M. Neumann, and Y. Chen, "An end-to-end deep learning architecture for graph classification," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [44] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," in *Advances in Neural Information Processing Systems*, 2018, pp. 4800–4810.
- [45] H. Gao and S. Ji, "Graph u-nets," *arXiv preprint arXiv:1905.05178*, 2019.
- [46] H. Gao, Y. Chen, and S. Ji, "Learning graph pooling and hybrid convolutional operations for text representations," in *WWW'19*. ACM, 2019, pp. 2743–2749.
- [47] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," *arXiv preprint arXiv:1904.08082*, 2019.
- [48] F. Diehl, "Edge contraction pooling for graph neural networks," *arXiv preprint arXiv:1905.10990*, 2019.
- [49] E. Ranjan, S. Sanyal, and P. P. Talukdar, "Asap: Adaptive structure aware pooling for learning hierarchical graph representations," *arXiv preprint arXiv:1911.07979*, 2019.
- [50] Z. Zhang, J. Bu, M. Ester, J. Zhang, C. Yao, Z. Yu, and C. Wang, "Hierarchical graph pooling with structure learning," *arXiv preprint arXiv:1911.05954*, 2019.
- [51] S. Rendle, "Factorization machines," in *ICDM'10*, 2010, pp. 995–1000.
- [52] F. D. Malliaros and M. Vazirgiannis, "Graph-based text representations: Boosting text mining, nlp and information retrieval with graphs," in *EMNLP'17: Tutorial Abstracts*, 2017.
- [53] R. Mihalcea and P. Tarau, "TextRank: Bringing order into text," in *EMNLP'04*, Jul. 2004, pp. 404–411.
- [54] D. Shen, G. Wang, W. Wang, M. R. Min, Q. Su, Y. Zhang, C. Li, R. Henao, and L. Carin, "Baseline needs more love: On simple word-embedding-based models and associated pooling mechanisms," in *ACL'18*, 2018, pp. 440–450.
- [55] R. He and J. McAuley, "Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering," in *WWW'16*, 2016, pp. 507–517.
- [56] A. Mnih and R. R. Salakhutdinov, "Probabilistic matrix factorization," in *NIPS'08*, 2008, pp. 1257–1264.
- [57] Y. Koren, "Factorization meets the neighborhood: a multifaceted collaborative filtering model," in *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2008, pp. 426–434.
- [58] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32*, 2019, pp. 8024–8035.
- [59] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.



Yong Liu is a Research Scientist at Alibaba-NTU Singapore Joint Research Institute and Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly (LILY), Nanyang Technological University, Singapore. Before that, he was a Data Scientist at NTUC Enterprise, Singapore from November 2017 to July 2018, and a Research Scientist at Data Analytics Department, Institute for Infocomm Research (I2R), A*STAR, Singapore from November 2015 to October 2017. He received his Ph.D. from the School of Computer Science and Engineering at Nanyang Technological University in 2016 and B.S. from the Department of Electronic Science and Technology at University of Science and Technology of China in 2008. His research areas include various topics in machine learning and data mining. His research papers appear in leading international conferences and journals. He has been invited as a PC member of major conferences such as KDD, SIGIR, IJCAI, AAAI, CIKM, ICDM, and reviewer for IEEE/ACM transactions.



Susen Yang received the B.S. degree from the China University of Petroleum, Tsingtao, China, in 2018. He is currently pursuing the M.S. degree from the University of Science and Technology of China, Hefei, China. His research interests include deep learning, recommendation system, and graph neural networks.



Yinan Zhang received the B.S. degree from Taishan College, Shandong University, China, in 2019. She is currently a Ph.D. student in the school of Computer Science and Engineering at Nanyang Technological University, Singapore. Her research interests include deep learning and recommendation systems.



Juyong Zhang is an associate professor in the School of Mathematical Sciences at University of Science and Technology of China. He received the BS degree from the University of Science and Technology of China in 2006, and the PhD degree from Nanyang Technological University, Singapore. His research interests include computer graphics, computer vision, and numerical optimization. He is an associate editor of The Visual Computer.



Chunyan Miao is a President's Chair Professor and the Chair of the School of Computer Science and Engineering (SCSE) at NTU Singapore. She received her PhD degree in Computer Engineering from NTU and was an NSERC Postdoctoral Fellow at Simon Fraser University (SFU), Canada. She was a founding faculty member of the Centre for Digital Media established by The University of British Columbia (UBC) and SFU. She was also a Tan Chin Tuan Engineering Fellow at Harvard and MIT. Dr. Miao has received

over 20 Best Paper/innovation awards in Artificial intelligence (AI) and real world AI applications for her impactful research in health, ageing, education and smart services. She is a recipient of the prestigious NRF Investigatorship Award 2018. She also holds major research funding including MOH National Innovation Challenge (NIC) on Ageing award 2018 and NRF AI.SG Health Grand Challenge Award 2019. She is the Founding Director of the Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly (LILY), Singapore's first centre focusing on AI empowered solutions to population aging challenges. She is also the Founding Director of the Alibaba-NTU Singapore Joint Research Institute (JRI), Alibaba's first and largest JRI outside China. She is an Editor/Associate Editor of leading international journals including IJIT, IEEE Big Data, IEEE IoT, IEEE Access and IEEE Service Computing and has served as Chair/TPC member of international conferences such as IEEE ICA, ICAA, ACM KDD. She serves on various national committees, including the MOH City for All Ages and Health Tech, the IMDA TechSkills Accelerator (TeSA) and is the Chair of the SCS AI Ethics Review Committee. She was awarded a Public Administration Medal (Bronze) from the President of Singapore in 2016.



Zaiqing Nie is currently a Named Professor of Tsinghua University and a Principal Researcher of the Institute for AI Industry Research. Previously he was a Principal Engineer of Alibaba AI Labs. He received a Ph.D. in Computer Science from Arizona State University in 2004, a Master of Engineering degree in Computer Applications from Tsinghua University in 1998, and a Bachelor of Engineering degree in Computer Science and Technology from Tsinghua University from in 1996. Before joining Alibaba, he was a Principal Researcher and led the Big Data Mining group at Microsoft Research. His research interests include natural language understanding, data mining, and machine learning.

Principal Researcher and led the Big Data Mining group at Microsoft Research. His research interests include natural language understanding, data mining, and machine learning.