

Knowledge Graph Convolutional Networks for Recommender Systems

Hongwei Wang*
Shanghai Jiao Tong University
Shanghai, China
wanghongwei55@gmail.com

Miao Zhao
The Hong Kong Polytechnic
University, Hong Kong, China
csmiaozhao@comp.polyu.edu.hk

Xing Xie
Microsoft Research Asia
Beijing, China
xingx@microsoft.com

Wenjie Li
The Hong Kong Polytechnic
University, Hong Kong, China
cswjli@comp.polyu.edu.hk

Minyi Guo†
Shanghai Jiao Tong University
Shanghai, China
guo-my@cs.sjtu.edu.cn

ABSTRACT

To alleviate sparsity and cold start problem of collaborative filtering based recommender systems, researchers and engineers usually collect attributes of users and items, and design delicate algorithms to exploit these additional information. In general, the attributes are not isolated but connected with each other, which forms a knowledge graph (KG). In this paper, we propose Knowledge Graph Convolutional Networks (KGCN), an end-to-end framework that captures inter-item relatedness effectively by mining their associated attributes on the KG. To automatically discover both high-order structure information and semantic information of the KG, we sample from the neighbors for each entity in the KG as their receptive field, then combine neighborhood information with bias when calculating the representation of a given entity. The receptive field can be extended to multiple hops away to model high-order proximity information and capture users' potential long-distance interests. Moreover, we implement the proposed KGCN in a minibatch fashion, which enables our model to operate on large datasets and KGs. We apply the proposed model to three datasets about movie, book, and music recommendation, and experiment results demonstrate that our approach outperforms strong recommender baselines.

KEYWORDS

Recommender systems; Knowledge graph; Graph convolutional networks

ACM Reference Format:

Hongwei Wang, Miao Zhao, Xing Xie, Wenjie Li, and Minyi Guo. 2019. Knowledge Graph Convolutional Networks for Recommender Systems. In *Proceedings of the 2019 World Wide Web Conference (WWW '19)*, May 13–17, 2019, San Francisco, CA, USA. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3308558.3313417>

*This work is done during the author's internship at Microsoft Research Asia.

†Minyi Guo is the corresponding author. This work was partially sponsored by the National Basic Research 973 Program of China (2015CB352403) and National Natural Science Foundation of China (61272291).

1 INTRODUCTION

With the advance of Internet technology, people can access a vast amount of online content, such as news [25], movies [5], and commodities [26]. A notorious problem with online platforms is that the volume of items can be overwhelming to users. To alleviate the impact of information overloading, recommender systems (RS) is proposed to search for and recommend a small set of items to meet users' personalized interests.

A traditional recommendation technique is collaborative filtering (CF), which assigns users and items ID-based representation vectors, then models their interactions by specific operation such as inner product [16] or neural networks [8]. However, CF-based methods usually suffer from sparsity of user-item interactions and the cold start problem. To address these limitations, researchers usually turn to feature-rich scenarios, where attributes of users and items are used to compensate for the sparsity and improve the performance of recommendation [3, 17].

A few recent studies [9, 18, 19, 22–24] have gone a step further than simply using attributes: They point out that attributes are not isolated but linked up with each other, which forms a *knowledge graph* (KG). Typically, a KG is a directed heterogeneous graph in which nodes correspond to *entities* (items or item attributes) and edges correspond to *relations*. Compared with KG-free methods, incorporating KG into recommendation benefits the results in three ways [18]: (1) The rich semantic relatedness among items in a KG can help explore their latent connections and improve the *precision* of results; (2) The various types of relations in a KG are helpful for extending a user's interests reasonably and increasing the diversity of recommended items; (3) KG connects a user's historically-liked and recommended items, thereby bringing explainability to recommender systems.

Despite the above benefits, utilizing KG in RS is rather challenging due to its high dimensionality and heterogeneity. One feasible way is to preprocess the KG by *knowledge graph embedding* (KGE) methods [20], which map entities and relations to low-dimensional representation vectors [9, 19, 23]. However, commonly-used KGE methods focus on modeling rigorous semantic relatedness (e.g., TransE [1] and TransR [12] assume $head + relation = tail$), which are more suitable for in-graph applications such as KG completion and link prediction rather than recommendation. A more natural and intuitive way is to design a graph algorithm directly to

exploit the KG structure [18, 22, 24]. For example, PER [22] and FMG [24] treat KG as a heterogeneous information network, and extract meta-path/meta-graph based latent features to represent the connectivity between users and items along different types of relation paths/graphs. However, PER and FMG rely heavily on manually designed meta-paths or meta-graphs, which are hardly to be optimal in reality. RippleNet [18] is a memory-network-like model that propagates users’ potential preferences in the KG and explores their hierarchical interests. But note that the **importance of relations is weakly characterized in RippleNet**, because the embedding matrix of a relation R can hardly be trained to capture the sense of importance in the quadratic form $\mathbf{v}^\top \mathbf{R} \mathbf{h}$ (\mathbf{v} and \mathbf{h} are embedding vectors of two entities). In addition, the size of ripple set may go unpredictably with the increase of the size of KG, which incurs heavy computation and storage overhead.

In this paper, we investigate the problem of KG-aware recommendation. Our design objective is to automatically capture both high-order structure and semantic information in the KG. Inspired by graph convolutional networks (GCN)¹ that try to generalize convolution to the graph domain, we propose Knowledge Graph Convolutional Networks (KGCN) for recommender systems. The key idea of KGCN is to aggregate and incorporate neighborhood information with bias when calculating the representation of a given entity in the KG. Such a design has two advantages: (1) Through the operation of **neighborhood aggregation**, the *local proximity structure* is successfully captured and stored in each entity. (2) Neighbors are **weighted by scores dependent on the connecting relation and specific user**, which characterizes both the *semantic information of KG* and *users’ personalized interests in relations*. Note that the size of an entity’s neighbors varies and may be prohibitively large in the worst case. Therefore, we sample a fixed-size neighborhood of each node as the receptive field, which makes the cost of KGCN predictable. The definition of neighborhood for a given entity can also be extended hierarchically to multiple hops away to model *high-order* entity dependencies and capture users’ potential *long-distance* interests.

Empirically, we apply KGCN to three datasets: MovieLens-20M (movie), Book-Crossing (book), and Last.FM (music). The experiment results show that KGCN achieves average AUC gains of 4.4%, 8.1%, and 6.2% in movie, book, and music recommendations, respectively, compared with state-of-the-art baselines for recommendation.

Our contribution in this paper are summarized as follows:

- We propose knowledge graph convolutional networks, an end-to-end framework that explores users’ preferences on the knowledge graph for recommender systems. By extending the receptive field of each entity in the KG, KGCN is able to capture users’ high-order personalized interests.
- We conduct experiments on three real-world recommendation scenarios. The results demonstrate the efficacy of KGCN-LS over state-of-the-art baselines.
- We release the code of KGCN and datasets (knowledge graphs) to researchers for validating the reported results and conducting further research. The code and the data are available at <https://github.com/hwwang55/KGCN>.

¹We will revisit GCN in related work.

2 RELATED WORK

Our method is conceptually inspired by GCN. In general, GCN can be categorized as spectral methods and non-spectral methods. Spectral methods represent graphs and perform convolution in the spectral space. For example, Bruna et al. [2] define the convolution in Fourier domain and calculates the eigendecomposition of the graph Laplacian, Defferrard et al. [4] approximate the convolutional filters by Chebyshev expansion of the graph Laplacian, and Kipf et al. [10] propose a convolutional architecture via a localized first-order approximation of spectral graph convolutions. In contrast, non-spectral methods operate on the original graph directly and define convolution for groups of nodes. To handle the neighborhoods with varying size and maintain the weight sharing property of CNN, researchers propose learning a weight matrix for each node degree [6], extracting locally connected regions from graphs [13], or sampling a fixed-size set of neighbors as the support size [7]. Our work can be seen as a non-spectral method for a special type of graphs (i.e., knowledge graph).

Our method also connects to PinSage [21] and GAT [15]. But note that both PinSage and GAT are designed for homogeneous graphs. The major difference between our work and the literature is that we offer a new perspective for recommender systems with the assistance of a heterogeneous knowledge graph.

3 KNOWLEDGE GRAPH CONVOLUTIONAL NETWORKS

In this section, we introduce the proposed KGCN model. We first formulate the knowledge-graph-aware recommendation problem. Then we present the design of a single layer of KGCN. At last, we introduce the complete learning algorithm for KGCN, as well as its minibatch implementation.

3.1 Problem Formulation

We formulate the knowledge-graph-aware recommendation problem as follows. In a typical recommendation scenario, we have a set of M users $\mathcal{U} = \{u_1, u_2, \dots, u_M\}$ and a set of N items $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$. The user-item interaction matrix $\mathbf{Y} \in \mathbb{R}^{M \times N}$ is defined according to users’ implicit feedback, where $y_{uv} = 1$ indicates that user u engages with item v , such as clicking, browsing, or purchasing; otherwise $y_{uv} = 0$. Additionally, we also have a knowledge graph \mathcal{G} , which is comprised of entity-relation-entity triples (h, r, t) . Here $h \in \mathcal{E}$, $r \in \mathcal{R}$, and $t \in \mathcal{E}$ denote the head, relation, and tail of a knowledge triple, \mathcal{E} and \mathcal{R} are the set of entities and relations in the knowledge graph, respectively. For example, the triple $(A \text{ Song of Ice and Fire}, \text{book}, \text{book.author}, \text{George Martin})$ states the fact that George Martin writes the book “A Song of Ice and Fire”. In many recommendation scenarios, an item $v \in \mathcal{V}$ corresponds to one entity $e \in \mathcal{E}$. For example, in book recommendation, the item “A Song of Ice and Fire” also appears in the knowledge graph as an entity with the same name.

Given the user-item interaction matrix \mathbf{Y} as well as the knowledge graph \mathcal{G} , we aim to predict whether user u has potential interest in item v with which he has had no interaction before. Our goal is to learn a prediction function $\hat{y}_{uv} = \mathcal{F}(u, v | \Theta, \mathbf{Y}, \mathcal{G})$, where \hat{y}_{uv} denotes the probability that user u will engage with item v , and Θ denotes the model parameters of function \mathcal{F} .

3.2 KGCN Layer

emmm 这就是原版的GAT吧。。。

KGCN is proposed to capture high-order structural proximity among entities in a knowledge graph. We start by describing a single KGCN layer in this subsection. Consider a candidate pair of user u and item (entity) v . We use $\mathcal{N}(v)$ to denote the set of entities directly connected to v ,² and r_{e_i, e_j} to denote the relation between entity e_i and e_j . We also use a function $g : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ (e.g., inner product) to compute the score between a user and a relation:

$$\pi_r^u = g(\mathbf{u}, \mathbf{r}), \quad (1)$$

where $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{r} \in \mathbb{R}^d$ are the representations of user u and relation r , respectively, d is the dimension of representations. In general, π_r^u characterizes the importance of relation r to user u . For example, a user may have more potential interests in the movies that share the same “star” with his historically liked ones, while another user may be more concerned about the “genre” of movies.

To characterize the topological proximity structure of item v , we compute the linear combination of v ’s neighborhood:

$$\mathbf{v}_{\mathcal{N}(v)}^u = \sum_{e \in \mathcal{N}(v)} \tilde{\pi}_{r_{v,e}}^u \mathbf{e}, \quad (2)$$

where $\tilde{\pi}_{r_{v,e}}^u$ is the normalized user-relation score

$$\tilde{\pi}_{r_{v,e}}^u = \frac{\exp(\pi_{r_{v,e}}^u)}{\sum_{e \in \mathcal{N}(v)} \exp(\pi_{r_{v,e}}^u)}, \quad (3)$$

and \mathbf{e} is the representation of entity e . User-relation scores act as *personalized filters* when computing an entity’s neighborhood representation, since we aggregate the neighbors with bias with respect to these user-specific scores.

In a real-world knowledge graph, the size of $\mathcal{N}(e)$ may vary significantly over all entities. To keep the computational pattern of each batch fixed and more efficient, we uniformly sample a fixed-size set of neighbors for each entity instead of using its full neighbors. Specifically, we compute the neighborhood representation of entity v as $\mathbf{v}_{\mathcal{S}(v)}^u$, where $\mathcal{S}(v) \triangleq \{e \mid e \sim \mathcal{N}(v)\}$ and $|\mathcal{S}(v)| = K$ is a configurable constant.³ In KGCN, $\mathcal{S}(v)$ is also called the (single-layer) *receptive field* of entity v , as the final representation of v is sensitive to these locations. Figure 1a gives an illustrative example of a two-layer receptive field for a given entity, where K is set as 2.

The final step in a KGCN layer is to aggregate the entity representation \mathbf{v} and its neighborhood representation $\mathbf{v}_{\mathcal{S}(v)}^u$ into a single vector. We implement three types of aggregators $agg : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ in KGCN:

- *Sum aggregator* takes the summation of two representation vectors, followed by a nonlinear transformation:

$$agg_{sum} = \sigma(\mathbf{W} \cdot (\mathbf{v} + \mathbf{v}_{\mathcal{S}(v)}^u) + \mathbf{b}), \quad (4)$$

where \mathbf{W} and \mathbf{b} are transformation weight and bias, respectively, and σ is the nonlinear function such *ReLU*.

- *Concat aggregator* [7] concatenates the two representation vectors first before applying nonlinear transformation:

$$agg_{concat} = \sigma(\mathbf{W} \cdot \text{concat}(\mathbf{v}, \mathbf{v}_{\mathcal{S}(v)}^u) + \mathbf{b}). \quad (5)$$

²The knowledge graph \mathcal{G} is treated undirected.

³Technically, $\mathcal{S}(v)$ may contain duplicates if $\mathcal{N}(v) < K$.

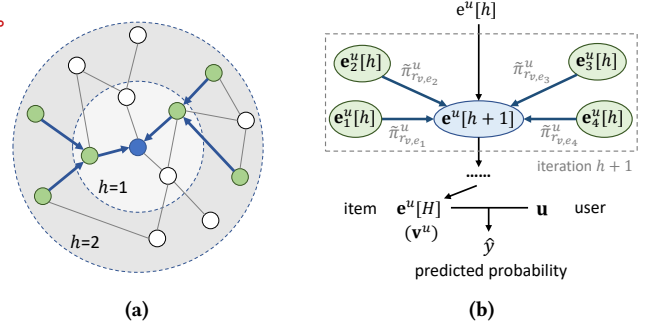


Figure 1: (a) A two-layer receptive field (green entities) of the blue entity in a KG. (b) The framework of KGCN.

- *Neighbor aggregator* [15] directly takes the neighborhood representation of entity v as the output representation:

$$agg_{neighbor} = \sigma(\mathbf{W} \cdot \mathbf{v}_{\mathcal{S}(v)}^u + \mathbf{b}). \quad (6)$$

Aggregation is a key step in KGCN, because the representation of an item is bound up with its neighbors by aggregation. We will evaluate the three aggregators in experiments.

3.3 Learning Algorithm

Through a single KGCN layer, the final representation of an entity is dependent on itself as well as its immediate neighbors, which we name 1-order entity representation. It is natural to extend KGCN from one layer to multiple layers to reasonably explore users’ potential interests in a broader and deeper way. The technique is intuitive: Propagating the initial representation of each entity (0-order representation) to its neighbors leads to 1-order entity representation, then we can repeat this procedure, i.e., further propagating and aggregating 1-order representations to obtain 2-order ones. Generally speaking, the h -order representation of an entity is a mixture of initial representations of itself and its neighbors up to h hops away. This is an important property for KGCN, which we will discuss in the next subsection.

The formal description of the above steps is presented in Algorithm 1. H denotes the maximum depth of receptive field (or equivalently, the number of aggregation iterations), and a suffix $[h]$ attached by a representation vector denotes h -order. For a given user-item pair (u, v) (line 2), we first calculate the receptive field \mathcal{M} of v in an iterative layer-by-layer manner (line 3, 13-19). Then the aggregation is repeated H times (line 5): In iteration h , we calculate the neighborhood representation of each entity $e \in \mathcal{M}[h]$ (line 7), then aggregate it with its own representation $\mathbf{e}^u[h-1]$ to obtain the one to be used at the next iteration (line 8). The final H -order entity representation is denoted as \mathbf{v}^u (line 9), which is fed into a function $f : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ together with user representation \mathbf{u} for predicting the probability:

$$\hat{y}_{uv} = f(\mathbf{u}, \mathbf{v}^u). \quad (7)$$

Figure 1b illustrates the KGCN algorithm in one iteration, in which the entity representation $\mathbf{v}^u[h]$ and neighborhood representations (green nodes) of a given node are mixed to form its representation for the next iteration (blue node).

Algorithm 1: KGCN algorithm

Input: Interaction matrix Y ; knowledge graph $\mathcal{G}(\mathcal{E}, \mathcal{R})$;
neighborhood sampling mapping $\mathcal{S} : e \rightarrow 2^{\mathcal{E}}$; trainable
parameters: $\{\mathbf{u}\}_{u \in \mathcal{U}}, \{\mathbf{e}\}_{e \in \mathcal{E}}, \{\mathbf{r}\}_{r \in \mathcal{R}}, \{\mathbf{W}_i, \mathbf{b}_i\}_{i=1}^H$;
hyper-parameters: $H, d, g(\cdot), f(\cdot), \sigma(\cdot), \text{agg}(\cdot)$
Output: Prediction function $\mathcal{F}(u, v | \Theta, Y, \mathcal{G})$

```

1 while KGCN not converge do
2   for  $(u, v)$  in  $Y$  do
3      $\{\mathcal{M}[i]\}_{i=0}^H \leftarrow \text{Get-Receptive-Field}(v)$ ;
4      $\mathbf{e}^u[0] \leftarrow \mathbf{e}, \forall e \in \mathcal{M}[0]$ ;
5     for  $h = 1, \dots, H$  do
6       for  $e \in \mathcal{M}[h]$  do
7          $\mathbf{e}_{\mathcal{S}(e)}^u[h-1] \leftarrow \sum_{e' \in \mathcal{S}(e)} \tilde{\pi}_{r_{e,e'}}^u \mathbf{e}^{u'}[h-1]$ ;
8          $\mathbf{e}^u[h] \leftarrow \text{agg}(\mathbf{e}_{\mathcal{S}(e)}^u[h-1], \mathbf{e}^u[h-1])$ ;
9      $\mathbf{v}^u \leftarrow \mathbf{e}^u[H]$ ;
10    Calculate predicted probability  $\hat{y}_{uv} = f(\mathbf{u}, \mathbf{v}^u)$ ;
11    Update parameters by gradient descent;
12 return  $\mathcal{F}$ ;

13 Function Get-Receptive-Field( $v$ )
14    $\mathcal{M}[H] \leftarrow v$ ;
15   for  $h = H-1, \dots, 0$  do
16      $\mathcal{M}[h] \leftarrow \mathcal{M}[h+1]$ ;
17     for  $e \in \mathcal{M}[h+1]$  do
18        $\mathcal{M}[h] \leftarrow \mathcal{M}[h] \cup \mathcal{S}(e)$ ;
19   return  $\{\mathcal{M}[i]\}_{i=0}^H$ ;
```

Note that Algorithm 1 traverses all possible user-item pairs (line 2). To make computation more efficient, we use a negative sampling strategy during training. The complete loss function is as follows:

$$\mathcal{L} = \sum_{u \in \mathcal{U}} \left(\sum_{v: y_{uv}=1} \mathcal{J}(y_{uv}, \hat{y}_{uv}) - \sum_{i=1}^{T^u} \mathbb{E}_{v_i \sim P(v_i)} \mathcal{J}(y_{uv_i}, \hat{y}_{uv_i}) \right) + \lambda \|\mathcal{F}\|_2^2, \quad (8)$$

where \mathcal{J} is cross-entropy loss, P is a negative sampling distribution, and T^u is the number of negative samples for user u . In this paper, $T^u = |\{v : y_{uv} = 1\}|$ and P follows a uniform distribution. The last term is the L2-regularizer.

4 EXPERIMENTS

In this section, we evaluate KGCN on three real-world scenarios: movie, book, and music recommendations.

4.1 Datasets

We utilize the following three datasets in our experiments for movie, book, and music recommendation, respectively:

- **MovieLens-20M**⁴ is a widely used benchmark dataset in movie recommendations, which consists of approximately 20 million explicit ratings (ranging from 1 to 5) on the MovieLens website.

⁴<https://grouplens.org/datasets/movielens/>

Table 1: Basic statistics and hyper-parameter settings for the three datasets (K : neighbor sampling size, d : dimension of embeddings, H : depth of receptive field, λ : L2 regularizer weight, η : learning rate).

| | MovieLens-20M | Book-Crossing | Last.FM |
|----------------|--------------------|--------------------|--------------------|
| # users | 138,159 | 19,676 | 1,872 |
| # items | 16,954 | 20,003 | 3,846 |
| # interactions | 13,501,622 | 172,576 | 42,346 |
| # entities | 102,569 | 25,787 | 9,366 |
| # relations | 32 | 18 | 60 |
| # KG triples | 499,474 | 60,787 | 15,518 |
| K | 4 | 8 | 8 |
| d | 32 | 64 | 16 |
| H | 2 | 1 | 1 |
| λ | 10^{-7} | 2×10^{-5} | 10^{-4} |
| η | 2×10^{-2} | 2×10^{-4} | 5×10^{-4} |
| batch size | 65,536 | 256 | 128 |

- **Book-Crossing**⁵ contains 1 million ratings (ranging from 0 to 10) of books in the Book-Crossing community.
- **Last.FM**⁶ contains musician listening information from a set of 2 thousand users from Last.fm online music system.

Since the three datasets are explicit feedbacks, we transform them into implicit feedback where each entry is marked with 1 indicating that the user has rated the item positively, and sample an unwatched set marked as 0 for each user. The threshold of positive rating is 4 for MovieLens-20M, while no threshold is set for Book-Crossing and Last.FM due to their sparsity.

We use **Microsoft Satori**⁷ to construct the knowledge graph for each dataset. We first select a subset of triples from the whole KG with a confidence level greater than 0.9. Given the sub-KG, we collect Satori IDs of all valid **movies/books/musicians** by matching their names with tail of triples (*head, film.film.name, tail*), (*head, book.book.title, tail*), or (*head, type.object.name, tail*). Items with multiple matched or no matched entities are excluded for simplicity. We then match the item IDs with the head of all triples and select all well-matched triples from the sub-KG. The basic statistics of the three datasets are presented in Table 1.

4.2 Baselines

We compare the proposed KGCN with the following baselines, in which the first two baselines are KG-free while the rest are all KG-aware methods. Hyper-parameter settings for baselines are introduced in the next subsection.

- **SVD** [11] is a classic CF-based model using inner product to model user-item interactions.⁸
- **LibFM** [14] is a feature-based factorization model in CTR scenarios. We concatenate user ID and item ID as input for LibFM.

⁵<http://www2.informatik.uni-freiburg.de/~cziegler/BX/>

⁶<https://grouplens.org/datasets/hetrec-2011/>

⁷<https://searchengineland.com/library/bing/bing-satori>

⁸We have tried NCF [8], i.e., replacing inner product with neural networks, but the result is inferior to SVD. Since SVD and NCF are similar, we only present the better one here.

Table 2: The results of AUC and F1 in CTR prediction.

| Model | MovieLens-20M | | Book-Crossing | | Last.FM | |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| | AUC | F1 | AUC | F1 | AUC | F1 |
| SVD | 0.963 (-1.5%) | 0.919 (-1.4%) | 0.672 (-8.9%) | 0.635 (-7.7%) | 0.769 (-3.4%) | 0.696 (-3.5%) |
| LibFM | 0.959 (-1.9%) | 0.906 (-2.8%) | 0.691 (-6.4%) | 0.618 (-10.2%) | 0.778 (-2.3%) | 0.710 (-1.5%) |
| LibFM + TransE | 0.966 (-1.2%) | 0.917 (-1.6%) | 0.698 (-5.4%) | 0.622 (-9.6%) | 0.777 (-2.4%) | 0.709 (-1.7%) |
| PER | 0.832 (-14.9%) | 0.788 (-15.5%) | 0.617 (-16.4%) | 0.562 (-18.3%) | 0.633 (-20.5%) | 0.596 (-17.3%) |
| CKE | 0.924 (-5.5%) | 0.871 (-6.5%) | 0.677 (-8.3%) | 0.611 (-11.2%) | 0.744 (-6.5%) | 0.673 (-6.7%) |
| RippleNet | 0.968 (-1.0%) | 0.912 (-2.1%) | 0.715 (-3.1%) | 0.650 (-5.5%) | 0.780 (-2.0%) | 0.702 (-2.6%) |
| KGCN-sum | 0.978 | 0.932* | 0.738 | 0.688* | 0.794 (-0.3%) | 0.719 (-0.3%) |
| KGCN-concat | 0.977 (-0.1%) | 0.931 (-0.1%) | 0.734 (-0.5%) | 0.681 (-1.0%) | 0.796* | 0.721* |
| KGCN-neighbor | 0.977 (-0.1%) | 0.932* | 0.728 (-1.4%) | 0.679 (-1.3%) | 0.781 (-1.9%) | 0.699 (-3.1%) |
| KGCN-avg | 0.975 (-0.3%) | 0.929 (-0.3%) | 0.722 (-2.2%) | 0.682 (-0.9%) | 0.774 (-2.8%) | 0.692 (-4.0%) |

* Statistically significant improvement by unpaired two-sample t -test with $p = 0.1$.

- **LibFM + TransE** extends LibFM by attaching an entity representation learned by TransE [1] to each user-item pair.
- **PER** [22] treats the KG as heterogeneous information networks and extracts meta-path based features to represent the connectivity between users and items.
- **CKE** [23] combines CF with structural, textual, and visual knowledge in a unified framework for recommendation. We implement CKE as CF plus a structural knowledge module in this paper.
- **RippleNet** [18] is a memory-network-like approach that propagates users’ preferences on the KG for recommendation.

4.3 Experiments Setup

In KGCN, we set functions g and f as inner product, σ as $ReLU$ for non-last-layer aggregator and $tanh$ for last-layer aggregator. Other hyper-parameter settings are provided in Table 1. The hyper-parameters are determined by optimizing AUC on a validation set. For each dataset, the ratio of training, evaluation, and test set is 6 : 2 : 2. Each experiment is repeated 3 times, and the average performance is reported. We evaluate our method in two experiment scenarios: (1) In click-through rate (CTR) prediction, we apply the trained model to predict each interaction in the test set. We use AUC and F1 to evaluate CTR prediction. (2) In top-K recommendation, we use the trained model to select K items with highest predicted click probability for each user in the test set, and choose Recall@K to evaluate the recommended sets. All trainable parameters are optimized by Adam algorithm. The code of KGCN-LS is implemented under Python 3.6, TensorFlow 1.12.0, and NumPy 1.14.3.

The hyper-parameter settings for baselines are as follows. For SVD, we use the unbiased version (i.e., the predicted rating is modeled as $r_{pq} = \mathbf{p}^T \mathbf{q}$). The dimension and learning rate for the four datasets are set as: $d = 8, \eta = 0.5$ for MovieLens-20M, Book-Crossing; $d = 8, \eta = 0.1$ for Last.FM. For LibFM, the dimension is $\{1, 1, 8\}$ and the number of training epochs is 50. The dimension of TransE is 32. For PER, we use manually designed user-item-attribute-item paths as features (i.e., “user-movie-director-movie”, “user-movie-genre-movie”, and “user-movie-star-movie” for MovieLens-20M; “user-book-author-book” and “user-book-genre-book” for Book-Crossing, “user-musician-date_of_birth-musician”

(date of birth is discretized), “user-musician-country-musician”, and “user-musician-genre-musician” for Last.FM). For CKE, the dimension of the three datasets are 64, 128, 64. The training weight for KG part is 0.1 for all datasets. The learning rate are the same as in SVD. For RippleNet, $d = 8, H = 2, \lambda_1 = 10^{-6}, \lambda_2 = 0.01, \eta = 0.01$ for MovieLens-20M; $d = 16, H = 3, \lambda_1 = 10^{-5}, \lambda_2 = 0.02, \eta = 0.005$ for Last.FM. Other hyper-parameters are the same as reported in their original papers or as default in their codes.

4.4 Results

The results of CTR prediction and top-K recommendation are presented in Table 2 and Figure 2, respectively (SVD, LibFM and other variants of KGCN are not plotted in Figure 2 for clarity). We have the following observations:

- In general, we find that the improvements of KGCN on book and music are higher than movie. This demonstrates that KGCN can well address sparse scenarios, since Book-Crossing and Last.FM are much sparser than MovieLens-20M.
- The performance of KG-free baselines, SVD and LibFM, are actually better than the two KG-aware baselines PER and CKE, which indicates that PER and CKE cannot make full use of the KG with manually designed meta-paths and TransR-like regularization.
- LibFM + TransE is better than LibFM in most cases, which demonstrates that the introduction of KG is helpful for recommendation in general.
- PER performs worst among all baselines, since it is hard to define optimal meta-paths in reality.
- RippleNet shows strong performance compared with other baselines. Note that RippleNet also uses multi-hop neighborhood structure, which interestingly shows that capturing proximity information in the KG is essential for recommendation.

The last four rows in Table 2 summarize the performance of KGCN variants. The first three (sum, concat, neighbor) correspond to different aggregators introduced in the preceding section, while the last variant KGCN-avg is a reduced case of KGCN-sum where

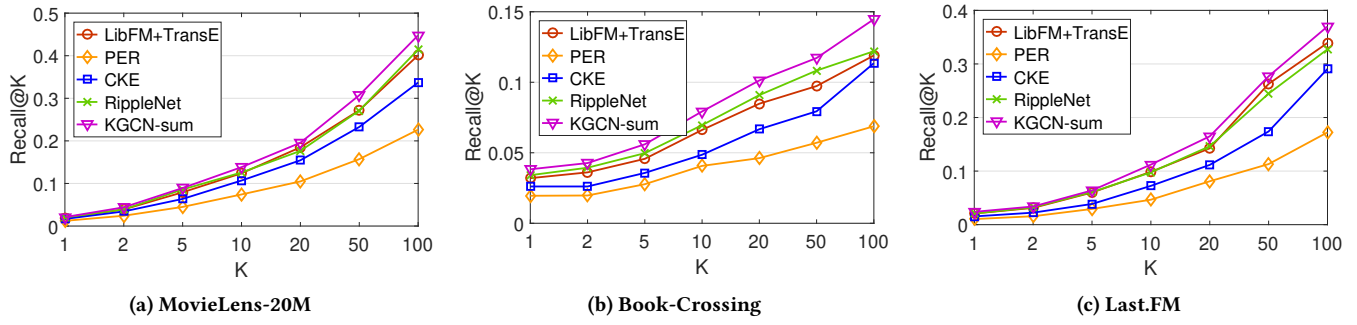


Figure 2: The results of Recall@K in top-K recommendation.

Table 3: AUC result of KGCN with different neighbor sampling size K .

| K | 2 | 4 | 8 | 16 | 32 | 64 |
|---------------|-------|--------------|--------------|-------|-------|-------|
| MovieLens-20M | 0.978 | 0.979 | 0.978 | 0.978 | 0.977 | 0.978 |
| Book-Crossing | 0.680 | 0.727 | 0.736 | 0.725 | 0.711 | 0.723 |
| Last.FM | 0.791 | 0.794 | 0.795 | 0.793 | 0.794 | 0.792 |

Table 4: AUC result of KGCN with different depth of receptive field H .

| H | 1 | 2 | 3 | 4 |
|---------------|--------------|--------------|-------|-------|
| MovieLens-20M | 0.972 | 0.976 | 0.974 | 0.514 |
| Book-Crossing | 0.738 | 0.731 | 0.684 | 0.547 |
| Last.FM | 0.794 | 0.723 | 0.545 | 0.534 |

neighborhood representations are directly averaged without user-relation scores (i.e., $\mathbf{v}_{N(v)}^u = \sum_{e \in N(v)} \mathbf{e}$ instead of Eq. (2)). Therefore, KGCN-avg is used to examine the efficacy of the ‘‘attention mechanism’’. From the results we find that:

- KGCN outperforms all baselines by a significant margin, while their performances are slightly distinct: KGCN-sum performs best in general, while the performance of KGCN-neighbor shows a clear gap on Book-Crossing and Last.FM. This may be because the neighbor aggregator uses the neighborhood representation only, thus losing useful information from the entity itself.
- KGCN-avg performs worse than KGCN-sum, especially in Book-Crossing and Last.FM where interactions are sparse. This demonstrates that capturing users’ personalized preferences and semantic information of the KG do benefit the recommendation.

4.4.1 Impact of neighbor sampling size. We vary the size of sampled neighbor K to investigate the efficacy of usage of the KG. From Table 3 we observe that KGCN achieves the best performance when $K = 4$ or 8. This is because a too small K does not have enough capacity to incorporate neighborhood information, while a too large K is prone to be misled by noises.

4.4.2 Impact of depth of receptive field. We investigate the influence of depth of receptive field in KGCN by varying H from 1 to 4. The results are shown in Table 4, which demonstrate that KGCN

Table 5: AUC result of KGCN with different dimension of embedding.

| d | 4 | 8 | 16 | 32 | 64 | 128 |
|---------------|-------|-------|--------------|--------------|--------------|-------|
| MovieLens-20M | 0.968 | 0.970 | 0.975 | 0.977 | 0.973 | 0.972 |
| Book-Crossing | 0.709 | 0.732 | 0.733 | 0.735 | 0.739 | 0.736 |
| Last.FM | 0.789 | 0.793 | 0.797 | 0.793 | 0.790 | 0.789 |

is more sensitive to H compared to K . We observe the occurrence of serious model collapse when $H = 3$ or 4, as a larger H brings massive noises to the model. This is also in accordance with our intuition, since a too long relation-chain makes little sense when inferring inter-item similarities. An H of 1 or 2 is enough for real cases according to the experiment results.

4.4.3 Impact of dimension of embedding. Lastly, we examine the influence of dimension of embedding d on performance of KGCN. The result in Table 5 is rather intuitive: Increasing d initially can boost the performance since a larger d can encode more information of users and entities, while a too large d adversely suffers from overfitting.

5 CONCLUSIONS AND FUTURE WORK

This paper proposes knowledge graph convolutional networks for recommender systems. KGCN extends non-spectral GCN approaches to the knowledge graph by aggregating neighborhood information selectively and biasedly, which is able to learn both structure information and semantic information of the KG as well as users’ personalized and potential interests. We also implement the proposed method in a minibatch fashion, which is able to operate on large datasets and knowledge graphs. Through extensive experiments on real-world datasets, KGCN is shown to consistently outperform state-of-the-art baselines in movie, book, and music recommendation.

We point out three avenues for future work. (1) In this work we uniformly sample from the neighbors of an entity to construct its receptive field. Exploring a non-uniform sampler (e.g., importance sampling) is an important direction of future work. (2) This paper (and all literature) focuses on modeling item-end KGs. An interesting direction of future work is to investigate whether leveraging user-end KGs is useful in improving the performance of recommendation. (3) Designing an algorithm to well combine KGs at the two ends is also a promising direction.

REFERENCES

- [1] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems*. 2787–2795.
- [2] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann Lecun. 2014. Spectral networks and locally connected networks on graphs. In *the 2nd International Conference on Learning Representations*.
- [3] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*. ACM, 7–10.
- [4] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*. 3844–3852.
- [5] Qiming Diao, Minghui Qiu, Chao-Yuan Wu, Alexander J Smola, Jing Jiang, and Chong Wang. 2014. Jointly modeling aspects, ratings and sentiments for movie recommendation (JMARS). In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 193–202.
- [6] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems*. 2224–2232.
- [7] Will Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*. 1024–1034.
- [8] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web*. 173–182.
- [9] Jin Huang, Wayne Xin Zhao, Hongjian Dou, Ji-Rong Wen, and Edward Y Chang. 2018. Improving Sequential Recommendation with Knowledge-Enhanced Memory Networks. In *the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. ACM, 505–514.
- [10] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *the 5th International Conference on Learning Representations*.
- [11] Yehuda Koren. 2008. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 426–434.
- [12] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning entity and relation embeddings for knowledge graph completion. In *AAAI*, Vol. 15. 2181–2187.
- [13] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. 2016. Learning convolutional neural networks for graphs. In *International Conference on Machine Learning*. 2014–2023.
- [14] Steffen Rendle. 2012. Factorization machines with libfm. *ACM Transactions on Intelligent Systems and Technology (TIST)* 3, 3 (2012), 57.
- [15] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. In *Proceedings of the 6th International Conferences on Learning Representations*.
- [16] Hongwei Wang, Jia Wang, Miao Zhao, Jiannong Cao, and Minyi Guo. 2017. Joint Topic-Semantic-aware Social Recommendation for Online Voting. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. ACM, 347–356.
- [17] Hongwei Wang, Fuzheng Zhang, Min Hou, Xing Xie, Minyi Guo, and Qi Liu. 2018. SHINE: signed heterogeneous information network embedding for sentiment link prediction. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining*. ACM, 592–600.
- [18] Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. 2018. RippleNet: Propagating User Preferences on the Knowledge Graph for Recommender Systems. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM.
- [19] Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. 2018. DKN: Deep Knowledge-Aware Network for News Recommendation. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*. 1835–1844.
- [20] Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. 2017. Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering* 29, 12 (2017), 2724–2743.
- [21] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [22] Xiao Yu, Xiang Ren, Yizhou Sun, Quanquan Gu, Bradley Sturt, Urvashi Khandelwal, Brandon Norick, and Jiawei Han. 2014. Personalized entity recommendation: A heterogeneous information network approach. In *Proceedings of the 7th ACM International Conference on Web Search and Data Mining*. ACM, 283–292.
- [23] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. 2016. Collaborative knowledge base embedding for recommender systems. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 353–362.
- [24] Huan Zhao, Quanming Yao, Jianda Li, Yangqiu Song, and Dik Lun Lee. 2017. Meta-graph based recommendation fusion over heterogeneous information networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 635–644.
- [25] Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang, Nicholas Jing Yuan, Xing Xie, and Zhenhui Li. 2018. DRN: A Deep Reinforcement Learning Framework for News Recommendation. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*. 167–176.
- [26] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 1059–1068.