

# Graph convolution machine for context-aware recommender system

Jiancan WU<sup>1</sup>, Xiangnan HE (✉)<sup>1</sup>, Xiang WANG<sup>2</sup>, Qifan WANG<sup>3</sup>, Weijian CHEN<sup>1</sup>,  
Jianxun LIAN<sup>4</sup>, Xing XIE<sup>4</sup>

<sup>1</sup> School of Information Science and Technology, University of Science and Technology of China, Hefei 230026, China

<sup>2</sup> 5 Prince George's Park, National University of Singapore, Singapore 118404, Singapore

<sup>3</sup> Google Research, Mountain View, CA 94043, USA

<sup>4</sup> Microsoft Research Asia, Beijing 100190, China

© Higher Education Press 2022

**Abstract** The latest advance in recommendation shows that better user and item representations can be learned via performing graph convolutions on the user-item interaction graph. However, such finding is mostly restricted to the collaborative filtering (CF) scenario, where the interaction contexts are not available. In this work, we extend the advantages of graph convolutions to context-aware recommender system (CARS, which represents a generic type of models that can handle various side information). We propose *Graph Convolution Machine* (GCM), an end-to-end framework that consists of three components: an encoder, graph convolution (GC) layers, and a decoder. The encoder projects users, items, and contexts into embedding vectors, which are passed to the GC layers that refine user and item embeddings with context-aware graph convolutions on the user-item graph. The decoder digests the refined embeddings to output the prediction score by considering the interactions among user, item, and context embeddings. We conduct experiments on three real-world datasets from Yelp and Amazon, validating the effectiveness of GCM and the benefits of performing graph convolutions for CARS.

**Keywords** context-aware recommender systems, graph convolution

## 1 Introduction

Recommendation has become a pervasive service in today's Web, serving as an important tool to alleviate information overload and improve user experience. The key data source for building a recommendation service is user-item interactions, e.g., clicks and purchases, which spawn wide research efforts on collaborative filtering (CF) [1–3] that leverage the interaction data only to predict user preference. Recently, inspired by the success of graph neural networks (GNNs) [4,5], researchers have attempted to employ GNNs on

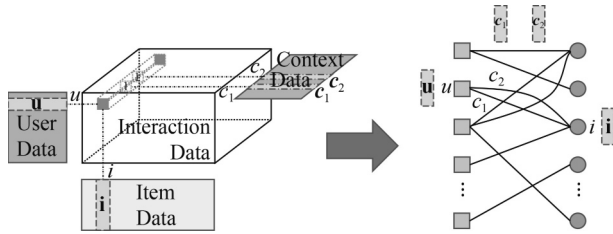
recommendation in which CF signals are exhibited as high-order connectivity [3,6–8]. While CF provides a universal solution for recommendation, it falls short in utilizing the side information of interaction contexts. In many scenarios, the current contexts could have a substantial impact on user choice. For example, in restaurant recommendation, the current time and location can effectively filter out unsuitable candidates; in E-commerce, the click behaviors in recent sessions provide strong signal about the user's next purchase. As such, it is important to develop context-aware recommender system (CARS) that can effectively integrate contexts (and possibly other side information like user profiles and item attributes) into user preference prediction [9].

Inspired by the matrix completion view of CF, early research naturally extended the problem of CARS to tensor completion [10], which however suffers from high complexity. Later on, Rendle proposed factorization machine (FM) [11], which addressed CARS from the view of standard supervised learning for the first time. Specifically, it converts all information related to an interaction to a feature vector via multi-hot encoding, modeling the second-order feature interactions to predict the interaction label. Due to its generality and effectiveness, FM soon becomes a prevalent solution for CARS and is followed by many work. For example, in the era of deep learning, Wide&Deep [12] and Deep Crossing [13] replaced the second-order interaction modeling with a neural network for implicit interaction modeling; recently, Neural FM [14], Attentional FM [15], xDeepFM [16], and Convolutional FM [17] extended FM with various kinds of neural networks to enhance its expressiveness.

Summarizing existing CARS models, we can find a common drawback: they follow the standard supervised learning scheme that ignores the relationship among data instances. This may limit the model's effectiveness in capturing the CF effect, since it needs to consider multiple interactions simultaneously to recognize the CF patterns. An evidence is from the neural graph collaborative filtering (NGCF) work [3], which demonstrates that connecting the interactions in the predictive

model significantly improves the embedding quality for CF. Since in CARS user-item interactions still play an important role by reflecting user preference, it is reasonable to believe that adequately modeling the relationship among interactions can improve the model quality. Moreover, the recent neural network-based methods like xDeepFM [16] and Convolutional FM [17] suffer from low efficiency in online serving, since each candidate item needs be scored separately with the deep model architecture that models complex feature interactions, which could be very time-consuming.

In this work, we aim to propose new CARS model by addressing the above-mentioned limitations. Firstly, we cast the data in CARS as an attributed user-item graph, where the side information of users and items are represented as node features, and the contexts are represented as edge features (Fig. 1). Secondly, we propose an end-to-end model that consists of three components: an encoder, graph convolution (GC) layers, and a decoder (Fig. 2). The encoder projects users, items, and contexts into embedding vectors; the GC layers then exploit the interactions to refine the embeddings via performing graph convolutions; lastly, the decoder models the interactions among embeddings via FM to output the prediction score. After the model is trained, the refined embeddings by GC layers can be pre-computed before



**Fig. 1** The data used for building a CARS. The mixture data of interaction tensor and user/item/context feature matrices are converted to an attributed user-item bipartite graph without loss of fidelity

serving. As such, the time complexity of online serving is the same as FM, being much more efficient than the recent neural network methods.

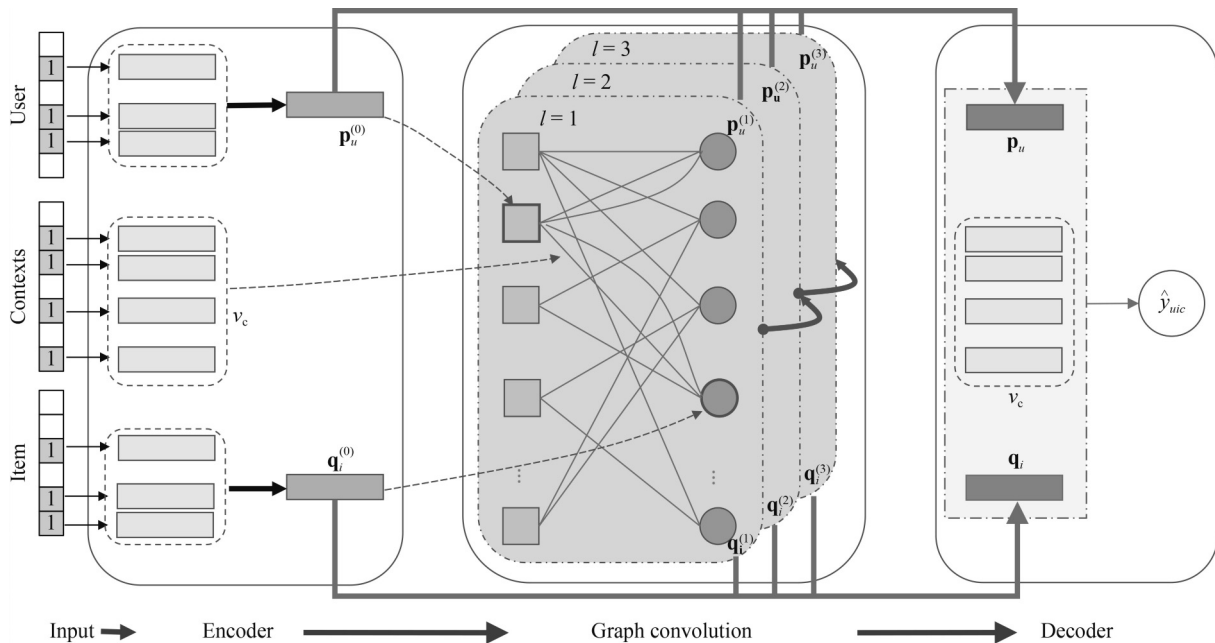
We summarize the contributions of this work as follows:

- We highlight the limitation of the mainstream supervised learning schemes and the necessity of exploiting the relationship among data instances in the predictive model of CARS.
- We propose a new model named Graph Convolution Machine (GCM), unifying the strengths of graph convolution network and factorization machine for CARS.
- We conduct extensive experiments on three real-world datasets which demonstrate the effectiveness and efficiency of GCM.

## 2 Related work

### 2.1 Context-aware recommendation

Extensive studies on context-aware recommender system (CARS) [11,14,16] have been conducted and achieved great success. Learning informative representations, based on user-item interactions (e.g., clicks, purchases) and contextual features (e.g., location, time, last purchase), has been a central theme of research on CARS. Towards this end, modeling interactions among different features is showing promise. Early, factorization machine (FM) [11] embeds each feature into a vector representation, and utilizes inner product to capture their pairwise relationships (e.g., the second-order feature interactions). Due to its generality and effectiveness, FM becomes a prevalent solution for CARS. Many works resort to this paradigm, such as FFM [18]. Recent works [12,14–16,19] leverage deep neural networks to model higher-order feature interactions, so as to generate better representations and enhance recommendation performance. For example, NFM [14] proposes a bilinear interaction operation which uses a sum pooling over the pair-wise dot-product of feature



**Fig. 2** The graph convolution machine model

vectors; AFM [15] learns the importance of each feature interaction via the attention mechanism; xDeepFM [16] extends the Cross Network [20] to the Compressed Interaction Network (CIN) which models high-order interactions explicitly at vector-wise level; while Convolutional FM [17] models second-order interaction with outer product, forming an interaction cube, then applying 3D convolution to learn high-order interactions. It is worth mentioning that another research line close to CARS is the CTR (Click Through Rate) prediction [21–24], which also focuses on modeling the complex feature interactions. The key difference lies in the evaluation protocol: most CARS models adopt top-k recommendation protocols, while CTR prediction models measure log loss or AUC metrics on positive/negative samples.

Despite effectiveness, we argue that present works treat user interactions as isolated data instances, while forgoing their relationships (e.g., user behaviors happened at the same time and location are highly likely to reflect user preferences). This would easily lead to suboptimal representations and limit the performance. We hence aim to explore relationships among user behaviors in this work.

## 2.2 Graph neural networks for recommendation

Another relevant research line is to leverage graph neural networks (GNNs) for recommendation. In particular, GNN models [4,5,25] exploit graph structure to guide the representation learning. The basic idea is the embedding propagation mechanism, which aggregates the embeddings of neighbors to update the target node’s embedding. By recursively performing such propagations, the information from multi-hop neighbors is encoded into the representation of the target node. GNN models have been widely used in many fundamental tasks due to their strong representation ability, spanning from node classification [26], link prediction [27], to graph classification [28], and achieved remarkable improvements.

Inspired by their success, researchers have attempted to employ GNNs on recommendation. Recent works on collaborative filtering (CF), such as NGCF [3], GC-MC [29], SpectralCF [8] and PinSage [30], reorganize historical user behaviors in the form of a user-item bipartite graph, exhibit CF signals as high-order connectivity, and encode such signals into representations. For CTR prediction task, Fi-GNN [31] takes multi-field features into consideration by constructing feature graph for each instance and converting the task of modeling feature interactions among fields into modeling node interactions on the feature graph; GIN [32] models implicit user intention by the multi-layered intention diffusion and aggregation on the co-occurrence click relationship graph; [33] builds the multi-relational item graph and applies GNN to capture complex transition relations between items in user behavior sequences. Moreover, GNN models have also been employed on other recommendation tasks, including social recommendation [17,34], sequential recommendation [35,36], and knowledge-aware recommendation [37,38]. As such, aggregating useful information from multi-hop neighbors is able to achieve better expressiveness, than single ID embed-

dings. Hence, it is reasonable to believe that graph learning is a promising solution to model the relationships among interactions adequately.

## 3 Problem definition

We divide the data used for CARS into four types: users, items, contexts, and interactions. Following [18], we define context as the information that is associated with an interaction, e.g., the current location, time, previous click, etc. Figure 1 illustrates the data in CARS, where the main data is the user-item-context interaction tensor. In the sparse tensor, each nonzero entry  $(u, i, c)$  denotes that the user  $u$  has interacted with the item  $i$  under the context  $c$ ; we give such entries a label of 1, i.e.,  $y_{uic} = 1$ . Each  $u, i, c$  is respectively associated with a multi-hot feature vector  $\mathbf{u}$ ,  $\mathbf{i}$ , and  $\mathbf{c}$ , which contain the features that describe the user, item, and context. For example,  $\mathbf{u}$  includes static user profiles like gender and interested tags,  $\mathbf{i}$  includes static item attributes like category and price, and  $\mathbf{c}$  includes dynamic contexts like the user’s current location and the time.

Given such data, we convert it to an attributed user-item bipartite graph with the same representation power. Specifically, each vertex represents a user or an item, and each edge represents the interaction between the connected user and item. Each vertex or edge is associated with a feature vector  $\mathbf{u}$ ,  $\mathbf{i}$ , or  $\mathbf{c}$ . Note that there may exist multiple edges between a user-item pair, since a user may interact with the same item multiple times under different contexts. We denote all edges in the graph as the set  $\mathcal{Y} = \{(u, i, c) | y_{uic} = 1\}$ , the neighbors of the user  $u$  as the set  $\mathcal{N}_u = \{(i, c) | y_{uic} = 1\}$ , and neighbors of the item  $i$  as the set  $\mathcal{N}_i = \{(u, c) | y_{uic} = 1\}$ .

We formulate the problem of CARS as:

**Input**: User-item-context interactions  $\{(u, i, c) | y_{uic} = 1\}$ , feature vectors of users  $\{\mathbf{u}\}$ , items  $\{\mathbf{i}\}$ , and contexts  $\{\mathbf{c}\}$ .

**Output**: Prediction function  $f : \mathbf{u}, \mathbf{i}, \mathbf{c} \rightarrow \mathbb{R}$ , which takes the feature vector of a user, an item, and a context as the input, and outputs a real value that estimates how likely the user will interact with the item under the context.

## 4 Graph convolution machine (GCM)

We present our method in this section. We first describe the predictive model, followed by the model complexity analyses and optimization details.

### 4.1 Predictive model

Figure 2 illustrates the model framework, which consists of three components: an encoder, graph convolution layers, and a decoder. We next describe each component one by one.

#### 4.1.1 Encoder

The input to the encoder has three fields: user-field features  $\mathbf{u}$ , item-field features  $\mathbf{i}$ , and the context-field features  $\mathbf{c}$ . We include the ID feature into the user-field and item-field features, since it helps to differentiate users (items) when their profiles (attributes) are the same<sup>1)</sup>. For each nonzero feature,

<sup>1)</sup> Note that there is no need to include ID into the context-field features, since a context  $c$  and its features  $\mathbf{c}$  are one-to-one mapping.

we associate it with an embedding vector, resulting in a set of embeddings to describe the input user, item, and context, respectively. We then pool the set of user (and item) field into a vector, so as to feed the vector into the following GC layers to refine the user (and item) representations. Specifically, we adopt average pooling, that is,

$$\mathbf{p}_u^{(0)} = \frac{1}{|\mathbf{u}|} \mathbf{P}^T \mathbf{u}, \quad (1)$$

where  $|\mathbf{u}|$  denotes the number of nonzero features in  $\mathbf{u}$ , and  $\mathbf{P} \in \mathbb{R}^{U \times D}$  is the embedding matrix for user features, where  $U$  denotes the number of total user features and  $D$  denotes the embedding size.  $\mathbf{p}_u^{(0)}$  denotes the initial representation vector for  $u$ . Similarly, we get the initial representation vector for item  $i$  as  $\mathbf{q}_i^{(0)}$ .

Note that other pooling mechanisms can be applied here, such as the attention-based pooling [17,39,40] which learns varying weights for feature embeddings. However, we tried that and found it does not improve the performance. Thus we keep the simplest average pooling and avoid introducing additional parameters. Since we do not update the context representation in the following GC layers, we do not perform pooling on the context field. We denote the set of context-field embeddings as  $\mathcal{V}_c = \{\mathbf{v}_s | s \in \mathbf{c}\}$ , where  $s \in \mathbf{c}$  denotes the nonzero feature in  $\mathbf{c}$  and  $\mathbf{v}_s$  stands for the embedding vector of context feature  $s$ . The encoder outputs  $\mathbf{p}_u^{(0)}$ ,  $\mathbf{q}_i^{(0)}$ , and  $\mathcal{V}_c$ , which are fed into the next component of GC layers.

#### 4.1.2 Graph convolution layers

This is the core component of GCM, designed to address the limitation of existing supervised learning-based CARS models. It refines  $\mathbf{p}_u^{(0)}$  and  $\mathbf{q}_i^{(0)}$  by exploiting holistic user-item interaction data, which can augment the user and item representations with explicit collaborative filtering signal [3]. The GC on the user-item graph is typically formulated as a message propagation framework:

$$\mathbf{p}_u^{(l+1)} = \sum_{i \in \mathcal{N}_u} g(\mathbf{p}_u^{(l)}, \mathbf{q}_i^{(l)}); \quad \mathbf{q}_i^{(l+1)} = \sum_{u \in \mathcal{N}_i} g(\mathbf{q}_i^{(l)}, \mathbf{p}_u^{(l)}), \quad (2)$$

where  $\mathbf{p}_u^{(l)}$  and  $\mathbf{q}_i^{(l)}$  denote the refined user representation and item representation of the  $l$ th GC layer, respectively, and  $g(\cdot)$  is a self-defined function. Recursively conducting such message propagation relates the representation of a user with her high-order neighbors, e.g., first-order for interacted items and second-order for co-interacted users, which is beneficial for collaborative filtering; and the same logic applies to item representation.

However, the standard GC does not consider the features on edges. In our constructed user-item graph, the edges between a user and an item carry the context features, which are important for understanding the context-dependent interaction patterns. For example, a user may prefer bars on Friday, and a restaurant is more popular at lunch time. As such, better user and item representations can be obtained if the context features can be adequately integrated into the GC.

To this end, we propose a new GC operation that incorporates the edge features of contexts:

$$\begin{aligned} \mathbf{p}_u^{(l+1)} &= \sum_{(i,c) \in \mathcal{N}_u} \frac{1}{\sqrt{|\mathcal{N}_u|}} (\mathbf{q}_i^{(l)} + \frac{1}{|\mathcal{V}_c|} \sum_{\mathbf{v}_s \in \mathcal{V}_c} \mathbf{v}_s), \\ \mathbf{q}_i^{(l+1)} &= \sum_{(u,c) \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i|}} (\mathbf{p}_u^{(l)} + \frac{1}{|\mathcal{V}_c|} \sum_{\mathbf{v}_s \in \mathcal{V}_c} \mathbf{v}_s). \end{aligned} \quad (3)$$

Next, we explain the rationality of the GC of the user side, since the item side can be interpreted in the same way. Here  $|\mathcal{N}_u|$  denotes the number of edges connected with the user  $u$ , and the coefficient  $\frac{1}{\sqrt{|\mathcal{N}_u|}}$  is a normalization term to avoid the scale of embedding values increasing with the GC. We incorporate the context features by averaging their embeddings and adding to the connected user embedding. This way, we build the connection between a user with both her interacted item and the interacted context. It is expected to capture the effect that if a user likes to choose an item under a certain context, then their representations are similar. Note that we have tried more complicated mechanisms like incorporating the pairwise interactions among  $\mathcal{V}_c$  and  $\mathbf{q}_i^{(l)}$ , and using a MLP to capture high-order interactions. However, these ways do not lead to performance improvements. Thus we use this simple average operation, which is easy to interpret and train (no additional parameters are introduced).

By stacking multiple such GC layers, a user (or an item) representation can be refined by its multi-hop neighbors. Since the representations of different layers carry different semantics, we next combine the representations of all layers to form a more comprehensive representation:

$$\mathbf{p}_u = \sum_{l=0}^L \alpha_l \mathbf{p}_u^{(l)}; \quad \mathbf{q}_i = \sum_{l=0}^L \alpha_l \mathbf{q}_i^{(l)}, \quad (4)$$

where  $\alpha_l$  denotes the weight of the  $l$ th layer representation, which can be treated as hyper-parameter and tuned via a grid search with the constraint that  $\alpha_l \geq 0$  and  $\sum_{l=0}^L \alpha_l = 1$ . However, the workload of tuning them increases exponentially, as the GCN goes deep. In our experiments, we find that setting  $\alpha_l$  to  $1/(L+1)$  generally leads to satisfactory performance. Therefore, we fix  $\alpha_l$  to  $1/(L+1)$  for simplicity. A possible extension is to learn  $\alpha_l$ , e.g., designing attention mechanisms or optimizing them on the validation data. We leave this extension as future work, since it is not the focus of this work.

In what follows, we provide the matrix form of GC layers for implementation. Let the user-item interaction matrix be  $\mathbf{R}_{ui} \in \mathbb{R}^{N \times M}$ , where  $N$  and  $M$  denotes the number of users and items. Each entry  $r_{ui} \in \mathbf{R}_{ui}$  is the number of times user  $u$  interacts with item  $i$ . Similarly, we utilize  $\mathbf{R}_{uc} \in \mathbb{R}^{N \times K}$  and  $\mathbf{R}_{ic} \in \mathbb{R}^{M \times K}$  to denote user-context interaction matrix and item-context interaction matrix, respectively, where  $K$  is the number of contexts. Then we define the adjacency matrix of the user-item-context graph as

$$\mathbf{A} = \begin{pmatrix} \mathbf{0} & \mathbf{R}_{ui} & \mathbf{R}_{uc} \\ \mathbf{R}_{ui}^T & \mathbf{0} & \mathbf{R}_{ic} \\ \mathbf{0} & \mathbf{0} & 2\mathbf{I} \end{pmatrix}, \quad (5)$$

where  $\mathbf{0}$  is all-zero matrix,  $\mathbf{I}$  is identity matrix. Let  $\mathbf{D}$  be

diagonal degree matrix of  $\mathbf{A}$ , that is, the  $t$ th diagonal element  $\mathbf{D}_{tt} = \sum_j \mathbf{A}_{tj}$ . The normalized adjacency matrix can be expressed as

$$\hat{\mathbf{A}} = \sqrt{2\mathbf{D}}^{-\frac{1}{2}} \mathbf{A}. \quad (6)$$

Then, we get the matrix form of the layer-wise propagation rule which is equivalent to Eq. (3):

$$\mathbf{E}^{(l)} = \hat{\mathbf{A}} \mathbf{E}^{(l-1)}, \quad (7)$$

where  $\mathbf{E}^{(l)} \in \mathbb{R}^{(N+M+K) \times D}$  is the concatenate of user, item and context embedding matrix.  $\mathbf{E}^{(0)}$  is set as the concatenate matrix of encoded embedding tables from Encoder, which can be expressed as

$$\mathbf{E}^{(0)} = [\underbrace{\mathbf{p}_{u_1}^{(0)}, \dots, \mathbf{p}_{u_N}^{(0)}}_{\text{user embeddings}}, \underbrace{\mathbf{q}_{i_1}^{(0)}, \dots, \mathbf{q}_{i_M}^{(0)}}_{\text{item embeddings}}, \underbrace{\mathbf{r}_{c_1}, \dots, \mathbf{r}_{c_K}}_{\text{context embeddings}}]^\top. \quad (8)$$

Lastly, we get the final embedding matrix

$$\begin{aligned} \mathbf{E} &= \alpha_0 \mathbf{E}^{(0)} + \alpha_1 \mathbf{E}^{(1)} + \alpha_2 \mathbf{E}^{(2)} + \dots + \alpha_L \mathbf{E}^{(L)} \\ &= \alpha_0 \mathbf{E}^{(0)} + \alpha_1 \hat{\mathbf{A}} \mathbf{E}^{(0)} + \alpha_2 \hat{\mathbf{A}}^2 \mathbf{E}^{(0)} + \dots + \alpha_L \hat{\mathbf{A}}^L \mathbf{E}^{(0)}. \end{aligned} \quad (9)$$

#### 4.1.3 Decoder

The GC layers output refined representation of user  $\mathbf{p}_u$  and item  $\mathbf{q}_i$ , and keep the embeddings of context features unchanged. The role of the decoder is to output the prediction score by taking in the representations. The standard choice of decoder is multi-layer perceptron (MLP), which, however, falls short here since it only models feature interactions in an implicit way. In CARS, explicitly modeling the interactions between features is known to be important for user preference estimation [14]. For example, the classic factorization machine (FM) models the pairwise interactions between feature embeddings and has long been a competitive model for CARS.

Inspired by the simplicity (linear model) and the effectiveness of FM, we adopt it as the decoder of GCM. The idea is to explicitly model the pairwise interactions between the (refined) representations of user, item, and contexts with inner product. Specifically, let the set of vectors  $\mathcal{V}$  be  $\mathcal{V}_c \cup \{\mathbf{p}_u, \mathbf{q}_i\}$ , the decoder outputs the prediction score as:

$$\hat{y}_{uic} = \frac{1}{2} \left( \sum_{\mathbf{v}_s \in \mathcal{V}} \sum_{\mathbf{v}_t \in \mathcal{V}} \mathbf{v}_s^\top \mathbf{v}_t - \sum_{\mathbf{v}_s \in \mathcal{V}} \mathbf{v}_s^\top \mathbf{v}_s \right). \quad (10)$$

Here the self-interactions  $\mathbf{v}_s^\top \mathbf{v}_s$  are excluded since they are useless for the prediction. The bias terms for each user, item, and context feature are omitted for clarity.

Note that our FM-based decoder slightly differs from the vanilla FM, which models the interactions between the embeddings of all input features. Here we project each user (item) into a vector, rather than retaining the embeddings of her (its) features. An advantage is that this way abandons the internal interactions of user-field (item-field) features, shedding more light on the interactions between user (item) and context features, which is as expected.

#### 4.2 Model complexity analyses

We analyze the complexity of GCM from two aspects: the number of trainable parameters and the time complexity.

All trainable parameters come from the encoder layer, i.e., the embeddings of input features, since the GC layers and the decoder layer introduce no parameters to train. Let the feature number for the user field, item field, and context field as  $U, I$ , and  $C$ , respectively, and the embedding size be  $D$ . Then the embedding layer costs  $(U + I + C) \times D$  parameters. This demonstrates the low model complexity of GCM, since the number of trainable parameters is the same as FM — the simplest embedding-based CARS model.

For model training, since the complexity of the encoder plus the decoder is the same as that of FM, we analyze the additional time complexity caused by the GC layers. We implement the training in the batch-wise matrix form. Assume a batch contains all interactions. Then performing one GC layer takes time  $O((|\mathcal{Y}| + N + M)D)$ , where  $N$  and  $M$  denote the number of users and items, respectively. This complexity increases linearly with the number of GC layers.

After the model is trained, we perform one pass of GC layers to obtain the refined representations of all users and items, which can be done offline before online serving. As such, during online serving, we only need to execute the decoder, which has the same time complexity of FM. This is much faster than the recently emerging deep neural network-based CARS models like xDeepFM [16] and Convolutional FM [17]. Table 1 shows the model inference time of evaluating 1000 Yelp-OH users in which each interaction has 10 nonzero features of embedding size 64 and batch size is 4000. The testing platform is GeForce GTX 1080Ti with 16GB memory CPU. As can be seen, GCM takes a similar time as FM, being 24.5 and 157.7 times faster than xDeepFM and Convolutional FM, respectively.

#### 4.3 Optimization

We opt for the pointwise log loss to optimize model parameters, which is a common choice in recommender system [1, 16]. In each training epoch, we randomly sample non-observed interactions for each instance in  $\mathcal{Y}$  to form the negative set  $\mathcal{Y}^-$ . That is, for each observed instance  $(u, i, c) \in \mathcal{Y}$  of Yelp (or Amazon) dataset, we randomly match 4 (or 2) items from the item pool that user  $u$  has not interacted under context  $c$ . Then we minimize the following objective function:

$$L = - \sum_{(u,i,c) \in \mathcal{Y}} \log \sigma(\hat{y}_{uic}) - \sum_{(u,i,c) \in \mathcal{Y}^-} \log(1 - \sigma(\hat{y}_{uic})) + \lambda \|\Theta\|_2^2, \quad (11)$$

where  $\sigma(\cdot)$  is the sigmoid function,  $\lambda$  controls the  $L_2$  regularization to prevent over-fitting. The optimization is done by mini-batch Adam [41].

## 5 Experiments

We evaluate experiments on three benchmark datasets, aiming to answer the following research questions:

**Table 1** Model inference time of evaluating 1,000 Yelp-OH users (14 million interactions and 10 nonzero features per interaction)

Model	FM	GCM	GIN	xDeepFM	Convolutional FM
Time/s	8.51	14.93	35.45	365.82	2354.25

- **RQ1:** Compared with the state-of-the-art models, how does GCM perform w.r.t. top- $k$  recommendation?
- **RQ2:** How do different settings (e.g., depth of layer, modeling of context features, design of decoder) affect GCM?
- **RQ3:** How do the representation learning benefit from multiple interactions among users, items and contexts for item cold start issue?

## 5.1 Experimental settings

### 5.1.1 Dataset description

To demonstrate the effectiveness of GCM, we conduct experiments on three datasets from Yelp and Amazon, which are publicly available and vary in domain and size. We summarize the statistics of datasets in [Table 2](#).

- **Yelp:** This dataset is released by Yelp and records users' reviews on local businesses like bars and restaurants. In particular, we extract records happened in two different areas of USA — North Carolina, Ohio States — to construct datasets, termed Yelp-NC and Yelp-OH respectively.
- **Amazon:** Amazon review data is widely used in recommendation [3]. We select book subset from the collection in this work, and term it Amazon-book.

In what follows, we briefly introduce the features of users, items, and contexts. Specifically, for Yelp-NC and Yelp-OH, each user profile includes *yelping\_since\_year*<sup>2)</sup> and *average\_stars*, while the pre-existing features of items are composed of three attributes: *city*, *stars* and *is\_open*. We treat each review record as an observed instance, and collect *city*<sup>3)</sup>, *month*, *hour*, *day\_of\_the\_week* and *last\_purchase* as its context feature. For Amazon-book, the static features of items are composed of two attributes: *price* and *brand*. Similarly, each review record is treated as an observed instance, and *year*, *month*, *day*, *day\_of\_the\_week* and *last\_purchase* are collected as its context feature. Moreover, for all datasets, the 10-core setting is adopted to ensure data quality, i.e., retaining users with at least ten interactions.

For each user, we select the last interaction record to constitute the test set, while the remains are served as the training set. To emphasize model capability in recommending novel items for a user, we further filter the training set if the user-item pairs have appeared in the test set.

### 5.1.2 Evaluation metrics

In the evaluation phase, for each user in the test set, we view

**Table 2** Statistics of the datasets. We omit the ID feature when counting the number of user and item features

Dataset	Yelp-NC	Yelp-OH	Amazon-book
#User	6,336	5,170	44,709
#Item	13,003	12,997	46,831
#Instance	185,408	143,884	1,174,785
#User Feature	24	24	—
#Item Feature	68	213	24,816
#Context Feature	13,209	13,347	46,900

<sup>2)</sup> We only keep the *year* of the *yelping\_since* field which indicates the time the user joined Yelp.

<sup>3)</sup> The context feature *city* means which city does the interaction happen on. It is set as the city of the interacted item.

all items that she has not consumed before as recommendation candidates. Each method outputs a ranking list over the candidates. We then adopt two widely-used protocols to evaluate the quality of ranking lists: Hit Ratio (HR) and Normalized Discounted Cumulative Gain (NDCG). In particular,  $HR@K$  measures whether the test item is in the top- $K$  positions of the recommended list, whereas  $NDCG@K$  assigns higher scores to the top-ranked items. In our experiments, we report the results of  $K = 10$  and  $K = 50$ .

### 5.1.3 Baselines

We compare our GCM with several methods as follows:

- **MF** [42]: This exploits the user-item interactions only to learn user and item embeddings, while forgoing the context features.
- **LightGCN** [43]: Such model is the state-of-the-art GNN-based CF recommender, which incorporates high-order connectivity in user-item interaction graph into embeddings, while neglecting context features.
- **FM** [11]: This takes into account all information related to an interaction by converting all information to a feature vector then modeling second-order feature interaction to predict user preference.
- **NFM** [14]: This model leverages a MLP to capture nonlinear and high-order interaction among user, item, and context features.
- **xDeepFM** [16]: This is a recent neural FM model which combines explicit and implicit high-order feature interactions.
- **GIN** [32]: This is a graph-based model which mines user intention by applying implicit intention propagation and attention mechanism on a commodity similarity graph.

Fi-GNN [31] is a recent work on click-through rate prediction with graph neural networks, which is highly relevant to our work. It differs from GCM in graph construction — it builds a feature graph for each interaction, rather than the user-item graph. As a graph needs to be built for each interaction to obtain its prediction, the method is very slow in evaluation since all recommendation candidates need to be scored. As such, this method is not suitable for our all-ranking CARS evaluation, and we do not further compare with it. The Convolutional FM is not compared for the same reason (see [Table 1](#) for model inference time).

### 5.1.4 Parameter settings

We implement our GCM model and all baselines in Tensorflow. We apply the mini-batch Adam to optimize all models; the learning rate and batch size are set to 0.001 and 2048, respectively. A grid search is conducted for confirming optimal hyperparameters: for LightGCN, the number of gcn layers is searched in {1,2,3,4}; for NFM, the number of hidden layers is set to 1, the dropout rate is tuned in {0.9,0.8,...,0.1} for bi-interaction layer and hidden layer respectively; for xDeepFM, the number of cross layers is

searched in  $\{1, 2, 3\}$  with neuron number per layer in  $\{10, 20, 50, 100, 200\}$ , the number of DNN layers is same to that of cross layers, while the neuron number per layer is set to 100; for GIN, the length of previous records is 1 since we only keep the last purchase information in the datasets, the depth parameter is searched in  $\{1, 2, 3, 4\}$ , the number of neighbor nodes is tuned in  $\{10, 20\}$ , the neighbor is selected by the Top-N function according to the edge weight (for nodes with few neighbors, we randomly sample from unconnected nodes as their potential neighbors), a 5-layer full-connection perception with ReLU activation is adopted as the setting in [32]; for the proposed GCM, we search the model depth  $L$  amongst  $\{1, 2, 3\}$  with  $\alpha_l = 1/(L+1)$ , and adopt average pooling to generate the final refined representations of GC layers. For all models, the coefficient of  $L_2$  regularization term is searched in  $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$ . Moreover, we set the embedding size to 64 and train all models from the scratch.

## 5.2 Performance comparison (RQ1)

We report the empirical results of all models in Table 3 and have the following observations:

- Clearly, MF achieves the worst performance on three datasets, indicating that modeling user-item pairs as isolated instances limits the representation ability severely. LightGCN obtains consistent improvements over MF. We attribute such improvements to the modeling of user-item connectivity. However, neither MF nor LightGCN considers the context features, ignoring important factors and being insufficient for CARS.
- FM, NFM and xDeepFM consistently outperform MF and LightGCN across all cases. This is reasonable since they incorporate context features into the representation learning, so as to achieve better expressiveness and help to solve the data sparsity issue; Among them, NFM and xDeepFM perform better than FM by a large margin since they model more complex feature interactions: NFM employs MLP on user, item, and context features to capture their nonlinear and complex interactions, while xDeepFM learns high-order feature interactions in a more explicit way through a CIN network. This verifies that simply linear functions (e.g., inner product adopted by MF and LightGCN) might limit the representation learning and interaction modeling.

- GIN is the strongest baseline in all cases except for NDCG@10 and NDCG@50 in Yelp-NC. Such improvements are mainly because of GIN’s capability to model user intention by applying message propagation in the commodity similarity graph, which also verifies the necessity of bridging the relationship among data instances.
- GCM consistently outperforms all baselines w.r.t. all measures. In particular, GCM achieves noticeable improvements over the strongest baselines w.r.t. HR@10 by 20.78%, 14.93%, and 3.08%, in Yelp-NC, Yelp-OH, and Amazon-book, respectively. From  $t$ -test, we can find  $p$ -value  $< 0.05$  across the board, indicating that the improvements of GCM over the strongest baseline are statistically significant. We attribute such improvements to: 1) GCM employs the embedding propagation over the attributed graph to distill useful information from neighbors and connected edges, thus improving the representation ability; 2) Comparing with GIN which only propagates item embedding in the graph, GCM integrates the representations of users, items and contexts into the graph for information propagation, which may results in a more unified representations; and 3) Having established the refined representations, GCM further adopts FM to explicitly model the feature interactions.

## 5.3 Study of GCM (RQ2)

We next report ablation studies to verify the rationality of some designs in GCM, i.e., analyzing the influence of model depth, context modeling, normalization term, and decoder.

### 5.3.1 Impact of model depth

As GC is the core of GCM and stacking more GC layers is expected to augment the user and item representations with information propagated from multi-hop neighbors, we investigate how the number of GC layers affects the performance. In particular, we search the number of GC layers,  $L$ , in the range of  $\{0, 1, 2, 3\}$  and report the empirical results in Fig. 3.

We use GCM-1 to represent the model with one GC layer, and similar notations for others. We have several findings:

- GCM-0 disables the embedding propagation over the user-item attributed graph and downgrades to a FM-like linear model, thereby achieving poor performance. This again justifies the importance of GC layers.

**Table 3** Overall performance comparison. The bold indicates the best result, while the second-best performance is underlined

	Yelp-NC				Yelp-OH				Amazon-book			
	HR		NDCG		HR		NDCG		HR		NDCG	
	@10	@50	@10	@50	@10	@50	@10	@50	@10	@50	@10	@50
MF	0.0384	0.1173	0.0175	0.0341	0.0429	0.1261	0.0206	0.0383	0.0402	0.1243	0.0203	0.0382
LightGCN	0.0499	0.1394	0.0241	0.0431	0.0518	0.1520	0.0249	0.0461	0.0543	0.1466	0.0274	0.0473
FM	0.0739	0.1804	0.0396	0.0624	0.1959	0.4201	0.1049	0.1538	0.0587	0.1477	0.0323	0.0514
NFM	0.0824	0.2110	0.0419	0.0695	0.2248	0.4836	0.1161	0.1725	0.0808	0.1954	0.0444	0.0692
xDeepFM	0.0851	0.2086	<u>0.0458</u>	<u>0.0723</u>	0.2296	0.4799	0.1218	0.1762	0.0886	0.2119	0.0481	0.0748
GIN	<u>0.0866</u>	<u>0.2175</u>	0.0449	0.0722	<u>0.2304</u>	<u>0.4965</u>	<u>0.1238</u>	<u>0.1818</u>	<u>0.0939</u>	<u>0.2189</u>	<u>0.0502</u>	<u>0.0774</u>
GCM	<b>0.1046</b>	<b>0.2421</b>	<b>0.0557</b>	<b>0.0854</b>	<b>0.2648</b>	<b>0.5166</b>	<b>0.1457</b>	<b>0.2008</b>	<b>0.0968</b>	<b>0.2232</b>	<b>0.0536</b>	<b>0.0810</b>
Improv./%	20.78	11.31	21.62	18.12	14.93	4.05	17.69	10.45	3.08	1.96	6.77	4.65
$p$ -value	3.35e-9	4.37e-7	6.75e-10	5.91e-10	5.36e-12	9.10e-6	8.22e-9	2.45e-8	1.86e-4	3.41e-4	1.70e-3	1.05e-3

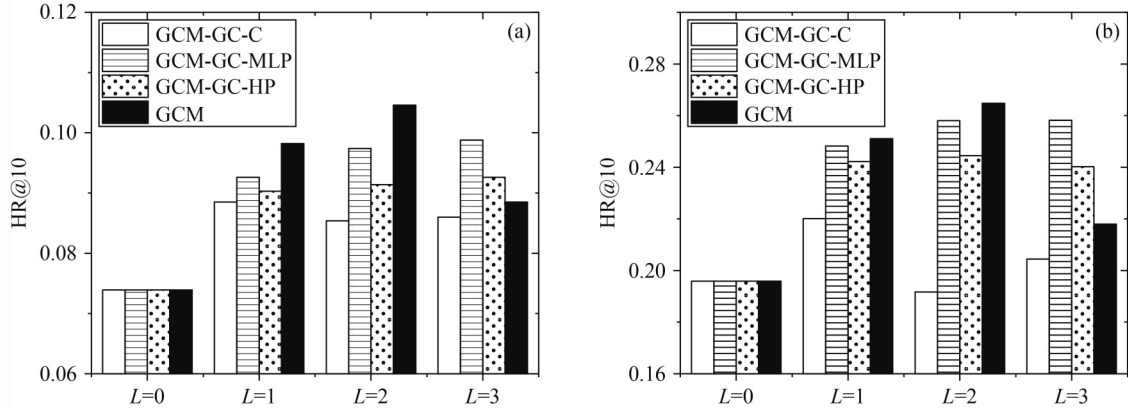


Fig. 3 The impact of depth and propagation rule in GC. (a) Yelp-NC; (b) Yelp-OH

- Obviously, increasing the number of GC layers results in better performance from  $L = 0$  to 2. In particular, GCM-2 performs better than GCM-1 in both datasets. It is reasonable since the signals passing from multi-hop neighbors (e.g., the second-order connectivity between behaviorally similar users or co-purchased items) are encoded into user and item representations of GCM-2, while GCM-1 only exploits personal history to enrich representations. This observation is consistent with that in NGCF [3]. We also tried to stack more GC layers (i.e., GCM-3), finding improvement degrades and over-smoothing issue. This suggests that GCM benefits most from the first- and second-order neighbors, but may suffer from degradation when higher-order neighbors are involved.

### 5.3.2 Impact of context modeling

One major contribution of GCM is to organize the context features as edges in the attributed user-item graph. We hence perform ablation studies, to demonstrate the rationality and effectiveness of this design. In particular, we build three different propagation rules for the GC layers of GCM: 1) GCM-GC-C removes the context features from the attributed graph and keeps the vanilla user-item interaction graph only; 2) GCM-GC-MLP first replaces the addition operation with concatenation in Eq. (3), then generates the message vector through a MLP; 3) GCM-GC-HP encodes the Hadamard product of the representations of neighboring node and their connected edge into the message during message passing. We show the comparison results in Fig. 3 and have the following observations.

- Modeling context features as the edges endows GCM with better generalization ability. In particular, GCM-GC-C performs worst among four competing methods in all cases, demonstrating the necessity of modeling context features when performing message propagation.
- GCM-GC-MLP consistently achieves better recommendation accuracy than GCM-GC-HP. One possible reason is that equipped with MLP, GCM-GC-MLP can model non-linear and high-order feature interactions, resulting in better representation ability.
- On both datasets, the best performance is always

achieved by 2-layers GCM, which again justifies the rationality of our context modeling strategy. Mathematically, GCM can be viewed as a special case of GCM-GC-MLP in which the learnt weights are the concatenation of two identity matrices. However, as more trainable parameters are involved, it would require more data to learn the function which may encounter difficulty in the actual learning process, especially for the notorious problem of data sparsity in recommendation system [44].

- Jointly analyzing Table 3 and Fig. 3, we find that GCM-GC-C without considering contexts achieves better performance than other baselines in Yelp-NC and comparable performance in Yelp-OH. This empirically suggests that propagating embeddings over interaction graphs is of importance to generate high-quality representations.

### 5.3.3 Impact of normalization term

For convenience, we only present the variants of the GC of the user side, since the same logic can be applied to the item side.

In GCM, we employ sqrt normalization term  $\frac{1}{\sqrt{|\mathcal{N}_u|}}$  on each neighbor embedding when performing neighborhood aggregation. To verify its rationality, we explore two different variants and report their empirical results here. The first variant uses symmetric normalization term, i.e.,  $\frac{1}{\sqrt{|\mathcal{N}_u|}\sqrt{|\mathcal{N}_i|}}$ , which is a common choice in GCN-based models [43], we term it GCM-sym. The other variant uses  $L_1$  normalization, i.e.,  $\frac{1}{|\mathcal{N}_u|}$ , we term it GCM- $L_1$ . Table 4 shows the results of the 2-layer GCM. We have the following observations:

- The best setting in general is using the sqrt normalization term on a single side (i.e., the current design of GCM). Adding additional regularization coefficients would significantly drop the performance.
- The smaller the normalization term, the worse the performance. To understand this observation, we can see the following inequalities:  $\frac{1}{\sqrt{|\mathcal{N}_u|}} > \frac{1}{\sqrt{|\mathcal{N}_u|}\sqrt{|\mathcal{N}_i|}} > \frac{1}{|\mathcal{N}_u|} = \frac{1}{|\mathcal{N}_u|}$ . The second inequality is due to



**Table 4** The variants of GCM with different normalization terms and decoders

	Yelp-NC		Yelp-OH	
	HR@10	NDCG@10	HR@10	NDCG@10
GCM	0.1046	0.0557	0.2648	0.1457
GCM- $L_1$	0.0810	0.0421	0.2373	0.1246
GCM-sym	0.0994	0.0527	0.2507	0.1383
GCM-APC	0.0947	0.0497	0.2321	0.1265
GCM-MLP	0.0892	0.0458	0.2263	0.1211
GCM-MF	0.0497	0.0253	0.0520	0.0251

$|\mathcal{N}_u| > |\mathcal{N}_i|$  on average in Yelp-NC and Yelp-OH.

### 5.3.4 Impact of decoder

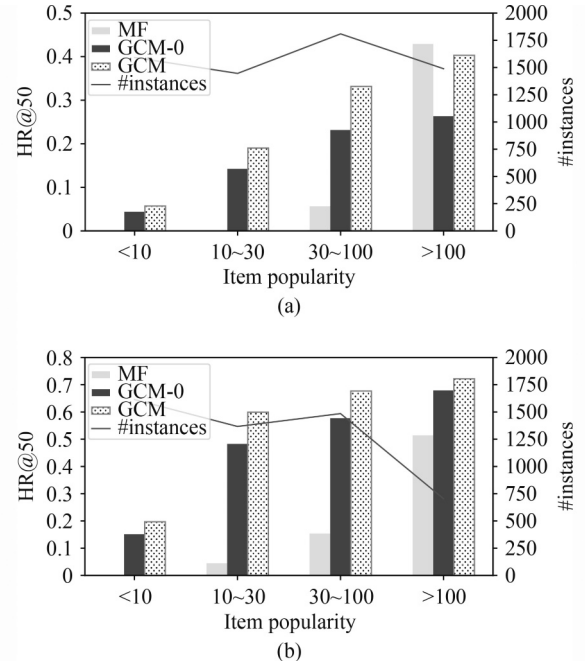
Having applied GC layers, we equip GCM with a decoder to model the pairwise interactions between the refined representations of users, items, and contexts. Here we investigate the role of such decoder. Towards this end, we compare GCM with three variants: 1) GCM-APC, which performs average pooling on context features before pair-wise inner product; 2) GCM-MLP, which replaces the decoder with MLP; and 3) GCM-MF, which replaces FM with inner product on user and item representations. Table 4 shows the comparison of results. There are several observations:

- Clearly, modeling feature interactions in the decoder enhances the predictive results. In particular, GCM, GCM-APC and GCM-MLP perform consistently better than GCM-MF, which relies only on the inner product of user and item representations.
- While having encoded context features into user and item representations via GC layer, GCM and GCM-APC highlight their influence in an explicit fashion, while GCM-MLP models the feature interactions in a rather implicit way. The better performances of GCM and GCM-APC again verify the rationality and effectiveness of FM-based decoder. In addition, after performing average pooling on context features, GCM-APC degenerates the performance by a noticeable margin; the reason is that GCM endows higher weights on feature interaction between context field and user (item) field, which is the core of CARS.

### 5.4 Performance w.r.t. item popularity (RQ3)

To alleviate the issue of the item cold start of CF models, taking side information into account is an auxiliary strategy that goes beyond modeling user-item interaction. In the proposed GCM, we apply gc layers to capture high-order connectivity on the user-item graph, which breaks down the independent interaction assumption of non-graph-based methods. We argue that such connectivity is a potential side information for the cold-start issue. To verify this viewpoint, we split the test set according to the popularity (the number of interaction records) of the target item, and report the performance of MF [42], GCM-0 and GCM in Fig. 4. We have the following observations:

- MF performs poorly at unpopular items, which indicates the item cold-start issue for CF models. GCM-0 has significant improvements in recommending uncommon items by introducing side information and modeling

**Fig. 4** Performances with respect to item popularity. (a) Yelp-NC; (b) Yelp-OH

feature interactions. Our GCM can further improve performance by 20%–30%. We attribute such improvements to modeling high-order connectivity since gnn increases the possibility of unpopular items being exposed through high-order links, thereby expanding the training data of unpopular items.

- For popular items, MF achieves comparable performance with GCM, even prevails over GCM in Yelp-NC. The possible reason is that the data of popular items occupies the majority of the training data, making MF adopt a cautious strategy — biased to recommending generally accepted items. Instead, GCM recommends items that are more niche but still consistent with user’s taste.
- The gain brought by gcn decreases as the popularity of items increases. This shows that as the number of neighbors increases, gcn may suffer from over-smoothing since these items have too many audiences, causing collecting information from neighbor nodes will also bring in noises.

## 6 Conclusion and future work

In this work, we emphasize the importance of exploiting multiple interactions in CARS. Towards this end, we first convert the features of users, items, and contexts into an attributed graph involving the contexts as edges between user and item nodes. We then develop a new model, GCM, which captures the interactions among multiple user behaviors via graph neural networks, and then models the interactions among features of individual behavior via factorization machine. To demonstrate the effectiveness of GCM, we test it on three public datasets, and it shows significant improvements over the state-of-the-art CF and CARS baselines. Extensive experiments also are conducted to verify the rationality of attributed graph and offer insights into how the

representations benefit from such graph learning.

Organizing user behaviors with contextual information in graphs is a promising direction to build an effective context-aware recommender. It helps build strong representations for users and items. However, GCM simply unifies all context features as an edge, neglecting the dynamic characteristics of some contexts (e.g., time) and hardly capturing the dynamic preference of users [45]. In the future, we plan to build dynamic graphs based on contextual information, instead of one static graph, and devise a dynamic graph neural network. Furthermore, rich side information is beneficial for explaining diverse intents behind user behaviors [46]. We hence plan to model user-item relationships at a granular level of user intents to generate disentangled representations [47].

**Acknowledgements** This work was supported by the National Key Research and Development Program of China (2020AAA0106000) and the National Natural Science Foundation of China (Grant Nos. 61972372, U19A2079, 62121002).

## References

1. He X, Liao L, Zhang H, Nie L, Hu X, Chua T S. Neural collaborative filtering. In: Proceedings of the 26th International Conference on World Wide Web. 2017, 173–182
2. Rendle S, Freudenthaler C, Gantner Z, Schmidt-Thieme L. BPR: Bayesian personalized ranking from implicit feedback. In: Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence. 2009, 452–461
3. Wang X, He X, Wang M, Feng F, Chua T. Neural graph collaborative filtering. In: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2019, 165–174
4. Kipf T N, Welling M. Semi-supervised classification with graph convolutional networks. In: Proceedings of the 5th International Conference on Learning Representations. 2017
5. Veličković P, Cucurull G, Casanova A, Romero A, Liò P, Bengio Y. Graph attention networks. In: Proceedings of the 6th International Conference on Learning Representations. 2018
6. Wei Y, Wang X, Nie L, He X, Hong R, Chua T. MMGCN: multi-modal graph convolution network for personalized recommendation of micro-video. In: Proceedings of the 27th ACM International Conference on Multimedia. 2019, 1437–1445
7. Wu L, Sun P, Fu Y, Hong R, Wang X, Wang M. A neural influence diffusion model for social recommendation. In: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2019, 235–244
8. Zheng L, Lu C T, Jiang F, Zhang J, Yu P S. Spectral collaborative filtering. In: Proceedings of the 12th ACM Conference on Recommender Systems. 2018, 311–319
9. Shi Y, Larson M, Hanjalic A. Collaborative filtering beyond the user-item matrix: a survey of the state of the art and future challenges. *ACM Computing Surveys*, 2014, 47(1): 1–45
10. Karatzoglou A, Amatriain X, Baltrunas L, Oliver N. Multiverse recommendation: n-dimensional tensor factorization for context-aware collaborative filtering. In: Proceedings of the 2010 ACM Conference on Recommender Systems. 2010, 79–86
11. Rendle S. Factorization machines. In: Proceedings of the 10th IEEE International Conference on Data Mining. 2010, 995–1000
12. Cheng H T, Koc L, Harmsen J, Shaked T, Chandra T, Aradhye H, Anderson G, Corrado G, Chai W, Ispir M, et al. Wide & deep learning for recommender systems. In: Proceedings of the 1st Workshop on Deep Learning for Recommender Systems. 2016, 7–10
13. Shan Y, Hoens T R, Jiao J, Wang H, Yu D, Mao J. Deep crossing: web-scale modeling without manually crafted combinatorial features. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016, 255–262
14. He X, Chua T. Neural factorization machines for sparse predictive analytics. In: Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2017, 355–364
15. Xiao J, Ye H, He X, Zhang H, Wu F, Chua T. Attentional factorization machines: learning the weight of feature interactions via attention networks. In: Proceedings of the 26th International Joint Conference on Artificial Intelligence. 2017, 3119–3125
16. Lian J, Zhou X, Zhang F, Chen Z, Xie X, Sun G. xdeepfm: combining explicit and implicit feature interactions for recommender systems. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2018, 1754–1763
17. Xin X, Chen B, He X, Wang D, Ding Y, Jose J. CFM: convolutional factorization machines for context-aware recommendation. In: Proceedings of the 28th International Joint Conference on Artificial Intelligence. 2019, 3926–3932
18. Rendle S, Gantner Z, Freudenthaler C, Schmidt-Thieme L. Fast context-aware recommendations with factorization machines. In: Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2011, 635–644
19. Guo H, Tang R, Ye Y, Li Z, He X. Deepfm: a factorization-machine based neural network for CTR prediction. In: Proceedings of the 26th International Joint Conference on Artificial Intelligence. 2017, 1725–1731
20. Wang R, Fu B, Fu G, Wang M. Deep & cross network for ad click predictions. In: Proceedings of the ADKDD'17. 2017, 1–7
21. Juan Y, Zhuang Y, Chin W, Lin C. Field-aware factorization machines for CTR prediction. In: Proceedings of the 10th ACM Conference on Recommender Systems. 2016, 43–50
22. Liu B, Tang R, Chen Y, Yu J, Guo H, Zhang Y. Feature generation by convolutional neural network for clickthrough rate prediction. In: Proceedings of the Web Conference. 2019, 1119–1129
23. Liu W, Tang R, Li J, Yu J, Guo H, He X, Zhang S. Fieldaware probabilistic embedding neural network for CTR prediction. In: Proceedings of the 12th ACM Conference on Recommender Systems. 2018, 412–416
24. Qu Y, Fang B, Zhang W, Tang R, Niu M, Guo H, Yu Y, He X. Product-based neural networks for user response prediction over multi-field categorical data. *ACM Transactions on Information Systems*, 2019, 37(1): 1–35
25. Hamilton W L, Ying R, Leskovec J. Inductive representation learning on large graphs. In: Proceedings of the 31st International Conference on Neural Information Processing. 2017
26. Gong L, Cheng Q. Exploiting edge features for graph neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019, 9211–9219
27. Zhang M, Chen Y. Link prediction based on graph neural networks. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. 2018, 5171–5181
28. Xu K, Hu W, Leskovec J, Jegelka S. How powerful are graph neural

- networks? In: Proceedings of the 7th International Conference on Learning Representations. 2019
29. Berg R, Kipf T N, Welling M. Graph convolutional matrix completion. 2017, arXiv preprint arXiv: 1706.02263
  30. Ying R, He R, Chen K, Eksombatchai P, Hamilton W L, Leskovec J. Graph convolutional neural networks for web-scale recommender systems. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2018, 974–983
  31. Li Z, Cui Z, Wu S, Zhang X, Wang L. Fi-gnn: modeling feature interactions via graph neural networks for ctr prediction. In: Proceedings of the 28th ACM International Conference on Information and Knowledge Management. 2019, 539–548
  32. Li F, Chen Z, Wang P, Ren Y, Zhang D, Zhu X. Graph intention network for click-through rate prediction in sponsored search. In: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2019, 961–964
  33. Wang W, Zhang W, Liu S, Liu Q, Zhang B, Lin L, Zha H. Beyond clicks: modeling multi-relational item graph for session-based target behavior prediction. In: Proceedings of the Web Conference. 2020, 3056–3062
  34. Fan W, Ma Y, Li Q, He Y, Zhao Y E, Tang J, Yin D. Graph neural networks for social recommendation. In: Proceedings of the Web Conference. 2019, 417–426
  35. Song W, Xiao Z, Wang Y, Charlin L, Zhang M, Tang J. Session-based social recommendation via dynamic graph attention networks. In: Proceedings of the 12th ACM International Conference on Web Search and Data Mining. 2019, 555–563
  36. Wu S, Tang Y, Zhu Y, Wang L, Xie X, Tan T. Sessionbased recommendation with graph neural networks. In: Proceedings of the 33rd AAAI Conference on Artificial Intelligence. 2019, 346–353
  37. Cao Y, Wang X, He X, Hu Z, Chua T. Unifying knowledge graph learning and recommendation: towards a better understanding of user preferences. In: Proceedings of the 2019 World Wide Web Conference. 2019, 151–161
  38. Wang X, He X, Cao Y, Liu M, Chua T. KGAT: knowledge graph attention network for recommendation. In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019, 950–958
  39. Mei L, Ren P, Chen Z, Nie L, Ma J, Nie J Y. An attentive interaction network for context-aware recommendations. In: Proceedings of the 27th ACM International Conference on Information and Knowledge Management. 2018, 157–166
  40. Wu C, Wu F, Qi T, Ge S, Huang Y, Xie X. Reviews meet graphs: enhancing user and item representations for recommendation with hierarchical attentive graph neural network. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing. 2019, 4886–4895
  41. Kingma D P, Ba J. Adam: a method for stochastic optimization. In: Proceedings of the 3rd International Conference on Learning Representations. 2015
  42. Koren Y, Bell R, Volinsky C. Matrix factorization techniques for recommender systems. *Computer*, 2009, 42(8): 30–37
  43. He X, Deng K, Wang X, Li Y, Zhang Y, Wang M. Lightgcn: simplifying and powering graph convolution network for recommendation. In: Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2020
  44. Rendle S, Krichene W, Zhang L, Anderson J R. Neural collaborative filtering vs. matrix factorization revisited. In: Proceedings of the 14th ACM Conference on Recommender Systems. 2020, 240–248
  45. Beutel A, Covington P, Jain S, Xu C, Li J, Gatto V, Chi E H. Latent cross: making use of context in recurrent recommender systems. In: Proceedings of the 11th ACM International Conference on Web Search and Data Mining. 2018, 46–54
  46. Wang X, He X, Feng F, Nie L, Chua T S. Tem: treeenhanced embedding model for explainable recommendation. In: Proceedings of the 2018 World Wide Web Conference. 2018, 1543–1552
  47. Ma J, Cui P, Kuang K, Wang X, Zhu W. Disentangled graph convolutional networks. In: Proceedings of the 36th International Conference on Machine Learning. 2019, 4212–4221



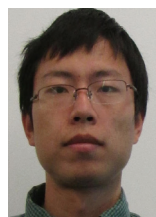
Jiancan Wu received his BS degree in Electronic Engineering and Information Science from the University of Science and Technology of China (USTC), China in 2017. He is currently a PhD student at USTC. His research interests focus on Recommender Systems, Machine Learning, and Graph Neural Networks.



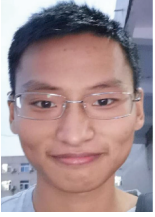
Xiangnan He is a professor at the University of Science and Technology of China (USTC), China. He received his PhD in Computer Science from the National University of Singapore (NUS), Singapore in 2016. His research interests span information retrieval, data mining, and multimedia analytics. He has over 70 publications that appeared in several top conferences such as SIGIR, WWW, and MM, and journals including TKDE, TOIS, and TMM. His work on recommender systems has received the Best Paper Award Honorable Mention in WWW 2018 and ACM SIGIR 2016. Moreover, he has served as the PC chair of CCIS 2019, area chair of MM (2019, 2020) ECML-PKDD 2020, and PC member for several top conferences including SIGIR, WWW, KDD, WSDM, etc., and the regular reviewer for journals including TKDE, TOIS, TMM, etc.



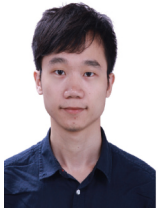
Xiang Wang is now a research fellow at National University of Singapore, Singapore. He received his PhD degree from National University of Singapore, Singapore in 2019. His research interests include recommender systems, graph learning, and deep learning techniques. He has published some academic papers on international conferences such as KDD, WWW, SIGIR, and AAAI. He serves as a program committee member for several top conferences such as SIGIR and WWW.



Qifan Wang received the BS and MS degrees from Tsinghua University, China, and the PhD degree from Purdue University, USA, all in computer science. He is a researcher in Google Research, and his research interests include machine learning, information retrieval, data mining, computer vision, and natural language processing.

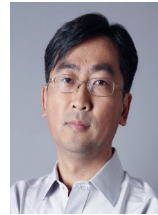


Weijian Chen is currently a PhD student at the University of Science and Technology of China (USTC), China. His research interests focus on User Profiling, Recommender Systems, and Graph Neural Networks. He has served as a research intern in the JD Data Science Laboratory and a project intern in the Kwai Multimedia Understanding Department.



Jianxun Lian is now a senior researcher at Microsoft Research Asia. He received his PhD degree from University of Science and Technology of China, China in 2018. His research interests include recommender systems and deep learning techniques. He has published some academic papers on international conferences such as KDD, IJCAI, WWW, SIGIR, and CIKM. He serves as a program

committee member for several top conferences such as AAAI, WWW, and IJCAI.



Xing Xie is currently a senior principal research manager at Microsoft Research Asia, and a guest PhD advisor at the University of Science and Technology of China, China. He received his BS and PhD degrees in Computer Science from the University of Science and Technology of China, China in 1996 and 2001, respectively. He joined Microsoft Research Asia in July 2001, working on data mining, social computing and ubiquitous computing. During the past years, he has published over 300 referred journal and conference papers, won the 10-year impact award in ACM SIGSPATIAL 2019, the best student paper award in KDD 2016, and the best paper awards in ICDM 2013 and UIC 2010.