# **Explainable Knowledge Graph-based Recommendation via Deep Reinforcement Learning**

Weiping Song\*1, Zhijian Duan1, Ziqing Yang1, Hao Zhu1, Ming Zhang1, Jian Tang2,3,4

1 Peking University, 2 Mila - Quebec AI Institute,

3 CIFAR AI Research Chair, 4 HEC Montréal

{weiping.song,zjduan,yangziqing,hzhu1998,mzhang\_cs}@pku.edu.cn

jian.tang@hec.ca

### **Abstract**

This paper studies recommender systems with knowledge graphs, which can effectively address the problems of data sparsity and cold start. Recently, a variety of methods have been developed for this problem, which generally try to learn effective representations of users and items and then match items to users according to their representations. Though these methods have been shown quite effective, they lack good explanations, which are critical to recommender systems. In this paper, we take a different path and propose generating recommendations by finding meaningful paths from users to items. Specifically, we formulate the problem as a sequential decision process, where the target user is defined as the initial state, and the walks on the graphs are defined as actions. We shape the rewards according to existing state-of-the-art methods and then train a policy function with policy gradient methods. Experimental results on three real-world datasets show that our proposed method not only provides effective recommendations but also offers good explanations .

# 1 Introduction

Recommender systems are essential to a variety of online applications such as e-Commerce Websites and social media platforms by providing the right items or information to the users. One critical problem of recommender systems is data sparsity, i.e., some items are purchased, rated, or clicked by only a few users or no users at all. Recently, there is an increasing interest in knowledge graph-based recommender systems, since knowledge graphs can provide complementary information to alleviate the problem of data sparsity and have been proved quite useful [3, 15, 20, 21, 30, 32].

Generally, existing knowledge graph-based recommendation methods try to learn effective representations of users and items according to the user-item interaction graphs and item-entity knowledge graphs, and then match the items to the users according to learned representations. For example, Zhang et al. [30] learn item representations by combining their representations in user-item graphs and knowledge-graphs. Zhang et al. [31] learn user and item representations on the integrated user-item-entity graphs based on knowledge graph embedding (KGE) method like TransE [2]. Wang et al. [22] and Cao et al. [3] jointly optimize the recommendation and knowledge graph embedding tasks in a multi-task learning setting via sharing item representations. These methods have been proved quite effective by integrating information from both user behaviors and knowledge graphs.

Although these methods are very effective, they lack good explanations. Intuitively, if the recommender system can give an explanation of a recommendation, the users would have more interest and trust in the recommended item [9, 14, 32, 33]. Indeed, there is some existing work that aims to

<sup>\*</sup>This work was done when the first author was visiting Mila.

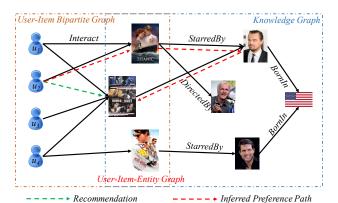


Figure 1: Explainable recommendation via reasoning over the integrated useritem-entity graph. Our Ekar generates the recommendation (i.e., "Romeo and Juliet") by inferring a preference path "u<sub>2</sub>  $\xrightarrow{Interact}$   $\xrightarrow{StarredBy}$  Titanic  $\xrightarrow{StarredIn}$  Romeo and Juliet".

provide such explanations for the recommendation results. For example, the RippleNet [20] aims to explain the recommendations by analyzing the attention scores. However, their method relies on the post analysis of the soft attention scores, which may not always be trustworthy.

In this paper, we take a different route and propose to generate a path from the target user to relevant items in the integrated user-item-entity graph. Take the movie recommendation in Figure 1 as an example. For the target user  $u_2$ , a path (the red dashed lines) is generated to the item "Romeo and Juliet" since 1) user  $u_2$  watched movie "Titanic"; 2) "Titanic" is starred by "Leonardo DiCaprio"; and 3) "Leonardo DiCaprio" also stars in "Romeo and Juliet". We can see that such a path offers good explanations of the recommendations in addition to provide meaningful recommendations.

However, finding meaningful paths on the large user-item-entity graph is challenging. One may enumerate all paths between user-item pairs and then use a classification/ranking model to select the most meaningful paths [24]. Nevertheless, enumerating paths between users and items is intractable due to the exponentially large path space in the user-item-entity graph. Although sampling a certain number of paths via breadth-first-search can be a practical substitution to enumerating, it has no assurance on the meaningfulness of sampled paths. In this paper, we formulate the generation of meaningful user-to-item paths as a sequential decision process. Specifically, the recommender agent starts from target users and extends its paths to relevant items by sequentially selecting walks on the user-item-entity graph. During training, we assign each path a positive reward if the starting user and terminal entity constitute an observation in recommendation data. Considering that the reward could be extremely sparse at the beginning due to the huge exploration space, we further augment the reward signals by reward shaping [17], where a soft reward function is first learned using state-of-the-art knowledge graph embedding methods. We use the REINFORCE [25] algorithm to maximize the expected rewards of our recommender agent. Finally, we verify the effectiveness and the explainability of the proposed method on three real-world datasets. Quantitative results demonstrate that our proposed Ekar: 1) significantly outperforms existing state-of-the-art KG-based recommendation methods, 2) offers clear and convincing explanations in the form of meaningful paths from users to recommended items.

## 2 Related Work

Our work is conceptually related to the explainable recommendation, knowledge graph-based recommendation, and recent advancements in applying reinforcement learning into relational reasoning.

**Explainable Recommendation**. As a widespread concern in the AI community, explainability has been widely discussed in recommender systems. According to Zhang and Chen [32], most of the existing explainable recommendation methods typically provide explanations via identifying users' preference on item features [1, 23, 33], understanding latent factors with topic modeling [16, 26, 35], or ranking over the user-item-aspect graph [8]. However, these methods require external information (e.g., reviews) about items, which may be difficult to collect. Some recent advancements utilize the attention mechanism [12, 19] to provide explanations, but they need extra efforts to explore attention scores. Since knowledge graphs (KG) provide common knowledge about our world, many recent works use knowledge graphs [3, 20, 21, 30] to provide explainable recommendations, which will be further discussed in the following paragraph.

**KG-based Recommendation**. Our work is closely related to KG-based recommendation, which utilizes general knowledge graphs (e.g., DBpedia, YAGO, and Satori) to improve recommender systems. Existing KG-based recommendation methods can be roughly divided into two classes: embedding-based methods and path-based methods. In embedding-based methods, users and items are represented by low-dimensional vectors, where entities' embeddings from the knowledge graph are used to enhance corresponding items' representations [3, 21, 22, 30]. Although these methods perform well, it's hard to explain the recommendation results because representations are in a latent space. In path-based methods, meta-paths and meta-graphs are commonly used to extract various semantic dependencies between users and items [29, 34]. However, it is almost computationally infeasible to enumerate all the useful meta-paths or meta-graphs. Moreover, the meta-paths and metagraphs need to be manually defined and cannot generalize to new datasets. Instead of pre-defining specific paths, the RippleNet [20] directly propagates users' preferences along edges in KG via the attention mechanism and then interprets the recommendations according to the attention scores, which however might not be trustworthy. The most recent work KPRN [24] uses LSTM to model the paths between users and items. However, sampling paths via breadth-first-search (BFS) is inefficient and may miss meaningful paths. Different from KPRN, our method defines the path finding problem as a sequential decision problem. We train an agent to automatically generate a meaningful path between a user and his/her relevant item via policy gradient methods.

**Relational Reasoning with Reinforcement Learning**. Our work is also related to recent work on knowledge graph reasoning with reinforcement learning [4, 13, 27], which aims to train an agent to walk on knowledge graphs to predict the missing facts. However, their goal is different from ours. We focus on the problem of recommendation with knowledge graphs and aim at finding meaningful paths for explaining the recommendation results while they focus on facts prediction.

# 3 Proposed Method: Ekar

In this section, we first formally define our problem and then introduce our proposed Explainable knowledge aware recommendation (Ekar) model in detail.

#### 3.1 Problem definition

We formally denote the interactions between users and items as a bipartite graph  $\mathcal{G}=(\mathcal{U},\mathcal{I})$ , where  $\mathcal{U}$  is the set of users, and  $\mathcal{I}$  is the set of items. Besides, we also have access to an open knowledge graph  $\mathcal{G}_k=(\mathcal{E}_k,\mathcal{R}_k)$ , where  $\mathcal{E}_k$  is the entity set and  $\mathcal{R}_k$  is the relation set. Each triplet  $\langle e_h,r,e_t\rangle\in\mathcal{G}_k$  indicates there exists a relation  $r\in\mathcal{R}_k$  from head entity  $e_h\in\mathcal{E}_k$  to tail entity  $e_t\in\mathcal{E}_k$ . For example,  $\langle$  *Titanic*, *DirectedBy, JamesCameron*  $\rangle$  reflects the fact that "Titanic" is directed by "James Cameron". As some items/entities are shared between  $\mathcal{I}$  and  $\mathcal{E}_k$ , we merge the user-item bipartite graph  $\mathcal{G}$  and the knowledge graph  $\mathcal{G}_k$  into an integrated user-item-entity graph  $\mathcal{G}'=(\mathcal{V}',\mathcal{R}')$ , where  $\mathcal{V}'=\mathcal{U}\cup\mathcal{I}\cup\mathcal{E}_k$ . For the user-item interaction graph  $\mathcal{G}$ , we assume all the edges belong to the relation "Interact", and therefore  $\mathcal{R}'=\{$  "Interact" $\}\cup\mathcal{R}_k$ .

Given a user u, our goal is to generate a path from u to a relevant item i on the user-item-entity graph  $\mathcal{G}'$ . Such a path not only allows to find the relevant item but also offers good explanations.

# 3.2 Formulating recommendation as a Markov Decision Process

There are some existing methods [24] trying to find meaningful paths between users and items. These methods first sample a collection of paths with breadth-first or depth-first search strategy and then measure the meaningfulness of the paths with a classification or ranking model. However, the number of possible paths between a user and an item could be exponentially large, and sampling a few of them could miss the meaningful ones. Moreover, the paths sampled through BFS or DFS strategy may not always be meaningful. In this paper, we take a different route and formulate the problem as a sequential decision making problem on the user-item-entity graph  $\mathcal{G}'$ . We aim to train an agent to walk on  $\mathcal{G}'$  to find relevant items. Starting from a target user u, the agent sequentially selects the next neighbor on the integrated user-item-entity graph  $\mathcal{G}'$  until it reaches the predefined maximum number of steps T. Formally, we define the states, actions, transition, and rewards of the Markov Decision Process as follows:

**States.** We represent the state as the sequence of traversed relations and entities so far, i.e.,  $s_t = (r_0, e_0, r_1, e_1, ..., r_t, e_t) \in \mathcal{S}_t$ , where  $r_t \in \mathcal{R}'$  and  $e_t \in \mathcal{V}'$  are relations and entities respectively. The initial state  $s_0 = (r_0, e_0)$  represents the target user, and  $r_0$  is an artificially introduced relation to be consistent with other  $(r_t, e_t)$  pairs.

**Actions**. When the agent is under state  $s_t$ , it can choose an outgoing edge of entity  $e_t$  as its next action. Formally, we define the possible actions under state  $s_t$  as  $\mathcal{A}_t = \{a = (r', e') | (e_t, r', e') \in \mathcal{G}'\}$ .

**Transition**. For the state transition  $\mathcal{P}(\mathcal{S}_{t+1} = s | \mathcal{S}_t = s_t, \mathcal{A}_t = a_t)$ , we adopt a deterministic strategy and simply extend the current state  $s_t$  by adding the new action  $a_t = (r_{t+1}, e_{t+1})$  as the next state, i.e.,  $s_{t+1} = (r_0, e_0, ..., r_t, e_t, r_{t+1}, e_{t+1})$ .

**Rewards**. No intermediate reward is provided for  $(s_{t < T}, a_{t < T})$ . The final reward depends on whether or not the agent correctly finds interacted items of the user u. Given the terminal entity  $e_T$ , the final reward  $R_T$  is +1 if user  $e_0$  has interacted with  $e_T$ , 0 if  $e_T$  is an item but user  $e_0$  has not interacted with it and -1 if  $e_T$  is not an item-type entity.

## 3.3 Solving recommendation MDP with policy gradient

We further parameterize the above MDP with deep neural networks and optimize it with policy gradient methods.

**Parameterizing MDP with Deep Neural Networks**. Since there are usually millions of entities and hundreds of relations in user-item-entity graph  $\mathcal{G}'$ , it is almost impossible to utilize discrete states and actions directly, the number of which is exponential to the number of symbolic atoms in  $s_t$  and  $a_t$  respectively. We, therefore, choose to represent entities and relations in  $\mathcal{G}'$  with low-dimensional embeddings. Each action a = (r, e) is represented as the concatenation of relation and entity embeddings, i.e.,  $\mathbf{a} = [\mathbf{r}'; \mathbf{e}']$ . The state  $s_t = (r_0, e_0, ..., r_t, e_t)$  is encoded by an LSTM [10]:

$$\mathbf{s}_0 = LSTM(\mathbf{0}, [\mathbf{r}_0; \mathbf{e}_0]),$$
  

$$\mathbf{s}_t = LSTM(\mathbf{s}_{t-1}, [\mathbf{r}_t; \mathbf{e}_t]), t > 0$$
(1)

where  $\mathbf{0}$  is a zero vector and  $\mathbf{s}_t$  is the low-dimensional representation of state  $s_t$ .

According to our initial definition of rewards, the agent gets a positive reward if and only if it successfully finds the target item. However, this might be problematic for a few reasons: First, for a large user-item-entity graph  $\mathcal{G}'$ , it is very difficult for the agent to reach the correct items due to the huge search space, especially at the beginning of training [27]. In other words, the rewards will be very sparse. As a result, the learning process of the agent could be very inefficient and take a long time to converge. Second, the goal of recommender systems is to infer new items that users are likely to interact with in the future, rather than repeating users' historical items. However, receiving positive rewards only from historical items discourages the agent to explore new paths and items, which should be the target of recommender systems. To accelerate the training process and meanwhile encourage the agent to explore items that have not been purchased or rated by the target user, we propose to shape the rewards [17] in the following way:

$$R_T = \begin{cases} 1, & \text{if } e_T \in \mathcal{I} \text{ and } < e_0, \bar{r}, e_T > \in \mathcal{G}', \\ \sigma(\psi(e_0, e_T)), & \text{if } e_T \in \mathcal{I} \text{ and } < e_0, \bar{r}, e_T > \notin \mathcal{G}', \\ -1, & \text{otherwise,} \end{cases}$$
 (2)

where  $\sigma(x) = \frac{1}{1+e^{(-x)}}$  is the sigmoid function and  $\bar{r}$  represents the relation "Interact".  $\psi(e_0, e_T)$  is the score function that measures the correlation between user  $e_0$  and the searched item  $e_T$ . In our study,  $\psi(e_0, e_T)$  is pre-trained by maximizing the likelihood of all triplets in graph  $\mathcal{G}'$  and can be the score function of any state-of-the-art knowledge graph embedding models [5, 28]. Different from the original rewards defined in Section 3.2, items that the target user has not interacted with now receive positive rewards, which are determined by the pre-trained knowledge graph embeddings.

**Policy Network**. In this paper, we use the policy gradient method to solve the proposed recommendation MDP. Based on parameterized state  $s_t$  and parameterized action a, we calculate the probability distribution over possible action  $A_t$  as follows:

$$\mathbf{y}_{t} = \mathbf{W}_{2} ReLU(\mathbf{W}_{1} \mathbf{s}_{t} + \mathbf{b}_{1}) + \mathbf{b}_{2},$$

$$\pi_{\theta}(a'|s_{t}) = \frac{\exp(\mathbf{a'}^{\mathsf{T}} \mathbf{y}_{t})}{\sum_{a \in \mathcal{A}_{t}} \exp(\mathbf{a}^{\mathsf{T}} \mathbf{y}_{t})},$$
(3)

where  $\{\mathbf{W}_1, \mathbf{W}_2\}$  and  $\{\mathbf{b}_1, \mathbf{b}_2\}$  are weight matrices and weight vectors of a two-layer fully-connected neural network, ReLU(x) = max(0, x) is the non-linear activation function and  $\pi_{\theta}(a'|s_t)$  is the probability of taken action a' under state  $s_t$ .

**Optimization**. During training, the agent starts with an initial state  $(r_0, e_0)$ , where  $e_0$  is the target user, and sequentially extends its path to a maximum length of T. We then use the reward function (i.e., Eq. 2) to assign the trajectory  $(s_0, a_0, s_1, a_1, ..., s_T)$  a final reward. Formally, we define the expected rewards over all traversed paths of all users as:

$$J(\theta) = \mathbb{E}_{e_0 \in \mathcal{U}}[\mathbb{E}_{a_1, a_2, \dots, a_T \sim \pi_{\theta}(a_t | s_t)}[R_T]]. \tag{4}$$

which is maximized via gradient ascent, and the gradients of all parameters  $\theta$  are derived by the REINFORCE [25] algorithm, i.e.,

$$\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \sum_{t} R_{T} \log \pi_{\theta}(a_{t}|s_{t}).$$
 (5)

## **4 Further Constraints on Actions**

The current action space under state  $s_t$  (i.e.,  $\mathcal{A}_t$ ) is defined as the set of outgoing edges of current entity  $e_t$ . This could be problematic for two reasons. First, for t < T, if entity  $e_t$  is already the correct item (i.e.,  $(e_0, e_t) \in \mathcal{G}$ ), the agent should stop and not continue to walk to other entities. Second, since the REINFORCE algorithm tends to encourage the agent to repeat historical experiences which receive high rewards [6], the algorithm may discourage the agent from exploring new paths and items, which could be relevant to the target user. We address the two problems in the following ways:

**Stop Action**. As the length of paths may vary for different user-item pairs, we should provide the agent an option to automatically terminate when it believes that it has found the right items ahead of T. Following Lin et al. [13], we add a special link from each node to itself. In this way, we allow the agent to stay at the ground truths, which can be understood as a stop action. We show the impact of using stop action in Section 5.4 by setting different path length T.

Action Dropout. To prevent the agent from repeating historical high-reward paths and encourage it to explore more possibilities, we propose to use action dropout [13] during the training. Specifically, instead of sampling an action from original  $\pi_{\theta}(a_t|s_t)$ , we use a mask upon  $\pi_{\theta}(a_t|s_t)$  to randomly drop some actions. In addition, action dropout can also help alleviate the problem of irrelevant paths between a user and an item since these paths may be found coincidentally at the beginning of training.

# 5 Experiments

In this section, we evaluate Ekar on three real-world datasets. Compared to other state-of-the-art methods, our proposed approach has the following advantages: 1) **Effectiveness**. Ekar significantly outperforms existing state-of-the-art KG-based recommendation methods in terms of recommendation accuracy. 2) **Explainability**. Case studies on generated paths demonstrate that Ekar can offer good explanations for recommended items.

## 5.1 Data and Experiment Settings

**Data.** We test Ekar on three benchmark datasets for KG-based recommendation: 1) *Last.FM*. This dataset contains a set of music artist listening information from a popular online music system Last.Fm. 2) *MovieLens-1M*. MovieLens-1M provides users' ratings towards thousands of movies. For these two datasets, we convert the explicit ratings into implicit feedback where each observed rating is treated as "1", and unobserved ratings are marked as "0"s. Following [22], we use Microsoft Satori to construct knowledge graphs for Last.FM and MovieLens-1M datasets respectively. 3) *DBbook2014*. This dataset provides users' reading history in the book domain. Its supporting knowledge graph is extracted from DBpedia. As we focus on KG-based recommendation, we remove items that have no matching entities in the corresponding knowledge graph. The statistics of processed datasets are presented in Table 1. Following [3, 20, 24], we randomly split the interactions of each user into training, validation, and test set with ratio 6:2:2.

Table 1: Statistics of evaluation datasets and corresponding knowledge graphs.

Data	User-Item Interaction				Knowledge Graph		
	# Users	# Items	# Events	Sparsity	# Entities	# Relations	# Triplets
Last.FM	1,872	3,846	21,173	99.71%	9,366	60	15,518
MovieLens-1M	6,040	2,347	656,462	95.37%	7,008	7	20,782
DBbook2014	5,576	2,598	65,445	99.55%	10,149	13	135,580

Baseline Methods. We compare Ekar with two kinds of methods: 1) Classical similarity-based methods including: ItemKNN, which recommends items that are most similar to target user's historical items; BPR-MF [29], which is a widely-used matrix factorization method using Bayesian Personalized Ranking (BPR) loss. 2) KG-based recommendation methods including: RippleNet [20], which propagates users' interests over knowledge graph with attention mechanism; CFKG [31], which learns users' and items' representations by applying TransE [2] on the graph  $\mathcal{G}'$ ; MKR [22], which learns both user-item matching task and knowledge graph embedding task under multi-task learning framework; KTUP [3], which is a state-of-the-art KG-based recommender that jointly learns translation-based recommendation [7] and translation-based knowledge graph embedding; ConvE-Rec, which learns users' and items' embeddings based on integrated graph  $\mathcal{G}'$  with ConvE [5]. As we use ConvE model for reward shaping, we treat ConvE-Rec as a special recommender.

**Evaluation Metrics**. Following [3, 8, 19], we adopt *Hit Ratio* (HR) and *Normalized Discounted Cumulative Gain* (NDCG) to evaluate the effectiveness of proposed Ekar and baseline methods. We use the same definition of HR and NDCG in [8], where HR measures whether test items are present in the recommendation list and NDCG assesses the ranking quality of test items respectively. In our study, we always report the averaged HR@K and NDCG@K scores across all users over five runs.

**Implementation Details.** First of all, we add  $\langle e_t, r^{-1}, e_h \rangle$  into  $\mathcal{G}'$  if a triplet  $\langle e_h, r, e_t \rangle$  exists to enhance the connectivity of graph, where  $r^{-1}$  is the inverse of relation r. Following [3], we only preserve those triplets that are directly connected to items in each supporting knowledge graph. We implement Ekar with Pytorch [18]. Entity and relation embeddings are pre-trained by applying ConvE <sup>1</sup> on graph  $\mathcal{G}'$ , and the embedding size is set to 32 for all methods except for ItemKNN, which has no latent representations. Meanwhile, we use the score function of ConvE to compute augmented rewards (i.e., Equation 2). From Figure 1, we can see path patterns "User->Item->Entity->item" and "User->Item->User->Item" are more probable to be meaningful, so we empirically set the maximum path length T to 3 as our default setting. We select action dropout rate from  $\{0.1-0.9\}$ , dropout rate for entity/relation embeddings from {0.1-0.9} using grid search. Meanwhile, grid search is also applied to select the optimal dropout rate for other baseline methods. For training, we use Adam [11] optimizer for all neural models with batch size of 512. During recommendation, we use beam search with beam size 64 to generate paths for target users. For duplicate paths leading to the same item, we keep the one with the highest probability. Finally, we adopt two different ranking strategies to generate final top-K recommendation list: (1) ranks the searched items according to the path probabilities and we denote it as Ekar, 2) ranks the searched items based on "rewards" defined by  $\sigma(\psi(e_0, e_T))$  in Equation 2 and we denote it as Ekar\*.

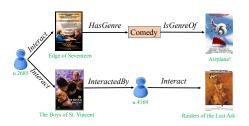
#### 5.2 Analysis of Recommendation Performance

We report the recommendation accuracy of different methods in Table 2. We can see that KG-based recommendation methods consistently outperform classical similarity-based methods, which indicates that knowledge graphs indeed help to alleviate the problem of data sparsity in the recommendation. Among KG-based recommendation methods, the RippleNet performs worst, which may be attributed to representing users with multi-hop away entities. KTUP performs strongly because it takes the advantages of both translation-based recommendation and multi-task learning. Note that the only difference between ConvE-Rec and CFKG is the used knowledge graph embedding methods; however, ConvE-Rec achieves much better performance. The reason behind this is that ConvE is a state-of-the-art knowledge graph embedding method, which outperforms TransE used in CFKG. By using pre-trained ConvE embeddings to augment rewards, our Ekar and Ekar\* significantly outperform

<sup>&</sup>lt;sup>1</sup>We also use DistMult [28] model to initialize entity/relation embeddings and to shape reward. It turns out that ConvE performs better. Please see the appendix for the comparison of using two models in detail.

Table 2: Recommendation	results of different models on three datasets.

Model	Last.FM		MovieLens-1M		DBbook2014	
Model	HR@10	NDCG@10	HR@10	NDCG@10	HR@10	NDCG@10
ItemKNN	0.0605	0.0511	0.0738	0.2273	0.0702	0.0665
BPR-MF	0.1199	0.0916	0.0895	0.1914	0.0829	0.0565
RippleNet	0.1008	0.0641	0.1269	0.2516	0.0763	0.0571
CFKG	0.1781	0.1226	0.1393	0.2512	0.1428	0.1036
MKR	0.1447	0.0850	0.1073	0.2245	0.0863	0.0575
KTUP	0.1891	0.1566	0.1579	0.3230	0.1761	0.1299
ConvE-Rec	0.2426	0.1742	0.1993	0.3676	0.1850	0.1357
Ekar	0.2201	0.1552	0.1889	0.3543	0.1716	0.1266
Ekar*	0.2483	0.1766	0.1994	0.3699	0.1874	0.1371
Gain over KTUP	31.31%	13.41%	26.28%	14.52%	6.41%	5.54%



Explainable preference paths	Pct. (%)
$U \xrightarrow{\mathit{Interact}} M \xrightarrow{\mathit{InteractedBy}} U \xrightarrow{\mathit{Interact}} M$	94.4
$U \xrightarrow{Interact} M \xrightarrow{HasGenre} G \xrightarrow{IsGenreOf} M$	1.5
$U \xrightarrow{Interact} B \xrightarrow{InteractedBy} U \xrightarrow{Interact} B$	29.6
$U \xrightarrow{Interact} B \xrightarrow{LinkedTo} WP \xrightarrow{Link} B$	27.4
$U \xrightarrow{Interact} B \xrightarrow{HasGenre} G \xrightarrow{IsGenreOf} B$	25.4

Figure 2: Recommendations and explanations for user #2685 in MovieLens-1M dataset. "Airplane!" is recommended because it shares the same movie genre (i.e., Comedy) with "Edge of Seventeen", which the target user watched before.

Table 3: Most frequent path patterns during recommendation on MovieLens-1M (top) and DB-book2014 (bottom) datasets. "U", "M", "G", "B" and "WP" represent User, Movie, Genre, Book and WikiPage respectively.

existing state-of-the-art KG-based recommendation methods in most cases and perform comparably to ConvE-Rec, which shows that our proposed methods are quite effective.

# 5.3 Analysis of Explainability

After demonstrating the effectiveness of Ekar, we now illustrate its explainability, which is the main contribution of this work to recommender systems. As introduced in previous sections, Ekar provides recommendations by generating meaningful paths from users to items, where paths serve as explanations for recommended items. To give you an intuitive example, we randomly select a real user from MovieLens-1M dataset and search preference paths for her/him with Ekar. As shown in Figure 2, we can easily understand that "Airplane!" is recommended because it shares the same genre (i.e., comedy) with "Edge of Seventeen", which the user watched before. The second recommendation "Raiders of the Lost Ark" can be explained with the well-known rule "Users who like A also like B".

Beyond explanations for an individual user, we are also interested in global preference path patterns discovered by Ekar. More specifically, we try to figure out what are the typical path patterns w.r.t. different datasets. From Table 3, we can see that Ekar heavily relies on the path pattern "User  $\xrightarrow{Interact}$  Movie  $\xrightarrow{Interact}$  Movie" on MovieLens-1M data. Additionally, Ekar learns more diverse path patterns such as "User  $\xrightarrow{Interact}$  Book  $\xrightarrow{LinkedTo}$  WikiPage  $\xrightarrow{Link}$  Book" and "User  $\xrightarrow{Interact}$  Book  $\xrightarrow{HasGenre}$  Genre  $\xrightarrow{IsGenreOf}$  book", which lead to new books that share the same WikiPage or Genre with users' historical books respectively. The reason behind the discrepancy of path patterns on two datasets may be two folds. First, note that the average number of interactions for each user in MovieLens-1M data is about 100 while this number is 12 in DBbook2014 data. Therefore it is easier to find a user sharing similar movie preference in MovieLens-1M data than to find a user with similar reading taste in DBbook2014 data. Second, the size of supporting knowledge graph for DBbook2014 data is much larger than the number of user-item interactions, so our Ekar learns to make more use of external knowledge in this case.

Table 4: Performance w.r.t. model variants, where [-] means removing that component from the Ekar.

Model	Last.FM		Movie	eLens-1M	DBbook2014	
	HR@10	NDCG@10	HR@10	NDCG@10	HR@10	NDCG@10
Ekar	0.2201	0.1552	0.1889	0.3543	0.1716	0.1266
Ekar-KG	0.2061	0.1466	0.1869	0.3489	0.0802	0.0525
Ekar-RS	0.0614	0.0349	0.0654	0.1132	0.1174	0.0867
Ekar-AD	0.1350	0.0827	0.1715	0.3217	0.1449	0.1083
Ekar (T=5)	0.2108	0.1505	0.1859	0.3500	0.1524	0.1125

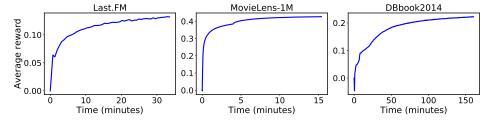


Figure 3: Results of running time with batch size of 512 and maximum path length of 3. The x-axis is the training time in minutes, and the y-axis is the average rewards over training samples.

## 5.4 Ablation Study

In this section, we compare different variants of Ekar to show the influences of some essential components such as KG, reward shaping, action dropout and maximum path length T, which are denoted as Ekar-KG, Ekar-RS, Ekar-AD and Ekar (T=5) respectively in Table 4. We find the influence of KG is very significant on Last.FM and DBbook2014 datasets. This is because these two datasets are extremely sparse, while the MovieLens-1M data is relatively dense. Removing reward shaping leads to a severe performance drop on all datasets because Ekar without reward shaping assigns zero rewards to all items that a user has not interacted with. In this way, the agent is penalized for exploring potential items that are of interest to users and therefore cannot effectively generate recommendations. Besides reward shaping, we also use action dropout to further encourage the agent to explore diverse paths, and we can see that Ekar-AD performs worse than the full model Ekar. At last, we try a larger maximum path length T to enable the agent to explore longer paths. We find that Ekar (T=5) performs worse than Ekar on all datasets. This is because long paths may introduce more noise and thus be less meaningful. However, thanks to the stop mechanism, the performance drop is not significant. The details about path patterns of Ekar (T=5), which further shows the role of stop action, are included in the appendix.

## 5.5 Convergence Analysis

We present the running time of Ekar in Figure 3. As can be seen, Ekar converges fast on MovieLens-1M dataset with less than ten minutes, while it takes longer to converge on the other two datasets. The reason for different convergence behaviors is that it is easier to walk to correct items on dense datasets (e.g., MovieLens-1M) than on sparse datasets (e.g., DBbook2014). Overall, our Ekar is efficient because we initialize entity/relation embeddings with pre-trained knowledge graph embeddings, and we use reward shaping to augment reward signals.

## 6 Conclusion

In this paper, we introduced a novel approach to provide explanations for recommendation with knowledge graphs. Our proposed Ekar generates meaningful paths from users to relevant items by learning a walk policy on the user-item-entity graph. Experimental results show that Ekar outperforms existing KG-based recommendation methods and is quite efficient. Furthermore, we demonstrate the explainability of Ekar via insightful case studies on different datasets. Future work includes incorporating domain knowledge to design proper reward functions for recommendation task and developing a distributed version of Ekar for even larger datasets.

#### Acknowledgments

The authors would like to thank Meng Qu and Zafarali Ahmed for providing useful feedback on initial versions of the manuscript. We also thank Yue Dong and Zhaocheng Zhu for editing the manuscript. WS and MZ are partially supported by Beijing Municipal Commission of Science and Technology under Grant No. Z181100008918005 as well as the National Natural Science Foundation of China (NSFC Grant Nos.61772039 and 91646202). WS also acknowledges the financial support by Chinese Scholarship Council. JT is supported by the Natural Sciences and Engineering Research Council of Canada, as well as the Canada CIFAR AI Chair Program.

## References

- [1] Konstantin Bauman, Bing Liu, and Alexander Tuzhilin. Aspect based recommendations: Recommending items with the most valuable aspects based on user reviews. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, pages 717–725, New York, NY, USA, 2017. ACM.
- [2] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in neural information processing systems*, pages 2787–2795, 2013.
- [3] Yixin Cao, Xiang Wang, Xiangnan He, Zikun Hu, and Tat-Seng Chua. Unifying knowledge graph learning and recommendation: Towards a better understanding of user preferences. In *The World Wide Web Conference*, WWW '19, pages 151–161, New York, NY, USA, 2019. ACM.
- [4] Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. In *International Conference on Learning Representations*, 2017.
- [5] Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [6] Kelvin Guu, Panupong Pasupat, Evan Liu, and Percy Liang. From language to programs: Bridging reinforcement learning and maximum marginal likelihood. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1051–1062, 2017.
- [7] Ruining He, Wang-Cheng Kang, and Julian McAuley. Translation-based recommendation. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*, RecSys '17, pages 161–169, New York, NY, USA, 2017. ACM.
- [8] Xiangnan He, Tao Chen, Min-Yen Kan, and Xiao Chen. Trirank: Review-aware explainable recommendation by modeling aspects. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, CIKM '15, pages 1661–1670, New York, NY, USA, 2015. ACM.
- [9] Jonathan L. Herlocker, Joseph A. Konstan, and John Riedl. Explaining collaborative filtering recommendations. In *Proceedings of the 2000 ACM Conference on Computer Supported Cooperative Work*, CSCW '00, pages 241–250, New York, NY, USA, 2000. ACM.
- [10] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [12] Jing Li, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Tao Lian, and Jun Ma. Neural attentive session-based recommendation. In *Proceedings of the 2017 ACM on Conference on Information* and Knowledge Management, CIKM '17, pages 1419–1428, New York, NY, USA, 2017. ACM.
- [13] Xi Victoria Lin, Richard Socher, and Caiming Xiong. Multi-hop knowledge graph reasoning with reward shaping. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3243–3253, 2018.

- [14] Yichao Lu, Ruihai Dong, and Barry Smyth. Why i like it: Multi-task learning for recommendation and explanation. In *Proceedings of the 12th ACM Conference on Recommender Systems*, RecSys '18, pages 4–12, New York, NY, USA, 2018. ACM.
- [15] Weizhi Ma, Min Zhang, Yue Cao, Woojeong Jin, Chenyang Wang, Yiqun Liu, Shaoping Ma, and Xiang Ren. Jointly learning explainable rules for recommendation with knowledge graph. In *The World Wide Web Conference*, WWW '19, pages 1210–1221, New York, NY, USA, 2019. ACM.
- [16] Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: Understanding rating dimensions with review text. In *Proceedings of the 7th ACM Conference on Recommender Systems*, RecSys '13, pages 165–172, New York, NY, USA, 2013. ACM.
- [17] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287, 1999.
- [18] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [19] Weiping Song, Zhiping Xiao, Yifan Wang, Laurent Charlin, Ming Zhang, and Jian Tang. Session-based social recommendation via dynamic graph attention networks. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, WSDM '19, pages 555–563, New York, NY, USA, 2019. ACM.
- [20] Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. Ripplenet: Propagating user preferences on the knowledge graph for recommender systems. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, CIKM '18, pages 417–426, New York, NY, USA, 2018. ACM.
- [21] Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. Dkn: Deep knowledge-aware network for news recommendation. In *Proceedings of the 2018 World Wide Web Conference*, WWW '18, pages 1835–1844, Republic and Canton of Geneva, Switzerland, 2018. International World Wide Web Conferences Steering Committee.
- [22] Hongwei Wang, Fuzheng Zhang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. Multi-task feature learning for knowledge graph enhanced recommendation. In *The World Wide Web Conference*, WWW '19, pages 2000–2010, New York, NY, USA, 2019. ACM.
- [23] Nan Wang, Hongning Wang, Yiling Jia, and Yue Yin. Explainable recommendation via multitask learning in opinionated text data. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '18, pages 165–174, New York, NY, USA, 2018. ACM.
- [24] Xiang Wang, Dingxian Wang, Canran Xu, Xiangnan He, Yixin Cao, and Tat-Seng Chua. Explainable reasoning over knowledge graphs for recommendation. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*, 2019.
- [25] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [26] Yao Wu and Martin Ester. Flame: A probabilistic model combining aspect based opinion mining and collaborative filtering. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, WSDM '15, pages 199–208, New York, NY, USA, 2015. ACM.
- [27] Wenhan Xiong, Thien Hoang, and William Yang Wang. Deeppath: A reinforcement learning method for knowledge graph reasoning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 564–573, 2017.
- [28] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. In *International Conference on Learning Representations*, 2015.
- [29] Xiao Yu, Xiang Ren, Yizhou Sun, Quanquan Gu, Bradley Sturt, Urvashi Khandelwal, Brandon Norick, and Jiawei Han. Personalized entity recommendation: A heterogeneous information network approach. In *Proceedings of the 7th ACM International Conference on Web Search and Data Mining*, WSDM '14, pages 283–292, New York, NY, USA, 2014. ACM.

- [30] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. Collaborative knowledge base embedding for recommender systems. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, pages 353–362, New York, NY, USA, 2016. ACM.
- [31] Yongfeng Zhang, Qingyao Ai, Xu Chen, and Pengfei Wang. Learning over knowledge-base embeddings for recommendation. In *Proceedings of the 41th International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '18, New York, NY, USA, 2018. ACM.
- [32] Yongfeng Zhang and Xu Chen. Explainable recommendation: A survey and new perspectives. *arXiv preprint arXiv:1804.11192*, 2018.
- [33] Yongfeng Zhang, Guokun Lai, Min Zhang, Yi Zhang, Yiqun Liu, and Shaoping Ma. Explicit factor models for explainable recommendation based on phrase-level sentiment analysis. In *Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '14, pages 83–92, New York, NY, USA, 2014. ACM.
- [34] Huan Zhao, Quanming Yao, Jianda Li, Yangqiu Song, and Dik Lun Lee. Meta-graph based recommendation fusion over heterogeneous information networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, pages 635–644, New York, NY, USA, 2017. ACM.
- [35] Kaiqi Zhao, Gao Cong, Quan Yuan, and Kenny Q Zhu. Sar: A sentiment-aspect-region model for user preference analysis in geo-tagged reviews. In 2015 IEEE 31st International Conference on Data Engineering, pages 675–686. IEEE, 2015.

# A Appendix

## A.1 Ekar with Different KGE Methods

Our proposed Ekar is independent of the knowledge graph embedding methods, which are utilized to pre-train entity/relation embeddings for initialization and reward shaping. Therefore we also test our model with another widely-used knowledge graph embedding method DistMult. The score functions of DistMult and ConvE are presented in Table 5. For a fair comparison, we use the same experimental settings and just substitute DistMult for ConvE in our experiments. Following ConvE-Rec, we denote recommendation with DistMult as DistMult-Rec.

The results of using different knowledge graph embedding methods are presented in Table 6. First, we can see that ConvE-Rec outperforms DistMult-Rec on all datasets. This is because ConvE has proven more effective than DistMult for knowledge graph embedding, as a result of which our Ekar with ConvE also outperforms Ekar with DistMult. Second, Ekar (DistMult) and Ekar\* (DistMult) outperform DistMult-Rec on Last.FM and MovieLens-1M datasets and perform comparably to DistMult-Rec on DBbook2014 dataset, which is consistent with the observations of Ekar (ConvE) and Ekar\* (ConvE).

Table 5: Score functions w.r.t. different KGE methods, where  $\langle \cdot \rangle$  denotes generalized inner product of three vectors,  $\bar{\cdot}$  denotes a 2D shaping of vectors, \* is the convolution operator,  $\omega$  denotes filters in convolutional layers,  $g(\cdot)$  is a non-linear activation function and vec( $\cdot$ ) converts a tensor to a vector.  $\mathbf{e}_0$ ,  $\mathbf{e}_T$  and  $\mathbf{r}$  are embeddings of user  $e_0$ , entity  $e_T$  and relation "Interact" respectively.

Model	Score function $\psi(e_0, e_T)$
DistMult ConvE	$g(\operatorname{vec}(g([\overline{\mathbf{e}_0}; \overline{\mathbf{r}}] * \omega))\mathbf{W})\mathbf{e}_T$

Table 6: Effectiveness comparison of using different knowledge graph methods for entity/relation initialization and reward shaping.

Model	Last.FM		MovieLens-1M		DBbook2014	
Model	HR@10	NDCG@10	HR@10	NDCG@10	HR@10	NDCG@10
ConvE-Rec	0.2426	0.1742	0.1993	0.3676	0.1850	0.1357
DistMult-Rec	0.1730	0.1174	0.1773	0.3341	0.1535	0.1090
Ekar (ConvE)	0.2201	0.1552	0.1889	0.3543	0.1716	0.1266
Ekar (DistMult)	0.1999	0.1397	0.1761	0.3352	0.1367	0.0958
Ekar* (ConvE)	0.2438	0.1766	0.1994	0.3699	0.1874	0.1371
Ekar* (DistMult)	0.1887	0.1265	0.1774	0.3343	0.1482	0.1061

## A.2 Further Analysis of Ekar (T=5)

In the previous experiments, we see that Ekar (T=5) performs a little worse than Ekarbecause long paths may not be meaningful and thus lead to bad recommendations. Now we present the path patterns of Ekar (T=5) on dense dataset MovieLens-1M and sparse dataset Last.FM.

We have two observations from Table 7. First, there are many paths of length five, which means that our agent is recommending items that have high-order similarity to users' historical items. Second, as the Last.FM dataset is very sparse, items with high-order similarity may be less relevant. Therefore our Ekar learns to take stop actions within 11.2% paths.

Table 7: Most frequent path patterns during recommendation on MovieLens-1M (top) and Last.FM (bottom) datasets. "U", "M", 'C', "G" and "A" represent User, Movie, Country, Genre and Artist respectively.

Explainable preference paths		
$U \xrightarrow{Interact} M \xrightarrow{InteractedBy} U \xrightarrow{Interact} M \xrightarrow{InteractedBy} U \xrightarrow{Interact} M$	60.6	
$U \xrightarrow{Interact} M \xrightarrow{Interacted By} U \xrightarrow{Interact} M \xrightarrow{Country} C \xrightarrow{CountryOf} M$	19.5	
$U \xrightarrow{Interact} M \xrightarrow{Interacted By} U \xrightarrow{Interact} M \xrightarrow{HasGenre} G \xrightarrow{IsGenreOf} M$	12.4	
$U \xrightarrow{Interact} A \xrightarrow{InteractedBy} U \xrightarrow{Interact} A \xrightarrow{InteractedBy} U \xrightarrow{Interact} A$	87.5	
$ U \xrightarrow{Interact} A \xrightarrow{InteractedBy} U \xrightarrow{Interact} A \xrightarrow{Self-loop} A \xrightarrow{Self-loop} A $	11.2	