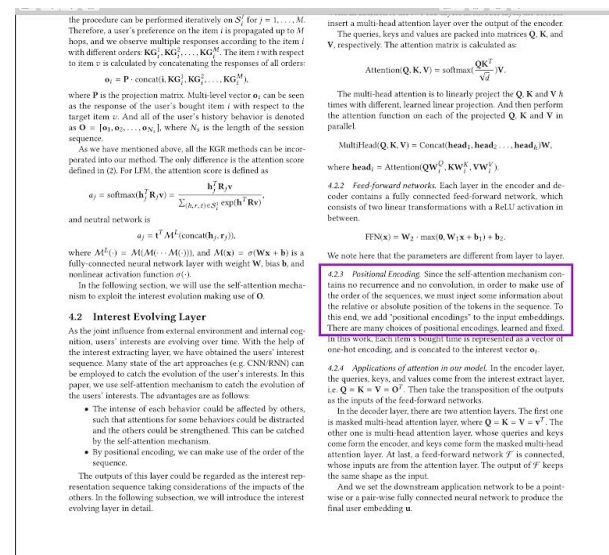


Reviewing paper is on the left, original paper is on the right

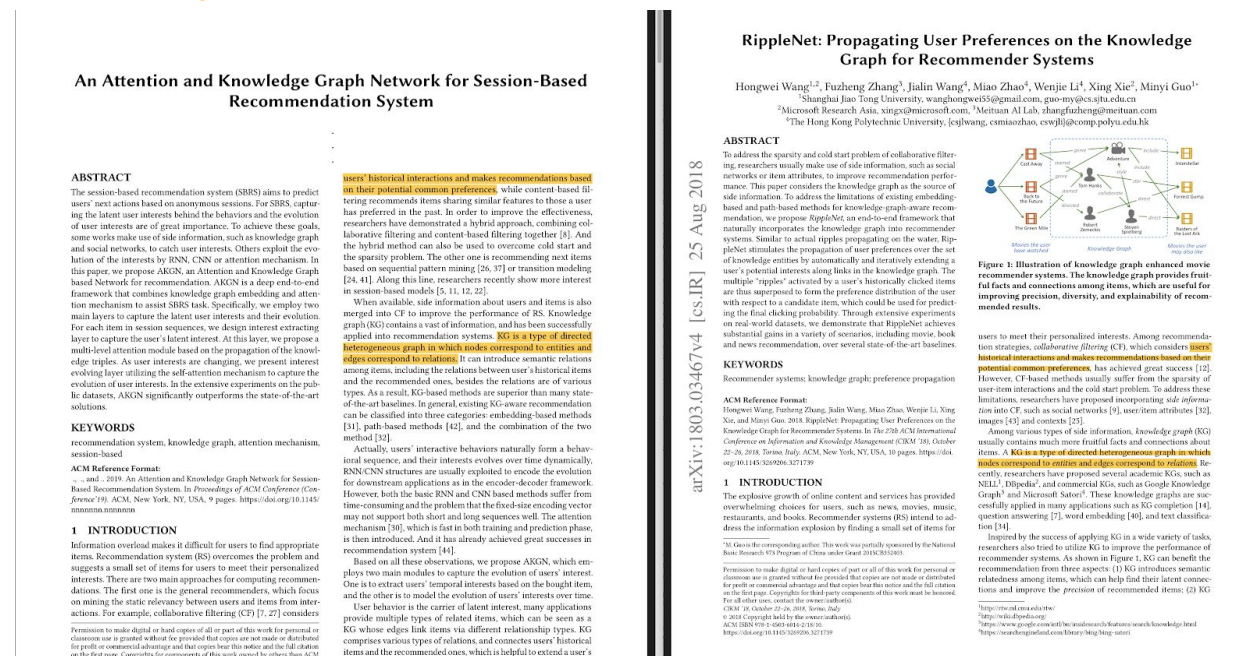
ATRank

<https://arxiv.org/pdf/1711.06632.pdf>



RippleNet

<https://arxiv.org/pdf/1803.03467.pdf>



- the proposed model uses a novel knowledge graph embedding and attention mechanism to assist SBSRS task.
- Given that the users' attention varies on different types of relations and the users' preferences evolve over time, we employ interest extracting layer and interest evolving layer to capture the latent user interests and their evolution.
- We conduct experiments on real-world recommendation scenarios, and the results prove the efficacy of AKGN over several state-of-the-art baselines.

2 RELATED WORKS

2.1 Knowledge Graph Representation

Our work connects to a large body of work in Knowledge Graph Representation (KGR) methods. KGR intends to embed entities and relations in a KG into low-dimensional continuous vector spaces while preserving the KG's structural information. Researchers have proposed many models to embed entities and relations in KGs [28], among which the most common way is translational distance models, such as TransE [37], TransH [36], TransR [37] and TransD [14] and so on. This kind of models exploit distance-based scoring functions when learning representations of entities and relations. For example, TransE [37] wants $h + r \approx t$ when (h, r, t) holds, where h, r, t are the corresponding representation vectors of head, relation and tail. Therefore, TransE assumes the score function $f_e(h, r, t) = \|h + r - t\|_2^2$ is low if (h, r, t) holds, and high otherwise. In TransH [36], each object is represented by two vectors. Given a triple (h, r, t) , its vectors are h_p, r_p, t_p . And the score function is $f_e(h, r, t) = \|h_p + r_p - t_p\|_2^2$. Another commonly used method is Latent Factor Model (LFM) [13, 29]. LFM employs a relation-specific bilinear form to consider the relationships between entities and relations, and the score function is defined as $f_e(h, r, t) = \mathbf{M}^T \mathbf{M} \mathbf{A}^T \text{conc}(\mathbf{h}, \mathbf{r})$, where $\mathbf{M}^T \mathbf{M} = \mathbf{M} \mathbf{A} \mathbf{M}^T = \mathbf{M} \mathbf{G}$, and $\mathbf{M} \mathbf{G} = \mathbf{e}(\mathbf{p}_h + \mathbf{p}_t)$ is a fully-connected neural network with weight \mathbf{W} , bias \mathbf{b} , and nonlinear activation function $\sigma(\cdot)$. We note here that all of the above KGR methods can be incorporated into our method.

2.2 Attention Mechanism

Attention mechanism was proposed in image classification [21] and machine translation [2], which aims to find the most relevant

information performance of all the tasks [38, 40]. All of the learning tasks or at least a subset of them are assumed to be related to each other, and learning these tasks jointly can lead to performance improvement compared with learning them individually. In general, MTL algorithms can be classified into several categories, including feature learning approach [35, 38], low-rank approach [16, 20], task clustering approach [46], task relation learning approach [45], and decomposition approach [3]. MTL can also be combined with other learning paradigms to improve the performance of learning tasks further, including semi-supervised learning, active learning, unsupervised learning and reinforcement learning. In this paper, we aim to utilize the connection of RS and KG to help improve their performance.

3 PROBLEM FORMULATION

Let $\mathcal{U} = \{u_1, u_2, \dots\}$ and $\mathcal{I} = \{i_1, i_2, \dots\}$ denote the sets of users and items, respectively. The user-item interaction matrix $\mathbf{Y} = \{y_{ui} | u \in \mathcal{U}, i \in \mathcal{I}\}$ is defined according to users' implicit feedback, where

$$y_{ui} = \begin{cases} 1 & \text{if there is an interaction between } u \text{ and } i, \\ 0 & \text{otherwise.} \end{cases}$$

In addition to the interaction matrix \mathbf{Y} , we also have the category id of each item and several sources of heterogeneous item-to-item relationships $\mathcal{R} = \{r_1, r_2, \dots, r_n\}$, which can be represented by a KG, denoted by \mathcal{G} .

Given interaction matrix \mathbf{Y} as well as sources of heterogeneous item-to-item relationships \mathcal{R} , we aim to predict whether user u has a potential interest in item i which he has had no interaction before. Our goal is to learn a prediction function $\hat{y}_{ui} = f(u, \mathbf{v}, \mathbf{0})$, where \hat{y}_{ui} denotes the probability that user u will click item i , and $\mathbf{0}$ denotes the model parameters of function f .

4 AKGN MODELING FRAMEWORK

In this section, we discuss the proposed method in detail. The framework of our method is illustrated in Figure 1. We take a user u and an item i as inputs, and output the predicted probability that user u will click item i .

For each item in session sequences, we design interest extracting layer to capture the user's latent interest. At this layer, we propose a multi-level attention module based on the propagation of the knowledge triples, and return a vector to denote the response of

Network (DKN) [33] treats entity embeddings and word embeddings as different channels, then designs a CNN framework to combine them together for news recommendation. Collaborative Knowledge base Embedding (CKE) [43] combines a CT module with knowledge embedding, text embedding, and image embedding of items in a unified Bayesian Framework. Signed Heterogeneous Information Network Embedding (SHINE) [42] designs deep autoencoders to embed sentiment networks, social networks and profile knowledge networks for celebrity recommendations. Embedding-based methods show high flexibility in utilizing KG to assist recommender systems, but the adopted KGR algorithms in these methods are usually more suitable for in-graph applications such as link prediction than for recommendation [35], thus the learned entity embeddings are less intuitive and effective to characterize inter-item relations.

The second category is path-based methods [42, 45], which explore the various patterns of connections among items in KG to provide additional guidance for recommendations. For example, Personalized Entity Recommendation (PER) [42] and Meta-Graph Based Recommendation [45] treat KG as a heterogeneous information network (HIN), and extract meta-path/meta-graph based latent features to represent the connectivity between users and items along different types of relation paths. Path-based methods make use of KG in a more natural and intuitive way but they rely heavily on manually designed meta-paths, which is hard to optimize in practice. Another concern is that it is impossible to design hand-crafted meta-paths in certain scenarios (e.g., news recommendation) where entities and relations are not within one domain.

To address the limitations of existing methods, we propose RipleNet, an end-to-end framework for knowledge-graph-aware recommendation. RipleNet is designed for click-through rate (CTR) prediction, which takes a user-item pair as input and outputs the probability of the user engaging (e.g., clicking, browsing) the item. The key idea behind RipleNet is preference propagation. For each user, RipleNet treats his historical interests as a seed set in the KG, then extends the user's interests iteratively along KG links to discover his hierarchical potential interests with respect to a candidate item. We analyze preference propagation with actual ripples created by random propagating on the water, in which multiple "ripples" suppose to be a result of preferential distribution of the user over the knowledge graph. The major difference between RipleNet and existing literature is that RipleNet combines the advantages of the two mentioned types of methods: (i) RipleNet incorporates the KG methods into recommendation

the knowledge graph.

In summary, our contributions in this paper are as follows:

- To the best of our knowledge, this is the first work to combine embedding-based and path-based methods in KG-aware recommendation.
- We propose RipleNet, an end-to-end framework utilizing KG to assist recommender systems. RipleNet automatically discovers user's hierarchical potential interests by iteratively propagating users' preferences in the KG.
- We conduct experiments on three real-world recommendation scenarios, and the results prove the efficacy of RipleNet over several state-of-the-art baselines.

2 PROBLEM FORMULATION

The knowledge-graph-aware recommendation problem is formulated as follows. In a typical recommender system, let $\mathcal{U} = \{u_1, u_2, \dots\}$ and $\mathcal{I} = \{i_1, i_2, \dots\}$ denote the sets of users and items, respectively. The user-item interaction matrix $\mathbf{Y} = \{y_{ui} | u \in \mathcal{U}, i \in \mathcal{I}\}$ is defined according to users' implicit feedback, where

$$y_{ui} = \begin{cases} 1 & \text{if interaction } (u, i) \text{ is observed;} \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

A value of 1 for y_{ui} indicates there is an implicit interaction between user u and item i , such as behaviors of clicking, watching, browsing, etc. In addition to the interaction matrix \mathbf{Y} , we also have a knowledge graph \mathcal{G} available, which consists of massive 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, respectively. \mathcal{E} and \mathcal{R} denote the set of entities and relations in the KG. For example, the triple *Jurassic Park, film, film director, Steven Spielberg* states the fact that Steven Spielberg is the director of the film "Jurassic Park". In many recommendation scenarios, an item $i \in \mathcal{I}$ may associate with one or more entities in \mathcal{E} . For example, the item "Jurassic Park" is linked with its namesake in KG, while news with title "France's Baby Panda Makes Public Debut" is linked with entities "France" and "Panda".

Given interaction matrix \mathbf{Y} as well as knowledge graph \mathcal{G} , we aim to predict whether user u has a potential interest in item i with which he has had no interaction before. Our goal is to learn a prediction function $\hat{y}_{ui} = f(u, \mathbf{v}, \mathbf{0})$, where \hat{y}_{ui} denotes the probability that user u will click item i , and $\mathbf{0}$ denotes the model parameters of function f .

Finally, the user embedding \mathbf{u} and the target item embedding \mathbf{v} are combined to predict the clicking probability:

$$\hat{y}_{ui} = \sigma(f(\mathbf{u}, \mathbf{v})),$$

where $\sigma(\cdot)$ is the sigmoid function, and f is a ranking function which can be either a dot-product form or a more complex deep neural network.

The loss function \mathcal{L} is the sigmoid cross entropy loss, which is defined as:

$$-\sum_{(u,i) \in \mathcal{Y}} (y_{ui} \log \sigma(f(\mathbf{u}, \mathbf{v})) + (1 - y_{ui}) \log(1 - \sigma(f(\mathbf{u}, \mathbf{v})))) \quad (3)$$

4.3 Multi-Task Learning

We use a multi-task learning framework [5] to jointly optimize all two tasks using shared variables.

$$\begin{aligned} \min \mathcal{L} = & \min (\mathcal{L}_{\text{AKGN}}) \\ = & \sum_{(u,i) \in \mathcal{Y}} -(y_{ui} \log \sigma(f(\mathbf{u}, \mathbf{v})) + (1 - y_{ui}) \log(1 - \sigma(f(\mathbf{u}, \mathbf{v})))) \\ & + \lambda_1 \left(\sum_{(h,r,t) \in \mathcal{G}} f_e(h, r, t) + \sum_{(u',r',t') \in \mathcal{G}} f_e(u', r', t') \right) \\ & + \lambda_2 (\|\mathbf{E}\|_F^2). \end{aligned} \quad (4)$$

In eq(4), the first term measures loss in the recommendation module, where \mathbf{u} and \mathbf{v} traverse the set of users and the items, respectively. The second term measures the squared error between the ground truths of the KG and the reconstructed one, $f_e(h, t)$ is the score function of TransE. And the third term is the regularizer for preventing over-fitting.

We employ a stochastic gradient descent (SGD) algorithm to iteratively optimize the loss function. In each training iteration, to make the computation more efficient, we randomly sample a minibatch of positive/negative interactions from \mathbf{Y} and true/false triples from \mathcal{G} . Then we calculate the gradients of the loss \mathcal{L} with respect to model parameters, and update all parameters by back-propagation based on the sampled minibatch. We will discuss the choice of hyper-parameters in the experiments section.

5 EXPERIMENT

In this section, we perform experiments on the real-world datasets to demonstrate the efficiency and quality of our proposed methods.

5.1 Experimental Settings

5.1.1 Datasets and Preprocessing. Amazon Dataset is composed of large corpora of users' reviews on products and multiple types of related items [19]. These data were crawled from Amazon.com and span May 1996 to July 2014. We use three subsets of Amazon dataset: Books, Electronics and Automotive, which have already been reduced to satisfy the 5-core property, such that each of the remaining users and items have 5 reviews each. The feature we use contains user id, item id, cate id, related items and timestamps. This set of data is not able for its high sparsity and variability. Amazon contains four types of item relationships that we use in our model: "also viewed", "also bought", "bought together", and "buy after viewing". To collect the related facts from Amazon, we only consider those triples that are directly related to the entities with mapped items. We then preprocess the four datasets by: filtering low-frequency and high-frequency users and items (i.e., lower than 5 and longer than 20 in both datasets), filtering out unrelated entities (i.e., lower than 5 in both datasets), cutting off infrequent relations and merging similar relations manually. In these datasets, we regard reviews as behaviors, and sort the reviews from one user by time. Assuming there are T behaviors of user u , our purpose is to use the $T-1$ behaviors to predict whether user u will write reviews that shown in T -th review. The statistics of the four datasets are shown in Table 1.

5.1.2 Baseline Algorithms. To evaluate the performance of the proposed method, we compare it with the following representative baselines:

- BPR-MF [25]: It is one of widely used matrix factorization methods, which optimizes the latent factor model with implicit feedback using a pairwise ranking loss in a Bayesian manner.
- Bi-LSTM: We implement a Bi-LSTM method to encode the user behavior history, whose difference with [39] is that we use a bidirectional LSTM as the encoder because of better performance in our experiments. Then we concatenate the two final hidden states to the user embedding. The stacked LSTM depth is set to be 1.
- CNN-Pooling: We use a CNN structure with max pooling to encode the user history as in [15, 43]. We use ten kinds of window sizes from one to ten for extracting different features, and all of the feature map have the same kernel size with 32. Then we apply a max-over-time pooling operation over the feature map, and pass all pooled features from different maps to a fully connected layer to produce final user embedding.

Based on the definition in Eq. (1), the scoring functions of entity-entity pairs in KG and item-entity pairs in preference propagation can thus be unified under the same calculation model. The last term in Eq. (9) is the likelihood function of the observed implicit feedback given Θ and the KG, which is defined as the product of Bernoulli distributions:

$$p(\mathbf{Y} | \Theta, \mathcal{G}) = \prod_{(u,i) \in \mathcal{Y}} \pi(u^T \mathbf{v}_i)^{y_{ui}} (1 - \pi(u^T \mathbf{v}_i))^{1-y_{ui}} \quad (12)$$

based on Eq. (2)-(7).

Taking the negative logarithm of Eq. (9), we have the following loss function for RipleNet:

$$\begin{aligned} \min \mathcal{L} = & -\log(p(\mathbf{Y} | \Theta, \mathcal{G})) - p(\mathcal{G} | \Theta) - p(\Theta) \\ = & \sum_{(u,i) \in \mathcal{Y}} (-y_{ui} \log \pi(u^T \mathbf{v}_i) + (1 - y_{ui}) \log(1 - \pi(u^T \mathbf{v}_i))) \\ & + \lambda_2 \sum_{r \in \mathcal{R}} (\|L_r - E^T RE\|_F^2 + \frac{\lambda_1}{2} (\|W_r\|_F^2 + \|E_r\|_F^2) + \sum_{t \in \mathcal{E}} \|R_{rt}^2\|_F^2) \end{aligned} \quad (13)$$

where \mathbf{W} and \mathbf{E} are the embedding matrices for all items and entities, respectively. L_r is the slice of the indicator tensor \mathbf{I} in KG for relation r , and \mathbf{R} is the embedding matrix of relation r . In Eq. (13), the first term measures the cross-entropy loss between ground truth of interactions \mathbf{Y} and predicted value by RipleNet, the second term measures the squared error between the ground truth of the KG \mathcal{G} and the reconstructed indicator matrix $E^T RE$, and the third term is the regularizer for preventing over-fitting.

It is intractable to solve the above objection directly, therefore, we employ a stochastic gradient descent (SGD) algorithm to iteratively optimize the loss function. The learning algorithm of RipleNet is presented in Algorithm 1. In each training iteration, to make the computation more efficient, we randomly sample a minibatch of positive/negative interactions from \mathbf{Y} and true/false triples from \mathcal{G} following the negative sampling strategy in [16]. Then we calculate the gradients of the loss \mathcal{L} with respect to model parameters Θ , and update all parameters by back-propagation based on the sampled minibatch. We will discuss the choice of hyper-parameters in the experiments section.

present a visualized example in the experiments section to directly demonstrate the explainability of RipleNet.

3.5.2 Ripple Superposition. A common phenomenon in KG is that a user's ripple sets may be large in size, which dilutes potential interests inevitably in preference propagation. I observe that relevant entities of different items in a user history often highly overlap. In other words, an entity reached by multiple paths in the KG starting from a user history. For example, "Saving Private Ryan" is connected with who has watched "The Terminal", "Jurassic Park" and "Ben Hur" through actor "Tom Hanks", director "Steven Spielberg" or "War", respectively. These parallel paths greatly increase interests in overlapped entities. We refer to the case as *ripple superposition*, as it is analogous to the interference phenomenon in which two waves superpose to form a resultant greater amplitude in particular areas. The phenomenon superposition is illustrated in the second KG in Figure 2, \mathbf{x} darker red around the two lower middle entities indicate strength of the user's possible interests. We will also discuss superposition in the experiments section.

3.6 Links to Existing Work

Here we continue our discussion on related work and comparisons with existing techniques in a greater scope.

3.6.1 Attention Mechanism. The attention mechanism was firstly proposed in image classification [18] and machine translation [1], which aims to learn where to find the most relevant input automatically as it is performing the task. The notion transplanted to recommender systems [4, 22, 33, 36, example, DADM [4] considers factors of specialty and de-assigning attention values to articles for recommendation [22] proposes an interpretable and dual attention-based CN that combines review text and ratings for product rating. DKN [33] uses an attention network to calculate the between a user's clicked item and a candidate item to dynamically aggregate the user's historical interests. RipleNet can be

Attention

<https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf?abstract=techstories.org>

by time. Assuming there are T behaviors of user u , our purpose is to use the $T-1$ behaviors to predict whether user u will write reviews that shown in T -th review. The statistics of the four datasets are shown in Table 1.

5.1.2 Baseline Algorithms. To evaluate the performance of the proposed method, we compare it with the following representative baselines:

- BPR-MF [25]: It is one of widely used matrix factorization methods, which optimizes the latent factor model with implicit feedback using a pairwise ranking loss in a Bayesian manner.
- Bi-LSTM: We implement a Bi-LSTM method to encode the user behavior history, whose difference with [39] is that we use a bidirectional LSTM as the encoder because of better performance in our experiments. Then we concatenate the two final hidden states to the user embedding. The stacked LSTM depth is set to be 1.
- CNN-Pooling: We use a CNN structure with max pooling to encode the user history as in [15, 43]. We use ten kinds of window sizes from one to ten for extracting different features, and all of the feature map have the same kernel size with 32. Then we apply a max-over-time pooling operation over the feature map, and pass all pooled features from different maps to a fully connected layer to produce final user embedding.

feature that represent the timestamp when that behavior happens. The statistics of each category is shown in table 2. This dataset is guaranteed that each user has at least 3 behaviors in each behavior group.

Competitors

- BPR-MF: Bayesian Personalized Ranking (Rendle et al. 2009) is a pairwise ranking framework. It trains on pairs of a user's positive and negative examples, maximizing the posterior probability of the difference given the same user. The item vector is a concatenation of item embedding and category embedding, each of which has a dimension of 64, the user vector has a dimension of 128.
- Bi-LSTM: We implement a Bi-LSTM method to encode the user behavior history, whose difference with (Zhang et al. 2014) is that we use a bidirectional LSTM as the encoder because of better performance in our experiments. Then we concatenate the two final hidden states to the user embedding. The stacked LSTM depth is set to be 1.
- Bi-LSTM+Attention: We add a vanilla attention on top of the Bi-LSTM method mentioned above.
- CNN-Pooling: We use a CNN structure with max pooling to encode the user history as in (Zhang, Necozi, and Ye 2017; Kim 2014). We use ten kinds of window sizes, from one to ten for extracting different features, and all of the feature map have the same kernel size with 32. Then we apply a max-over-time pooling operation over the feature map, and pass all pooled features from different maps to a fully connected layer to produce final user embedding.

Table 4 illustrates the average vanilla-attention score for different time buckets over the whole amazon dataset, which can be inferred that time encoding via self-attention can