

# Reinforced Path Reasoning for Counterfactual Explainable Recommendation

Xiangmeng Wang  
Data Science and Machine Intelligence Lab  
University of Technology Sydney  
Sydney, Australia  
xiangmeng.wang@student.uts.edu.au

Qian Li <sup>\*†</sup>  
School of Elec Eng, Comp and Math Sci  
Curtin University  
Perth, Australia  
qli@curtin.edu.au

Dianer Yu  
Data Science and Machine Intelligence Lab  
University of Technology Sydney  
Sydney, Australia  
Dianer.Yu-1@student.uts.edu.au

Guandong Xu<sup>†</sup>  
Data Science and Machine Intelligence Lab  
University of Technology Sydney  
Sydney, Australia  
guandong.xu@uts.edu.au

**Abstract**—Counterfactual explanations interpret the recommendation mechanism via exploring how minimal alterations on items or users affect the recommendation decisions. Existing counterfactual explainable approaches face huge search space and their explanations are either action-based (e.g., user click) or aspect-based (i.e., item description). We believe item attribute-based explanations are more intuitive and persuadable for users since they explain by fine-grained item demographic features (e.g., brand). Moreover, counterfactual explanation could enhance recommendations by filtering out negative items.

In this work, we propose a novel *Counterfactual Explainable Recommendation (CERec)* to generate item attribute-based counterfactual explanations meanwhile to boost recommendation performance. Our CERec optimizes an explanation policy upon uniformly searching candidate counterfactuals within a reinforcement learning environment. We reduce the huge search space with an adaptive path sampler by using rich context information of a given knowledge graph. We also deploy the explanation policy to a recommendation model to enhance the recommendation. Extensive explainability and recommendation evaluations demonstrate CERec’s ability to provide explanations consistent with user preferences and maintain improved recommendations. We release our code at <https://github.com/Chrystalii/CERec>.

**Index Terms**—Explainable Recommendation; Counterfactual Explanation; Counterfactual Reasoning; Reinforcement Learning

## I. INTRODUCTION

Modern recommendation models become sophisticated and opaque by modeling complex user/item context, e.g., social relations [1] and item profiles [2]. Hence the pressing need for faithful explanations to interpret user preferences meanwhile enable model transparency. Explainable Recommendation Systems (XRS) aim to provide personalized recommendations complemented with explanations that answer why particular

items are recommended [3]. It is fairly well-accepted that high-quality explanations can help improve users’ satisfactions and recommendation persuasiveness [4], [5]. Explanations can also facilitate system designers in tracking the decision-making of recommendation models for model debugging [6].

Existing XRS approaches can be categorized to model-intrinsic methods [7], [8] and model-agnostic methods [9]–[12]. Model-intrinsic methods seek recommendation models that are inherently interpretable, e.g., decision trees [7] and association rules [8], thus are limited in their applications to prominent deep learning models [3]. Recently, model-agnostic methods attract increasing attention [13] to allow the recommendation model to be black-box models, e.g., neural networks [10]. Model-agnostic approaches identify correlations between input user-item interactions and output recommendations from black-box models [13], to help users understand how their behaviors [10] (e.g., click) or which item features [11], [14] (e.g., brand) contribute to the recommendation. However, existing model-agnostic approaches suffer from major limitations: 1) They should firstly identify influential factors (e.g., item features) that cause positive user interactions and then construct explanations, while they seldom consider what recommendations would be if we directly intervene on explanations to alternative ones. 2) They do not consider explanation complexity by constructing explanations with a fixed number of influential factors [9], while totally ignoring to search for a minimal set of influential factors that well-explain reasons for recommendations.

Recently, counterfactual explanation [15] has emerged as a favorable opportunity to solve the above questions. Typically, *counterfactual explanation* is defined as a minimal set of influential factors that, if applied, flip the recommendation decision. Counterfactual explanation interprets the recommendation mechanism via exploring how minimal alterations on items or users affect the recommendation decisions. To build counterfactual explanations, we have to fundamentally

This work is supported by the Australian Research Council (ARC) under Grant No. DP200101374, LP170100891, DP220103717 and LE220100078.

\*: Contributing equally with the first author

†: Corresponding author

address: *what the recommendation result would be if a minimal set of factors (e.g., user behaviors/item features) had been different* [16]. With counterfactual explanations, users can understand how minimal changes affect recommendations and conduct counterfactual thinking under “what-if” scenario [17].

Existing counterfactual explanation based XRS can be summarized into two main categories. The first line of search-based approaches [18], [19] conducts greedy search on counterfactual instances given by perturbing user interactions or item features. Those counterfactual instances with the highest scores measured by heuristic quality metrics are considered counterfactual explanations. For example, Kaffes et al. [18] perturb counterfactual instances upon removal of items from users’ interactions, with normalized length and candidate impotence measuring counterfactuals to guide the search. Xiong et al. [19] perform constrained feature perturbations on item features and search perturbed features as counterfactual explanations. But search-based approaches suffer from high computational burden for the large-scale search space [16]. Another line of optimization-based approaches [15], [20]–[22] formulates counterfactual reasoning as an optimization problem to alleviate the computation burden. Typically, the optimization goal is to find simple but essential user actions or item aspects that directly cause user preference changes. Ghazimatin et al. [21] perform random walks over a user-item interaction graph and calculate PageRank scores after removing user actions edges from the graph. User actions that change PageRank scores are considered as counterfactual explanations. Tan et al. [15] modify item aspect scores to observe user preference changes based on a pre-defined user-aspect preference matrix. In summary, existing optimization-based approaches either focus on user action [20]–[22] or item aspect explanations [15], leaving the item attribute-based counterfactual explanations largely unexplored.

We claim that item attribute-based counterfactual explanations could benefit both users’ trust and recommendation performance. On the one hand, item attribute-based counterfactual explanations are usually more intuitive and persuadable to seek users’ trust. This is mainly because users prefer to know detailed information, e.g., which item attribute caused the film “Avatar” not be recommended anymore? Is it the director of “Avatar”? Besides, it is essential to search for a minimal change on item attributes that would alter the recommendation. We take Figure 1 as an example.  $u_1$  gave negative feedback on  $i_2$  with attributes  $\{p_1, p_2, p_3, p_4\}$ . We could infer attributes  $\{p_1, p_2, p_3, p_4\}$  reveal negative preferences of  $u_1$ . However, item  $i_1$  with attributes  $\{p_1, p_2\}$  was liked by  $u_1$ , which somehow reflect positive user preferences on  $\{p_1, p_2\}$ . Thus, merely using attributes  $\{p_1, p_2, p_3, p_4\}$  as explanations for  $u_1$ ’s recommendations would be controversial and misunderstand users’ preferences. In fact, the slightly changes between  $i_1$ ’s attributes and  $i_2$ ’s attributes, i.e.,  $\{p_3, p_4\}$ , could be the true determinant factors for explaining  $u_1$ ’s dislike. On the other hand, counterfactual explanations could boost recommendation performance since they offer high-quality negative signals of user preferences. Particularly, if we know the slightly

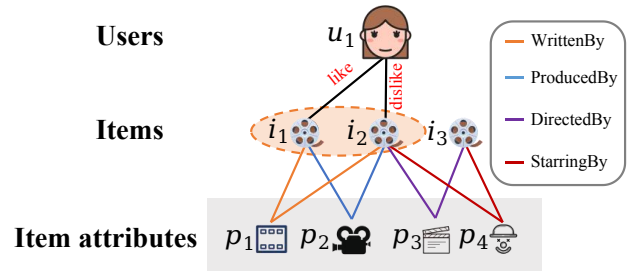


Fig. 1: Toy example of inferring item attribute-based counterfactual explanations from knowledge graphs.

changes  $\{p_3, p_4\}$  explain  $u_1$ ’s dislike, we could infer that item  $i_3$  with  $\{p_3, p_4\}$  would be disliked by  $u_1$  as well. To the end, the counterfactual explanation helps generate more precise recommendations by filtering out negative items.

In this work, we leverage knowledge graphs (KGs) that represent relations among real-world entities of users, items, and attributes to infer attribute-based counterfactual explanations meanwhile boosting recommendation performance. We propose a new *Counterfactual Explainable Recommendation (CERec)* that crafts the counterfactual optimization problem as a reinforcement learning (RL) task. The RL agent optimizes an explanation policy upon uniformly searching candidate counterfactuals. To reduce the search space, we propose an adaptive path sampler over a KG with a two-step attention mechanism and select counterfactual item attributes as the explanation candidates. With explanation candidates, the RL agent optimizes the explanation policy to find optimal counterfactual attribute explanations. We also deploy the explanation policy to a recommendation model to enhance the recommendation. The contributions of our work are:

- To the best of our knowledge, we are the first to leverage the rich attribute information in knowledge graphs to provide attribute-based counterfactual explanations for recommendations.
- We propose an RL-based framework to find the optimal counterfactual explanations, driven by an adaptive path sampler and a counterfactual reward function.
- We design an adaptive path sampler with a two-step attention mechanism to reduce the search space of counterfactual explanations.
- We train the recommendation model uniformly with the counterfactual explanations with the aim of boosting the recommendation.
- Extensive results on explainability and recommendation performance demonstrate the effectiveness of CERec.

## II. RELATED WORK

### A. Explainable Recommendation

Explainable recommendation is proven to improve user satisfaction [23] and system transparency [24]. Researches in explainable recommendation can be divided into two main categories [3], [25]. The first category of model-intrinsic

approaches [7], [8], [26]–[29] designs interpretable recommendation models to facilitate explanations. These model-intrinsic methods train explainable models to identify influential factors for users’ preferences, and then construct explanations accordingly. For example, Shulman et al. [7] propose a per-user decision tree model to predict users’ ratings on items. Each decision tree is built with a single user as a root and items as leaves, while a linear regression is applied to learn leaf node values to build decision rules. The explanation is formed as a sequence of decisions along the decision tree. Lin et al. [8] propose a personalized association rule mining method targeted at a specific user. Explanations are generated by identifying association rules between users and items. Other works [26]–[31] explore content-based approach [26] and neighborhood-based collaborative filtering [27], [28] to generate explanations based on users’ neighbors [27], [30], [31], item similarity [28], user/item features [26], or combinations thereof [29].

As modern recommendation models become complicated and black-box, e.g., embedding-based models [32]–[34] and deep neural networks [35], [36], developing model-intrinsic approaches [3] is not practical anymore. Thus, another category of model-agnostic methods [10], [14], [37] aims to explain black-box models in a post-hoc manner. For instance, Ghazimatin et al. [10] perform graph learning on a heterogeneous information network that unifies users’ social relations and historical interactions. The learnt graph embeddings are used to train an explainable learning-to-rank model to output explanation paths. Singh et al. [14] train a Learning-to-rank latent factor model to produce ranking labels. These ranking labels are then used for training an explainable tree-based model to generate interpretable item features for ranking lists.

### B. Counterfactual Explainable Recommendation

Counterfactual explanations have a long history in philosophy, psychology, and social sciences [16], [38], [39]. They have been considered as satisfactory explanations [40]–[42] and elicit causal reasoning in humans [43]–[45]. The counterfactual explanation has recently emerged as a hot topic in recommendation systems. Successful works can be categorized into search-based and optimization-based approaches. The first category of search-based approaches [18], [19] performs greedy search for counterfactual explanations. For instance, Kaffes et al. [18] perturb user-item interactions by deleting items from user interaction queues to generate perturbed user interactions as counterfactual explanations. Then a breadth-first search is used to find counterfactual explanations who achieve the highest normalized length and candidate impotence scores. Xiong et al. [19] propose a constrained feature perturbations on the features of candidate items and consider the perturbed item features as counterfactual explanations. The second category of optimization-based approaches [15], [20]–[22] optimize explanation models to find counterfactual explanations with minimal changes. For instance, Ghazimatin et al. [21] perform random walks over a heterogeneous information network and calculate the PageRank scores after removing user actions edges from the graph. Those minimal

sets of user actions that change PageRank scores are deemed counterfactual explanations. Tran et al. [22] identify a minimal set of user actions that updates parameters of neural models. Tan et al. [15] modify item aspect scores to observe user preference changes based on a user-aspect preference matrix.

Existing optimization-based approaches either focus on user action [20]–[22] or item aspect explanations [15]. In contrast, our counterfactual explanation is defined on the item attribute-side based on real-world item demographic features. Besides, we combine the search-based and optimization-based methods into an integrated mechanism by searching candidate counterfactuals to optimize an explanation policy within a reinforcement learning environment.

## III. PRELIMINARY

In this section, we first describe the collaborative knowledge graph that unifies user-item interactions and an external knowledge graph. Then, we introduce the concept of attribute-based counterfactual explanation for recommendation and give our task formulation.

### A. Collaborative Knowledge Graph

The collaborative knowledge graph (CKG) encodes user-item interactions and item knowledge as a unified relational graph. Let  $\mathcal{U} \in \mathbb{R}^M$ ,  $\mathcal{I} \in \mathbb{R}^N$  denote the sets of users and items, respectively. The historical user-item interaction matrix  $\mathcal{Y} = \{y_{ui} \mid u \in \mathcal{U}, i \in \mathcal{I}\}$  is defined according to users’ implicit feedback, where each entry  $y_{ui} = 1$  indicates there is an interaction between  $u$  and  $i$  and otherwise,  $y_{ui} = 0$ . In addition to user-item interactions, we also have an external knowledge graph (KG) that profiles item knowledge, which consists of real-world item and attribute entities and relations among them. Inspired by [46], we integrate rich item knowledge and the interaction data into a collaborative knowledge graph. Let  $\mathcal{E}$  and  $\mathcal{R}$  denote the entity and relation sets in the external KG. We first establish an item-entity alignment function  $\psi : \mathcal{I} \rightarrow \mathcal{E}$  to map items in the interaction data into KG entities, in which each item  $i \in \mathcal{I}$  is mapped to a corresponding item entity  $e \in \mathcal{E}$  in the KG. Then, CKG can be formulated as  $\mathcal{G} = \{(h, r, t) \mid h, t \in \mathcal{E}', r \in \mathcal{R}'\}$ , where  $\mathcal{E}' = \mathcal{E} \cup \mathcal{U}$ , and  $\mathcal{R}' = \mathcal{R} \cup \{\mathbb{I}(y_{ui})\}$ .  $\mathbb{I}(\cdot)$  is an edge indicator that denotes the observed edge between user  $u$  and item  $i$  when  $y_{ui} = 1$ , i.e.,  $\mathbb{I}(y_{ui}) = 1$  when  $y_{ui} = 1$ , otherwise  $\mathbb{I}(y_{ui}) = 0$ . Each triplet  $(h, r, t)$  describes that there is a link  $r$  from head entity  $h$  to tail entity  $t$ , which contains rich semantic relations among diverse types of real-world entities. For example,  $(\textit{Michael Jackson}, \textit{SingerOf}, \textit{Beat it})$  states the fact that *Michael Jackson* is the singer of the song *Beat it*.

### B. Counterfactual Explanation for Recommendation

The counterfactual explanation is defined as a minimal perturbation set of original sample that, if applied, would result in the opposite prediction of the sample [47]. In this work, we focus on finding the attribute-based counterfactual explanations, in which the minimal perturbation set is defined as a set of item attributes, e.g., *genre*, *brand*. Specifically,

suppose  $f_R$  is a Top- $K$  recommendation model that gives the recommendation list  $Q_u$  with length  $K$  for a user  $u$ , and we say  $i \in Q_u$  if item  $i$  is recommended. For a user-item pair  $(u, i)$ , we aim to search for a minimal item attributes set  $\Delta_{ui} = \{a_{ui}^1, \dots, a_{ui}^r\}$ . Each item attribute  $a_{ui}^i \in \Delta_{ui}$  is an attribute entity, which contains real-world semantic that describes the item. With the recommendation model  $f_R$ , and the attributes set  $\Delta_{ui}$ , our optimization goal is to estimate whether applying the  $\Delta_{ui}$  on the recommended item  $i$  results in replacing the  $i$  with a different item  $j$  for user  $u$ . If the optimization goal meet, the attributes set  $\Delta_{ui}$  is termed a *counterfactual explanation* that flips the recommendation result and  $j$  is the *counterfactual item* for the original item  $i$ . Meanwhile, the counterfactual explanation  $\Delta_{ui}$  is *minimal* such that there is no smaller set  $\Delta'_{ui} \in \mathcal{E}$  satisfying  $|\Delta'_{ui}| < |\Delta_{ui}|$  when  $\Delta'_{ui}$  also meets the optimization goal. With the optimized  $\Delta_{ui}$ , we can generate the counterfactual explanation for recommending item  $i$  to user  $u$ , which takes the following form:

Had a minimal set of attributes  $[\Delta_{ui}(s)]$  been different for item  $i$ , the recommend item would have been  $j$  instead.

### C. Task Formulation

Given the collaborative knowledge graph  $\mathcal{G}$ , we aim to construct a counterfactual explanation model that exploits rich relations among  $\mathcal{G}$  to produce counterfactual items for original recommend items. We formulate our counterfactual explanation model as:

$$j \sim \pi_E(u, i, f_R, \mathcal{G} | \Theta_E), \quad (1)$$

where  $\pi_E(\cdot)$  is the counterfactual explanation model parameterized with  $\Theta_E$ ,  $f_R$  is the Top- $K$  recommendation model that produces the recommendation results. The counterfactual explanation model generates empirical distribution over items among  $\mathcal{G}$  to yield counterfactual item  $j$  for the recommended item  $i$ , which is expected to meet the counterfactual goal with a minimal set of item attributes that reverses the recommendation result.

We also aim to use counterfactual items produced by  $\pi_E$  to improve the recommendation model  $f_R$ . Conceptually, a counterfactual item can offer high-quality negative signals to the recommendation, since it owns attributes that make the positive item not match the user's preference anymore. Inspired by pairwise ranking learning [48], we pair one positive user-item interaction with one counterfactual item to aggregate counterfactual signals into the recommendation model  $f_R$ . Formally,

$$\min_{\Theta_R} - \ln \sigma(f_R(u, i | \Theta_R) - f_R(u, j | \Theta_R)), \forall j \sim \pi_E(\Theta_E) \quad (2)$$

where  $f_R(\cdot)$  is the recommendation model parameterized with  $\Theta_R$  and  $\sigma(\cdot)$  is the sigmoid function.  $(u, i)$  is the positive interaction that satisfies  $\exists r \in \mathcal{R}', s.t. (u, r, i) \in \mathcal{G}$ . This formulation encourages positive items to receive higher prediction scores from the target user than counterfactual items. As such, the recommendation model not only produces recommendation

results for counterfactual explanation model training, but is also co-trained with the counterfactual explanation model by interactively aggregating counterfactual signals to enhance the recommendation.

## IV. MODEL FRAMEWORK

We now introduce our counterfactual explainable recommendation (*CERec*) framework that leverages rich relations among a collaborative knowledge graph (CKG) to generate attribute-based counterfactual explanations meanwhile providing improved recommendations. Our *CERec* consists of three modules - one recommendation model, one graph learning module and the proposed counterfactual explanation model. The recommendation model generates ranking scores and is co-trained with the counterfactual explanation model by interactively aggregating the produced counterfactual items. The graph learning module embeds users, items, and attributes entities among a given CKG as embedding vectors. The counterfactual explanation model conducts effective path sampling based on entity embeddings and ranking scores to discover high-quality counterfactual items. Two main parts are performed in our counterfactual explanation model: 1) counterfactual path sampler: uses entity embeddings to sample paths as actions for reinforcement learning agent; 2) reinforcement learning agent: learns explanation policy by optimizing the cumulative counterfactual rewards of deployed actions from the sampler. We depict the framework of *CERec* in Figure 2. We introduce the recommendation model and the graph learning module and give an overview of our counterfactual explanation model as below.

### A. Recommendation Model

We now present the recommendation model  $f_R$  that uses user and item latent factors for Top- $K$  recommendation. Here we employ the pairwise learning-to-rank model CliMF [49] as the recommendation model. The CliMF initializes IDs of users and items as latent factors, and updates latent factors by directly optimizing the Mean Reciprocal Rank (MRR) to predict ranking scores of items for users. It is worth noting that  $f_R$  can be any model as long as it takes users' and items' embeddings as part of the input and produce ranking results, which makes our counterfactual explanation framework applicable to a broad scope of models, e.g., neural networks [33]. Specifically, we firstly map users and items into latent factors with the recommendation model,

$$f_R(u, i) = U_u^\top V_i \quad (3)$$

where  $U_u$  denotes  $d$ -dimensional latent factors for user  $u$ , and  $V_i$  denotes  $d$ -dimensional latent factors for item  $i$ . We use the pairwise Mean Reciprocal Rank loss [49] to define our objective function to optimize parameters  $\Theta_R$  as below.

$$\mathcal{L}_R = \min_{\Theta_R} \sum_{(u, i) \in \mathcal{Y}} [\ln \sigma(f_R(u, i)) + \sum_{j=1} \ln(1 - \sigma(f_R(u, j) - f_R(u, i)))] \quad (4)$$

where  $\mathcal{Y} \in \mathbb{R}^{M \times N}$  denotes the historical user-item interaction matrix,  $\sigma(\cdot)$  is the sigmoid function and  $j \sim \pi_E(\Theta_E)$

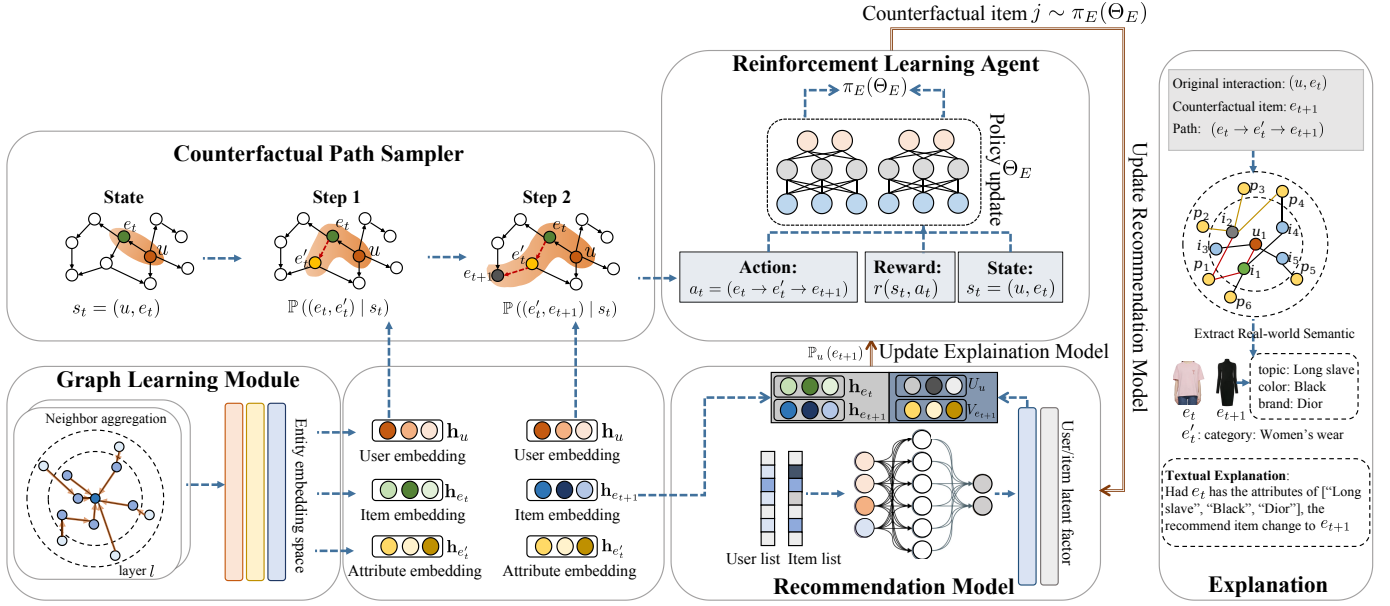


Fig. 2: Framework of *CERec*: *graph learning module* learns embeddings of users, items, and attributes entities from a CKG; *recommendation model* generates ranking scores of items for users; *counterfactual path sampler* uses entity embeddings to sample paths as actions for reinforcement learning agent; *reinforcement learning agent* learns the explanation policy by optimizing the cumulative counterfactual rewards of deployed actions from the sampler. The learnt explanation policy outputs counterfactual items for original user-item interactions. Meanwhile, paths connecting original items and counterfactual items from path sampler are saved to retrieve counterfactual item attributes. Finally, counterfactual explanations are generated by abstracting real-world semantics of counterfactual items and counterfactual item attributes.

is the counterfactual item generated from our counterfactual explanation model. By optimizing this loss function, we can get the ranking score for each item  $i$  from user  $u$  based on  $U$  and  $V$ . Formally,

$$\mathbb{P}_u(i) = \frac{\exp(U_u^\top V_i)}{\sum_{k=1}^K \exp(U_u^\top V_k)} \quad (5)$$

where  $K$  is the length of user  $u$ 's recommendation list  $\mathbf{Q}_u$ ,  $\mathbb{P}_u(i)$  indicates the ranking score of an item  $i$  in  $u$ 's recommendation list (i.e., with  $K$  items).

### B. Graph Learning Module

The graph learning module (GLM) learns users, items and attributes representations (i.e., embeddings) from the given collaborative knowledge graph  $\mathcal{G}$ . The learnt embeddings are deployed into our counterfactual explanation model to: (1) calculate the importance scores of user, item and attribute entities to form sampling paths as actions for counterfactual path sampler; (2) calculate the similarities among item entities to define the counterfactual rewards of deployed actions for reinforcement learning agent.

Inspired by recent advances in Graph Neural Networks (GNNs) [33], [34] for graph data representation, we employ GraphSage [32] in our GLM to learn representations for users, items and attributes entities. For an entity  $e \in \mathcal{G}$ , the GraphSage aggregates the information propagated from its neighbors  $\mathcal{N}_e$  to learn  $e$ 's representation. As a user entity would connect with entities whose type is the item, i.e.,

$\mathcal{N}_e \in \mathcal{I}$ , the learnt user embeddings capture the influence of historical user-item interactions. Analogously, item entities connect with item attribute entities such that the learnt item embeddings absorb context information from item attributes. In particular, we firstly initialize entity representations at the 0-th layer of GraphSage with Multi-OneHot [50] by mapping entity IDs into embeddings, where the embedding for an entity  $e$  is denoted by  $\mathbf{h}_e^0$ . Then, at the  $l$ -th graph convolutional layer, an entity  $e$  receives the information propagated from its neighbors to update its representation, as:

$$\mathbf{h}_e^{(l)} = \delta\left(\mathbf{W}^{(l)}\left(\mathbf{h}_e^{(l-1)} \parallel \mathbf{h}_{\mathcal{N}_e}^{(l-1)}\right)\right) \quad (6)$$

where  $\mathbf{h}_e^{(l)} \in \mathbb{R}^{d_l}$  is the embedding of an entity  $e$  at layer  $l$  and  $d_l$  is the embedding size;  $\mathbf{W}^{(l)} \in \mathbb{R}^{d_l \times 2d_{l-1}}$  is the weight matrix,  $\parallel$  is the concatenation operation and  $\delta(\cdot)$  is a nonlinear activation function as LeakyReLU [51].  $\mathbf{h}_{\mathcal{N}_e}^{(l-1)}$  is the information propagated from  $e$ 's neighbors set  $\mathcal{N}_e$ , as:

$$\mathbf{h}_{\mathcal{N}_e}^{(l-1)} = \sum_{e' \in \mathcal{N}_e} \frac{1}{\sqrt{|\mathcal{N}_e| |\mathcal{N}_{e'}|}} \mathbf{h}_{e'}^{(l-1)} \quad (7)$$

where  $\mathcal{N}_e = \{e' \mid (e, e') \in \mathcal{G}\}$  is a set of entities connected with  $e$ .

Having obtained the representations  $\mathbf{h}_e^{(l)}$  at each graph convolutional layer  $l \in \{1, \dots, L\}$ , we adopt layer-aggregation mechanism [52] to concatenate embeddings at all layers into a single vector, as follows:

$$\mathbf{h}_e = \mathbf{h}_e^{(1)} + \dots + \mathbf{h}_e^{(L)} \quad (8)$$

where  $\mathbf{h}_e^{(i)}$  is the embedding for an entity  $e$  at the  $i$ -th layer. By performing layer-aggregation, we can capture higher-order propagation of entity pairs across different graph convolutional layers. After stacking  $L$  layers, we obtain the final representation for each entity among the CKG. Note that in the following, we use  $\mathbf{h}_u$  to denote the embedding of user entity  $u$ , while  $\mathbf{h}_{e_t}$  is the item embedding for item entity  $e_t$  and  $\mathbf{h}_{e'_t}$  is the embedding for item attribute entity  $e'_t$ .

### C. Counterfactual Explanation Model

Our counterfactual explanation model contains two main parts: the counterfactual path sampler and the reinforcement learning agent. The counterfactual path sampler performs two-step attention on users, items, and attributes embeddings from GLM to search for a high-quality path as action  $a_t$  for each state  $s_t$ . Then action  $a_t$  and state  $s_t$  are fed into our reinforcement learning agent to learn the explanation policy. Based on ranking scores produced by recommendation model and item embeddings, the reinforcement learning agent learns the reward  $r(s_t, a_t)$  at state  $s_t$  and updates the explanation policy  $\pi_E(\Theta_E)$  accordingly. Finally, based on the learnt explanation policy  $\pi_E$  and path histories, our *CERec* generates attribute-based counterfactual explanations for recommendations. We detail our counterfactual explanation model in the next section.

## V. REINFORCED LEARNING FOR COUNTERFACTUAL EXPLANATION MODEL

In this section, we introduce our counterfactual explanation model assisted by GLM and recommendation model to generate explanation policy  $\pi_E$  over reinforcement depth  $T$ . Our counterfactual explanation model contains two main parts: *counterfactual path sampler* (CPS) performs attention mechanisms on entities among CKG to sample paths as actions for reinforcement learning agent; *reinforcement learning agent* learns the explanation policy  $\pi_E$  by optimizing cumulative counterfactual rewards of the sampled actions from CPS. We now introduce each part in our counterfactual explanation model as follows.

### A. Counterfactual Path Sampler

The counterfactual path sampler (CPS) conducts path exploration over CKG to sample paths as actions for the latter explanation policy learning. The basic idea is to, conditioned on the target user, start from the recommended item, learn to navigate to its item attribute, and then yield the potential counterfactual item along the sampling paths. In practice, to conduct such higher-order path sampling, large-scale CKGs are required since they encode rich relations (i.e., more than one-hop connectivity) among users, items, and item attributes. As a result, counterfactual items are sampled from higher-hop neighbors of target user-item interactions to form the candidate action space. However, learning counterfactual explanations from the whole candidate action space is infeasible since the space would cover potentially enormous paths. Thus, our CPS is designed to reduce the candidate action space by filtering out irrelevant paths and selecting important paths for

later explanation policy learning. Here, we employ attention mechanisms to calculate the importance of the visited entity condition on the source entity to generate sampling paths. We now introduce our CPS  $\mathbb{P}(a_t|s_t)$  that samples paths as actions  $a_t$  at each state  $s_t$  using attention mechanisms.

At each state  $s_t$ , the CPS produces a path of the given CKG as action. Formally,  $a_t \sim \mathbb{P}(a_t|s_t) = (e_t \rightarrow e'_t \rightarrow e_{t+1})$  is the path that roots at an item  $e_t$  towards another item proposal  $e_{t+1}$ , where  $(e_t, e'_t), (e'_t, e_{t+1}) \in \mathcal{G}$  and  $e_t, e_{t+1} \in I$  are connected via the item attribute  $e'_t$ . The action  $a_t = (e_t \rightarrow e'_t \rightarrow e_{t+1})$  is generated by the two-step path sampling: 1) choose an outgoing edge from  $e_t$  to the internal item attribute entity  $e'_t$ ; 2) determine the third entity  $e_{t+1}$  conditioned on  $e'_t$ . We separately model the confidence of each exploration step into two attention mechanisms, i.e.,  $\mathbb{P}((e_t, e'_t) | s_t)$  and  $\mathbb{P}((e'_t, e_{t+1}) | s_t)$ .

The first attention mechanism  $\mathbb{P}((e_t, e'_t) | s_t)$  specifies the importance of item attributes for  $e_t$ , which are sensitive to the current state  $s_t = (u, e_t)$ , i.e., user  $u$  and item  $e_t$ . Formally, for each outgoing edge from  $e_t$  to its item attribute  $e'_t \in \mathcal{N}_{e_t}$ , we first obtain the embeddings of item  $e_t$ , attribute  $e'_t$  and user  $u$  from Eq (8), which are denoted by  $\mathbf{h}_{e_t}$ ,  $\mathbf{h}_{e'_t}$  and  $\mathbf{h}_u$ . Then, the importance score of item attribute  $e'_t$  is calculated by:

$$\alpha^{(1)}(e_t, e'_t) = \mathbf{h}_u^\top \delta(\mathbf{h}_{e_t} \odot \mathbf{h}_{e'_t}) \quad (9)$$

where  $\alpha^{(1)}$  is the attention score at the first attention mechanism.  $\odot$  is the element-wise product and  $\delta(\cdot)$  is LeakyReLU [51]. Thereafter, we employ a softmax function to normalize the scores of all neighbors of  $e_t$  as:

$$\mathbb{P}((e_t, e'_t) | s_t) = \frac{\exp(\alpha^{(1)}(e_t, e'_t))}{\sum_{e''_t \in \mathcal{N}_{e_t}} \exp(\alpha^{(1)}(e_t, e''_t))} \quad (10)$$

where  $\mathcal{N}_{e_t}$  is the neighbor item attributes set for  $e_t$ .

Having selected item attribute  $e'_t$ , we employ another attention mechanism  $\mathbb{P}((e'_t, e_{t+1}) | s_t)$  to decide yield which item from its neighbors  $\mathcal{N}_{e'_t}$  as the counterfactual item proposal. We firstly calculate the attention score of  $e_{t+1} \in \mathcal{N}_{e'_t}$  based on attribute embedding of  $e'_t$ , item embedding of  $e_{t+1}$  and user embedding of  $u$ , as:

$$\alpha^{(2)}(e'_t, e_{t+1}) = \mathbf{h}_u^\top \delta(\mathbf{h}_{e'_t} \odot \mathbf{h}_{e_{t+1}}) \quad (11)$$

where  $\alpha^{(2)}$  is the attention score at the second attention mechanism.  $\mathbf{h}_u$ ,  $\mathbf{h}_{e'_t}$  and  $\mathbf{h}_{e_{t+1}}$  are embeddings of user  $u$ , attribute  $e'_t$  and item  $e_{t+1}$ . Then we normalize attention scores for all item neighbors in  $\mathcal{N}_{e'_t}$  to generate the selection probability of item  $e_{t+1}$ . Moreover, since we care for generating valid counterfactual items as proposals, we filter out irrelevant items that do not meet the counterfactual goal, i.e., those being recommended for  $u$  in the recommendation list  $\mathcal{Q}_u$ . Formally,

$$\mathbb{P}((e'_t, e_{t+1}) | s_t) = \begin{cases} \frac{\exp(\alpha^{(2)}(e'_t, e_{t+1}))}{\sum_{e''_{t+1} \in \mathcal{N}_{e'_t}} \exp(\alpha^{(2)}(e'_t, e''_{t+1}))}, & e_{t+1} \notin \mathcal{Q}_u \\ 0, & e_{t+1} \in \mathcal{Q}_u \end{cases} \quad (12)$$

Finally, the two attention mechanisms are aggregated into the CPS  $\mathbb{P}(a_t | s_t)$  to yield the path  $a_t = (e_t \rightarrow e'_t \rightarrow e_{t+1})$  as action for each state  $s_t$ :

$$\mathbb{P}(a_t | s_t) = \mathbb{P}((e_t, e'_t) | s_t) \cdot \mathbb{P}((e'_t, e_{t+1}) | s_t) \quad (13)$$

where  $\mathbb{P}((e_t, e'_t) | s_t)$  represents the probability of stepping from  $e_t$  to  $e'_t$  is derived from Eq. (10);  $\mathbb{P}((e'_t, e_{t+1}) | s_t)$  derived from Eq. (12) is the probability of selecting  $e_{t+1}$  as the counterfactual item proposal. With  $\mathbb{P}(a_t | s_t)$ , we can generate the action  $a_t$  for each state  $s_t$  for explanation policy learning.

## B. Reinforcement Learning Agent

We formulate the counterfactual explanation model as the path-based reinforcement learning (RL) agent to discover high-quality counterfactual items. Each action  $a_t$  is a path towards candidate counterfactual item. With the counterfactual reward  $r(s_t, a_t)$  measures whether the action  $a_t$  returns a valid counterfactual item for current state  $s_t$ .

1) *Agent*: Formally, the counterfactual explanation model is formulated as a Markov Decision Process (MDP)  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}\}$ , where  $s_t \in \mathcal{S}$  is the state absorbing the current user and the visited entity,  $a_t \in \mathcal{A}$  is the action deposited to the current state.  $\mathcal{P}$  is the transition of states, and  $\mathcal{R}$  is the reward function. In the policy learning, the explanation policy  $\pi_E(a_t | s_t)$  selects an action  $a_t \in \mathcal{A}$  to take conditioning on the current state  $s_t \in \mathcal{S}$ , and the counterfactual explanation model receives counterfactual reward  $r(s_t, a_t) \in \mathcal{R}$  for this particular state-action pair. The final explanation policy is learnt to maximize the expected cumulative counterfactual rewards. We introduce these key elements for RL as follows.

- $\mathcal{S}$ : is a continuous state space describing a target user and the current visited item entity among the CKG. Formally, for a user  $u$ , at step  $t$ , the user state  $s_t$  is defined as  $s_t = (u, e_t)$ , where  $u \in \mathcal{U}$  is a specific user and  $e_t \in \mathcal{I}$  is the item entity the agent visits currently. The initial state  $s_0$  is  $(u, i)$  and  $i$  is the positive item of  $u$ , i.e.,  $y_{ui} \in \mathcal{Y}$ .
- $\mathcal{A}$ : is a discrete space containing actions available for policy learning. The action  $a_t \in \mathcal{A}$  is a path sampled from the CPS  $\mathbb{P}(a_t | s_t)$ .
- $\mathcal{P}$ : is the state transition which absorbs transition probabilities of the current states to the next states. Given action  $a_t$  at state  $s_t$ , the transition to the next state  $s_{t+1}$  is determined as  $\mathbb{P}(s_{t+1} | s_t, a_t) \in \mathcal{P} = 1$ , where  $s_{t+1} = (u, e_{t+1})$ ,  $s_t = (u, e_t)$  and  $a_t = (e_t \rightarrow e'_t \rightarrow e_{t+1})$ .
- $\mathcal{R}$ :  $r(s_t, a_t) \in \mathcal{R}$  is the counterfactual reward measures whether the visited item  $e_{t+1}$  is a valid counterfactual item by deploying action  $a_t = (e_t \rightarrow e'_t \rightarrow e_{t+1})$  at state  $s_t$ , which is defined based on the two criteria: 1) *Rationality* [16]:  $e_{t+1}$  should receive a high confidence of being removed from the current user  $u$ 's recommendation list compared with the original item  $e_t$ ; 2) *Similarity* [38]: as a counterfactual explanation requires the minimal change of item attributes between counterfactual item and original item,  $e_{t+1}$  should be as similar as possible with

the original item  $e_t$ . The formal definition of the reward  $r(s_t, a_t)$  is given by:

$$r(s_t, a_t) = \begin{cases} 1 + \cos(\mathbf{h}_{e_t}, \mathbf{h}_{e_{t+1}}), & \text{if } \mathbb{P}_u(e_t) - \mathbb{P}_u(e_{t+1}) \geq \epsilon \\ \cos(\mathbf{h}_{e_t}, \mathbf{h}_{e_{t+1}}), & \text{otherwise} \end{cases} \quad (14)$$

where  $\epsilon$  is the recommendation threshold determines *Rationality*.  $\epsilon$  is defined as the margin between ranking scores of the original item  $e_t$  and the  $K$ -th item (i.e.,  $Q_u^K$ ) in  $u$ 's recommendation list  $Q_u$ , i.e.,  $\epsilon = \mathbb{P}_u(e_t) - \mathbb{P}_u(Q_u^K)$ .  $\cos(\mathbf{h}_{e_t}, \mathbf{h}_{e_{t+1}})$  is the cosine similarity between item embeddings  $\mathbf{h}_{e_t}$  and  $\mathbf{h}_{e_{t+1}}$  and is used to measure *Similarity*, i.e.,  $\cos(\mathbf{h}_{e_t}, \mathbf{h}_{e_{t+1}}) = \frac{\langle \mathbf{h}_{e_t}, \mathbf{h}_{e_{t+1}} \rangle}{\|\mathbf{h}_{e_t}\| \|\mathbf{h}_{e_{t+1}}\|}$ . Note that  $\mathbf{h}_{e_t}$  and  $\mathbf{h}_{e_{t+1}}$  are obtained by Eq.(8).

2) *Objective Function*: Using the trajectories  $\{\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}\}$  from the agent, our counterfactual explanation model seeks a counterfactual explanation policy  $\pi_E$  that maximizes the cumulative rewards  $R(\pi_E)$  over reinforcement depth  $T$ :

$$R(\pi_E) = \mathbb{E}_{s_0 \sim \mathcal{S}, a_t \sim \mathbb{P}(a_t | s_t), s_{t+1} \sim \mathbb{P}(s_{t+1} | s_t, a_t)} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right] \quad (15)$$

where  $T$  is the terminal step determines reinforcement depth,  $\gamma^t$  is a decay factor at current step  $t \in T$ .  $\pi_E$  is the explanation policy that produces counterfactual items for users' recommend items. The final counterfactual explanations are formed by distilling the attributes of counterfactual items and their real-world semantics.

## VI. MODEL OPTIMIZATION

We adopt the iteration optimization [53] to optimize the recommendation model and the counterfactual explanation model. The recommendation model is initialized to provide ranking scores for counterfactual explanation model training; then the counterfactual explanation model is optimized to output high-quality counterfactual items for positive user-item interactions. Thereafter, the recommendation model is updated by pairing one positive user-item interaction with one counterfactual item generated through the counterfactual explanation model. Next, we detail the recommendation model optimization and the explanation policy optimization.

### A. Recommendation Model Optimization

We first initialize recommendation model by pairing one positive interaction  $y_{ui} \in \mathcal{Y}$  with one unobserved item  $v \in \mathcal{I}$  sampled from uniform sampling [54]. Then the recommendation model is optimized by training together with the counterfactual explanation model. At each iteration, the counterfactual explanation model outputs a counterfactual item  $j \sim \pi_E(\Theta_E)$ , such counterfactual item  $j$  is then fed into recommendation model to update user and item latent factors  $U$  and  $V$  in Eq.(3). Finally, the recommendation model loss  $\mathcal{L}_R$  in Eq.(4) w.r.t. model parameters  $\Theta_R$  is optimized by using stochastic gradient descent (SGD) [55].

## B. Explanation Policy Optimization

Since our counterfactual explanation model involves discrete sampling steps, i.e., the counterfactual path sampler, which block gradients when performing differentiation. Conventional policy gradient methods such as stochastic gradient descent [55] fail to calculate such hybrid policy gradients. Thus, we solve the optimization problem through REINFORCE with baseline [56]. Having obtained the policy optimization function from Eq. (15), and the sampling function from Eq. (13), the optimization goal is to combine the two functions and optimize them together with a REINFORCE policy gradient w.r.t.  $\Theta_E$ :

$$\begin{aligned} \mathcal{L}_E &= \nabla_{\Theta_E} R(\pi_E) \\ &= \nabla_{\Theta_E} \mathbb{E}_{s_0 \sim \mathcal{S}, a_t \sim \mathbb{P}(a_t | s_t), s_{t+1} \sim \mathbb{P}(s_{t+1} | s_t, a_t)} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right] \\ &\simeq \frac{1}{T} \sum_{t=0}^T [\gamma^t r(s_t, a_t) \nabla_{\Theta_E} \log \mathbb{P}(a_t | s_t)] \end{aligned} \quad (16)$$

where  $\Theta_E$  are learnable parameters for our counterfactual explanation model.

## VII. EXPERIMENTS

We thoroughly evaluate the proposed CERec for counterfactual explainable recommendation on three publicly available datasets. We evaluate our CERec in terms of recommendation performance and explainability performance to particularly answer the following research questions:

- Whether CERec with a counterfactual explanation model could improve the recommendation performance compared with state-of-the-art recommendation models?
- Are the counterfactual explanations generated by CERec appropriate to explain users' preferences?
- How the reinforcement depth  $T$  in CERec impacts the recommendation performance?
- How does the training epoch impact the stability of our recommendation model?

### A. Experimental Setup

1) *Dataset*: We use three publicly available datasets: Last-FM, Amazon-book and Yelp2018. The statistics of these datasets are presented in Table I, which depicts their recorded historical user-item interactions and the underlying knowledge graph. The Amazon-book<sup>1</sup> dataset is a widely used book recommendation dataset, and the Last-FM<sup>2</sup> is a music listening dataset. For Amazon-book and Last-FM, we first map their recorded items into Freebase [57] entities. Then the knowledge graphs for Amazon-book and Last-FM are built by extracting knowledge-aware facts for each item from the Freebase. Yelp2018<sup>3</sup> is a business recommendation dataset. For Yelp2018, we extract its item

knowledge from the local business information network to construct the knowledge graph. Each dataset is processed by the following settings: for user-item interactions, we adopt a 10-core setting, i.e., retaining users and items with at least ten interactions. Each knowledge-aware fact (i.e.,  $\langle \text{user}, \text{item} \rangle$ ,  $\langle \text{item}, \text{attribute} \rangle$ ) is represented as a conceptual edge among the collaborative knowledge graph (CKG). Moreover, to ensure CKG quality, we filter out infrequent entities (i.e., lower than 10 in both datasets) and retain the relations appearing in at least 50 triplets.

TABLE I: Statistics of the datasets. Density is computed by  $\#Interactions / (\#Users \cdot \#Items)$ .

Dataset		Last-FM	Amazon-book	Yelp2018
User-Item Interaction	#Users	23,566	70,679	45,919
	#Items	48,123	24,915	45,538
	#Interactions	3,034,796	847,733	1,185,068
	#Density	0.268%	0.048%	0.057%
Knowledge Graph	#Entities	58,266	88,572	90,961
	#Relations	9	39	42
	#Triplets	464,567	2,557,746	1,853,704

2) *Evaluation*: To evaluate the recommendation performance, all models are evaluated with three representative Top- $K$  recommendation metrics: Recall@ $K$ , Normalized Discounted Cumulative Gain (NDCG)@ $K$  and Hit Ratio (HR)@ $K$ . The  $K$  is set as 20 by default. The average results w.r.t. the metrics over all users are reported in Table II while a Wilcoxon signed-rank test [58] is performed in Table II to evaluate the significance. We evaluate the quality of explanations in terms of *consistency*. The *consistency* measures to what extent the attributes in counterfactual explanations are appropriate to explain users' preferences. We adopt the Precision, Recall and  $F_1$  protocols follows CountER [15], which require ground-truth attribute sets for evaluations. As we aim to search for a minimal item attributes set that can flip the positive user-item interaction to the negative. The ground-truth set therefore absorbs item attributes that cause the user to dislike the item. Formally, the ground-truth set is defined as  $\mathcal{O}_{ui} = \{o_{ui}^1, \dots, o_{ui}^p\}$ , where  $o_{ui}^p = 1$  if user  $u$  has negative preference on the  $p$ -th attribute of item  $i$ ; otherwise,  $o_{ui}^p = 0$ . The implementation of negative item attributes extraction is detailed in Section VII-A4. On the other hand, our model produces the attributes set  $\Delta_{ui} = \{a_{ui}^1, \dots, a_{ui}^r\}$ , which constitutes the counterfactual explanation for user-item pair  $(u, i)$ . Then, for each user-item pair, the Precision, Recall and  $F_1$  of  $\Delta_{ui}$  with regard to  $\mathcal{O}_{ui}$  are calculated by:

$$\begin{aligned} \text{Precision} &= \frac{\sum_{p=1}^r o_{ui}^p \cdot I(a_{ui}^p)}{\sum_{p=1}^r I(a_{ui}^p)}, \quad \text{Recall} = \frac{\sum_{p=1}^r o_{ui}^p \cdot I(a_{ui}^p)}{\sum_{p=1}^r o_{ui}^p}, \\ F_1 &= 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \end{aligned} \quad (17)$$

where  $I(a_{ui}^p)$  is an identity function such that  $I(\cdot) = 1$  when  $a_{ui}^p \neq 0$  and otherwise,  $I(\cdot) = 0$ .

3) *Baselines*: To evaluate the recommendation and explainability performance, we compare our CERec with five popular recommendation models as well as three attribute-aware baselines. Each baseline is aligned with different evaluation

<sup>1</sup><http://jmcauley.ucsd.edu/data/amazon>

<sup>2</sup><https://grouplens.org/datasets/hetrec-2011/>

<sup>3</sup><https://www.yelp.com/dataset>



tasks either for (I) recommendation evaluation or for (II) explainability evaluation. The baseline models are listed in the following:

- **NeuMF** [36] (I): extends Matrix Factorization (MF) to Deep Neural Network (DNN) for modeling user-item interactions.
- **MCreC** [2] (I): leverages DNN to model meta path-based context from Heterogeneous Information Network (HIN). It propagates the context to user/item embeddings with the co-attention mechanism.
- **KGpolicy** [59] (I): uses a reinforcement learning agent to explore negative items over a knowledge graph (KG). It trains a Bayesian Personalized Ranking model with sampled negatives for the recommendation.
- **KGQR** [60] (I): models KG information with Graph Convolution Network(GCN). It learns target recommendation policy with KG-enhanced state representations through the deep Q-network.
- **KGAT** [46] (I): introduces collaborative knowledge graph (CKG) to learn node embeddings by aggregating CKG neighbors to enhance the recommendation.
- **NeuACF** [61] (II): uses a DNN to learn attribute-based latent factors of users and items through meta-path based similarity matrices. We calculate users’ preference scores on item attributes with the learnt attribute-based latent factors. Then we select the last-10 important attributes for each user-item pair to construct the explanation.
- **CaDSI** [30] (II): disentangles user embeddings as separated user intent chunks. It learns attribute-aware intent embeddings by assigning item context trained from a HIN to each user intent chunk. We use attribute-aware embeddings to predict users’ preference on item attributes and select the last-10 important attributes as explanations.
- **RDExp** (II): We randomly select 10 attributes from item attribute space for each user-item interaction and generate explanations based on the selected attributes.

4) *Implementation Details:* We train the recommendation model by the CliMF<sup>4</sup> with the train/test/validate sets, which are split from user-item interactions with a proportion of 60%/20%/20% of the original dataset. We optimize the CliMF using stochastic gradient descent (SGD) [55]. The same data splitting and gradient descent method are also applied in all baselines. The three datasets for training our counterfactual explanation model are preprocessed into collaborative knowledge graphs, and the REINFORCE [56] policy gradient is calculated to update the parameters. The REINFORCE is also applied to KGpolicy and KGQR, and same collaborative knowledge graphs are used in KGpolicy and KGAT. The embedding size for all baselines and our CERec is fixed as  $d = 64$ . Two graph convolutional layers with  $\{32, 64\}$  output dimensions are performed for the graph learning in our model. All neural networks-based (i.e., DNN, GCN) baselines also keep 2 layers. For MCreC, NeuACR and CaDSI, we use meta-paths with the path scheme of user-item-attribute-item

to ensure the model compatibility, e.g., *user-book-author-book* for Amazon-book dataset. For counterfactual explanation, we freeze the parameters of the trained counterfactual explanation model to produce one counterfactual item for each positive user-item interaction. The final explanation comprises item attributes that connect with the counterfactual item in the CKG but not with the positive user-item interaction. For explanation consistency evaluation, we construct ground-truth attributes sets using dynamic negative sampling (DNS) [62]. We first train an attribute-aware MF model by feeding positive user-item interactions and random negative item attributes. The DNS then uniformly draws one item attribute from the attribute space and feeds its latent factors into the MF model to predict the preference score. Then item attributes with the highest scores are selected as negative samples to train the MF model recursively. Among multiple rounds of sampling, the DNS generates negative item attributes that match users’ negative preferences.

The hyper-parameters of all models are chosen by the grid search, including learning rate,  $L_2$  norm regularization, discount factor  $\gamma$ , etc. The maximum epoch for all methods is set as 400, while an early stopping strategy is performed (i.e., if the loss stops to increase, then terminate the model training). We also search the reinforcement depth  $T$  in  $\{1, 2, 3, 4, 5\}$  (cf. Eq 15) for our model and report its effect in Section VII-D1.

### B. Counterfactual Enhanced Recommendation

In this section, we present the recommendation performance evaluation to answer the RQ1. At each iteration, the counterfactual explanation model in our CERec produces one counterfactual item for each positive user-item interaction to recursively train the recommendation model. Thus, we are interested to know whether incorporating counterfactual items could boost our recommendation performance. We present the recommendation results of our CERec and baselines in Table II and discuss the main findings below.

Our proposed CERec equipped with a counterfactual explanation model consistently outperforms all baselines across three datasets on both evaluation metrics. For example, CERec obtains 16.67%, 0.33% and 17.47% improvements for Recall@20, NDCG@20 and HR@20 respectively over the best baseline on Last-FM dataset. By designing a counterfactual explanation model, CERec is capable of exploring the high-order connectivity among the CKG to generate counterfactual items for recommendation model training. This verifies the effectiveness of our counterfactual explanation model in producing high-quality counterfactual items that can boost the recommendation performance. Counterfactual items provide high-quality negative signals of user preference. By pairing one counterfactual item with one positive user-item interaction for recommendation model training, our recommendation model learns to distinguish between positive items and negative items to generate more precise recommendations.

Among the knowledge-aware models, our CERec consistently outperforms MCreC, KGpolicy, KGQR, and KGAT. This is mostly because these knowledge-aware baselines might

<sup>4</sup><https://github.com/salvacarrion/orange3-recommendation>

TABLE II: Recommendation evaluation: bold numbers are the best results, best baselines are marked with underlines.

Top- $K$ Recommendation Evaluation										
Model	Dataset	Last-FM			Amazon-book			Yelp2018		
	$K$	Recall@ $K$	NDCG@ $K$	HR@ $K$	Recall@ $K$	NDCG@ $K$	HR@ $K$	Recall@ $K$	NDCG@ $K$	HR@ $K$
NeuMF	20	0.0651	0.0767	0.2921	0.0799	0.0611	0.1010	0.0279	0.0380	0.1008
	40	0.0807	0.0848	0.3409	0.1377	0.0769	0.2660	0.0357	0.0411	0.1241
	60	0.0913	<u>0.0955</u>	0.3788	0.1581	0.0888	0.3007	0.0389	0.0448	0.1552
	80	0.1005	<u>0.1009</u>	0.4267	0.1888	0.0974	0.3332	0.0411	0.0477	0.2008
MCRec	20	0.0770	0.0821	0.2811	0.0879	0.0666	0.1420	0.0327	0.0412	0.1234
	40	0.0825	0.0845	0.3124	0.1041	0.0712	0.2177	0.0338	0.0424	0.1406
	60	0.0919	0.0919	0.3888	0.1801	0.0844	0.2805	0.0429	0.0436	0.1721
	80	0.1112	0.0945	0.4282	0.2257	0.0957	0.3672	0.0457	0.0437	0.2205
KGpolicy	20	0.0761	0.0689	0.3197	<u>0.1242</u>	0.0684	<u>0.2211</u>	<u>0.0557</u>	0.0435	0.1528
	40	<u>0.1018</u>	0.0763	0.4091	<u>0.1716</u>	0.0805	<u>0.2947</u>	<u>0.0649</u>	0.0460	0.2588
	60	<u>0.1179</u>	0.0815	0.4639	<u>0.2048</u>	0.0879	<u>0.3417</u>	<u>0.0832</u>	<u>0.0547</u>	<u>0.3300</u>
	80	<u>0.1302</u>	0.0855	0.5006	0.2315	0.0935	0.3762	<u>0.0883</u>	<u>0.0618</u>	<u>0.3841</u>
KGQR	20	0.0821	0.0856	0.3162	0.1076	0.0787	0.2182	0.0417	0.0458	0.1621
	40	0.0932	<u>0.0923</u>	<u>0.4229</u>	0.1652	<u>0.0891</u>	0.2358	0.0456	0.0499	0.1899
	60	0.0999	0.0939	<u>0.5228</u>	0.1987	0.0912	0.3077	0.0512	0.0535	0.2172
	80	0.1132	0.0947	0.5323	0.2212	0.0942	0.3531	0.0557	0.0587	0.2566
KGAT	20	<u>0.0870</u>	<u>0.0897</u>	<u>0.3292</u>	0.0801	<u>0.0791</u>	0.2006	0.0444	<u>0.0472</u>	<u>0.2341</u>
	40	0.0962	0.0918	0.3762	0.1435	0.0812	0.2634	0.0469	<u>0.0511</u>	<u>0.2666</u>
	60	0.1083	0.0937	0.4515	0.1766	<u>0.0922</u>	0.3244	0.0501	0.0546	0.3015
	80	0.1232	0.0939	<u>0.5464</u>	<u>0.2378</u>	<u>0.0985</u>	<u>0.3921</u>	0.0544	0.0577	0.3655
CERec	20	<b>0.1015</b>	<b>0.0900</b>	<b>0.3867</b>	<b>0.1406</b>	<b>0.0803</b>	<b>0.2451</b>	<b>0.0650</b>	<b>0.0495</b>	<b>0.2515</b>
	40	<b>0.1304</b>	<b>0.0993</b>	<b>0.4801</b>	<b>0.1881</b>	<b>0.0926</b>	<b>0.3174</b>	<b>0.1025</b>	<b>0.0555</b>	<b>0.3549</b>
	60	<b>0.1478</b>	<b>0.1052</b>	<b>0.5295</b>	<b>0.2203</b>	<b>0.0999</b>	<b>0.3624</b>	<b>0.1317</b>	<b>0.0639</b>	<b>0.4218</b>
	80	<b>0.1603</b>	<b>0.1094</b>	<b>0.5628</b>	<b>0.2462</b>	<b>0.1054</b>	<b>0.3948</b>	<b>0.1572</b>	<b>0.0706</b>	<b>0.4739</b>
Improv.%	20	+ 16.67	+ 0.33	+ 17.47	+ 13.20	+1.52	+ 10.85	+ 16.70	+ 4.87	+ 7.43
	40	+ 28.09	+ 7.58	+ 13.53	+ 9.62	+ 3.93	+ 6.55	+ 57.94	+ 8.61	+ 33.12
	60	+ 25.36	+ 10.16	+ 1.28	+ 7.57	+ 8.35	+ 6.06	+ 58.29	+ 16.82	+ 27.82
	80	+ 23.11	+ 8.42	+ 3.00	+ 3.53	+ 7.01	+ 0.67	+ 78.03	+ 14.24	+ 23.38

not fully utilize the item knowledge by lacking the counterfactual reasoning ability as in our CERec. All knowledge-aware baselines employ various attention mechanisms to capture users' preferences on item knowledge w.r.t. different item attributes. However, they fail to consider the underlying mechanism that really triggered users' interactions. On the contrary, our CERec learns to discover the item attributes that cause the change of users' interactions, thus resulting in a promoted recommendation performance.

Our CERec achieves significant improvements over baselines on the Yelp2018 dataset, e.g., 78.03%, 14.24%, and 23.38% improvements for Recall@80, NDCG@80, and HR@80, respectively. The Yelp2018 dataset records a large number of inactive users that have few interactions with items. It is well acknowledged that inferring user preference for inactive users is challenging and limit the performance of recommendation systems [63]. However, our CERec can still achieve favorable recommendation performance on the Yelp2018. We attribute the improvements of our CERec to the augmenting of counterfactual items that could infer negative user preferences for inactive users and thus boost the recommendation performance.

### C. Counterfactual Explanation Quality

To answer RQ2, we evaluate the explainability of our framework by reporting the quantitative results of explanation evaluation metrics. Then we discuss the interpretability of our generated counterfactual explanations by distilling their real-world semantics.

1) *Quantitative Results:* We report the Precision, Recall and  $F_1$  (cf. (17)) of explanations generated by our CERec and baseline models in Table III to test the *consistency*. The three evaluation metrics reflect to what extent the item attributes in explanations are appropriate to explain users' preferences. By analyzing Table III, we have several observations:

Firstly, the RDExp method performs very poorly on both Precision, Recall, and  $F_1$ . The RDExp generates explanations by randomly sampling item attributes from the knowledge graph provided. The poor performance of RDExp shows that randomly choosing item attributes as explanations can barely reveal the reasons for recommendations. This demonstrates that high-quality explanations require the appropriate choice of item attributes. Secondly, the counterfactual explanations generated by our CERec consistently show superiority against all baselines in finding item attributes that are consistent with users' preferences. This indicates that our CERec achieves su-

TABLE III: Explainability evaluation w.r.t the consistency. All numbers in the table are percentage numbers with ‘%’ omitted. The number after  $\pm$  is the standard error.

		Explanation Consistency					
Model	Dataset	Amazon-book			Yelp2018		
		Precision%	Recall%	$F_1\%$	Precision%	Recall%	$F_1\%$
	RDEP	0.0196 $\pm$ 0.00	0.0090 $\pm$ 0.00	0.0117 $\pm$ 0.00	0.4595 $\pm$ 0.01	0.0987 $\pm$ 0.01	0.1320 $\pm$ 0.03
	NeuACF	4.7336 $\pm$ 1.21	0.6171 $\pm$ 0.06	1.0226 $\pm$ 0.08	8.5603 $\pm$ 1.51	8.2619 $\pm$ 1.47	7.0589 $\pm$ 1.38
	CaDSI	13.2632 $\pm$ 1.98	1.0897 $\pm$ 0.09	2.8721 $\pm$ 0.06	17.7322 $\pm$ 2.07	19.0801 $\pm$ 2.38	18.3096 $\pm$ 2.60
	CERec	<b>23.8969</b> $\pm$ 2.84	<b>3.1020</b> $\pm$ 0.92	<b>5.1411</b> $\pm$ 0.83	<b>45.5013</b> $\pm$ 3.11	<b>41.5931</b> $\pm$ 3.74	<b>37.2509</b> $\pm$ 3.09

perior explainability and can generate more relevant attribute-based explanations that truly match user preference. Unlike NeuACF and CaDSI that generate explanations by selecting a fixed number of item attributes indicated by their importance scores, our CERec searches a minimal set of item attributes that would flip the recommendation result. We hence conclude that inferring item attributes with minimal changes is more appropriate to generate explanations, since they reflect simple but essential attributes that directly cause user preference changes. This further sheds light on the importance of conducting counterfactual explanations for recommendations.

2) *Interpretability of Counterfactual Explanations:* We present a case study on the interpretability of counterfactual paths that root at a positive user-item interaction and end at its counterfactual item. Each freebase entities among paths are mapped into real-world entities. We show the real-world cases in Figure 3 and give our observations in the following.

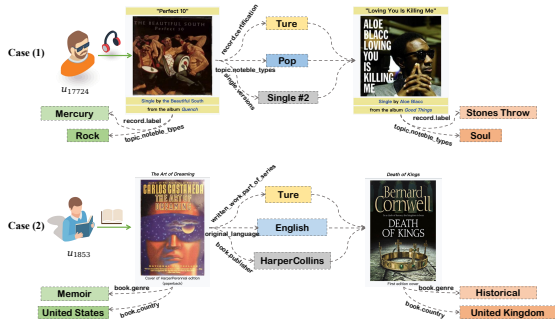


Fig. 3: Real cases of counterfactual reasoning paths.

The first example (Case 1) comes from the Last-FM dataset, where user  $u_{17724}$  positively interacted with “Perfect 10” in his listening record. Our CERec picks the “Loving you is killing me” as the counterfactual item, which shares three overlapping item attributes with the “Perfect 10”. Conventional explainable recommendation models would pick these three common item attributes to generate explanations as “you like Perfect 10 is because you like pop music that was certificated for its single#2 version”. However, this explanation cannot answer why the “Loving you is killing me” with the same item attributes is actually being removed from the user’s recommendation list. By contrast, our CERec captures the counterfactual item attributes (i.e., *Stones Throw* and *Soul*) and generates the counterfactual explanation as “if Perfect 10 has the label as *Stones Throw* and the type as *Soul*, then it will not

be recommended anymore.” This counterfactual explanation provides us with deeper insights of users’ real preferences. For example, analyzing the attribute differences, we find that although “Perfect 10” shares the same type *Pop* as “Loving you is killing me”, it is actually the *Rock* music published by *Mercury*. On the contrary, the “Loving you is killing me” that holds the opposite music type *Soul* and different polisher would potentially be disliked by the same user.

Analogously, in Case 2 from the Amazon-book dataset, our CERec picks the “Death of Kings” as the counterfactual item for the positive interaction between  $u_{1853}$  and “The Art of Dreaming”, and uses item attributes *Historical* and *United Kingdom* as the counterfactual explanation to reflect  $u_{1853}$ ’s preferences on book genre and country. These counterfactual explanations are more rational than conventional explanations, since they discard tedious associations of item attributes and expand explanations to counterfactual attributes.

#### D. Parameter Analysis

To answer RQ3, we first study how the reinforcement depth  $T$  in Eq. (15) affects the recommendation performance. We then investigate how the training epoch impacts the stability of our recommendation model.

1) *Impact of Reinforcement Depth:* The reinforcement depth  $T$  determines the search space, with  $T = 1$  denoting that at most 1-hop neighbors of starting entities are visited as proposals for the policy learning. We search the number of  $T$  in  $\{1, 2, 3, 4, 5\}$  and report the recommendation performance of our model in Figure 4. We have the following observations:

Firstly, our CERec with  $T = 4$  yields the best performance on the Amazon-book and Yelp2018 dataset, while  $T = 2$  gives the best results on the Last-FM. The Amazon-book and Yelp2018 are presented with 0.048% and 0.057% density and are much sparser compared with the Last-FM with 0.26% density. That means more inactive users with few item interactions are recorded in the Amazon-book and Yelp2018. To achieve the optimal recommendation performance, the Amazon-book and Yelp2018 require larger reinforcement depth (i.e.,  $T = 4$ ) compared with the Last-FM (i.e.,  $T = 2$ ). This is because diverse counterfactual items are retrieved by increasing the reinforcement depth to help the recommendation model achieve optimal performance on the two datasets. This finding indicates that augmenting diverse counterfactual items could provide precise recommendations for inactive users to enhance the recommendation.

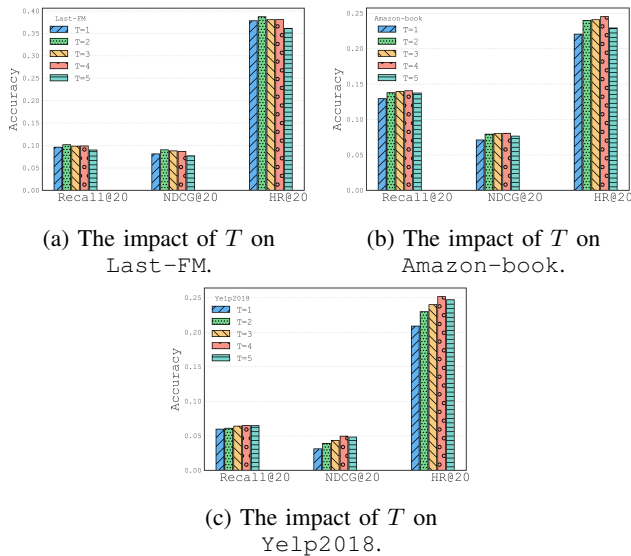


Fig. 4: Impact of reinforcement depth  $T$  on three datasets.

Counterfactual items offer high-quality negative signals for user preferences, and thus could help filter out negative items when providing recommendations for inactive users. Secondly, increasing the reinforcement depth enhances the recommendation performance before reaching the peaks on both datasets. We attribute such consistent improvements to the improved diversity of counterfactual items. This is because higher-hop item neighbors naturally cover more items beyond those unexposed but are actually the counterfactual ones than lower-hop neighbors. Thirdly, after peaks, increasing reinforcement depth leads to downgraded performance. This is because performing too many counterfactual item explorations introduces less-relevant items that may bias the recommendation results.

2) *Reward Gradient*: To evaluate the stability and the robustness of our model, we plot the cumulative rewards while training our model on the three datasets in Figure 5. Here are our observations.

First, our counterfactual explanation model gains stable cumulative rewards on both datasets along with training, i.e., the cumulative rewards increase as the epoch increases and finally reach stable states without suffering any drastic fluctuations. This indicates that our counterfactual explanation model enjoys stable training without losing reward gradients (e.g., gradients vanishing). This further verifies the rationality and robustness of our proposed model. Second, the rates of reward convergences are different across different datasets. For example, the cumulative rewards on Amazon-book and Yelp2018 start to increase drastically at the beginning and converge at around epoch 40, while the counterpart on Last-FM first converges slowly and then becomes stable at around epoch 120. This indicates that our model on Amazon-book and Yelp2018 can quickly reach stable states using a small number of iterations; even Amazon-book and Yelp2018 are much more sparse than Last-FM. This is because we incorporate counterfactual items while training

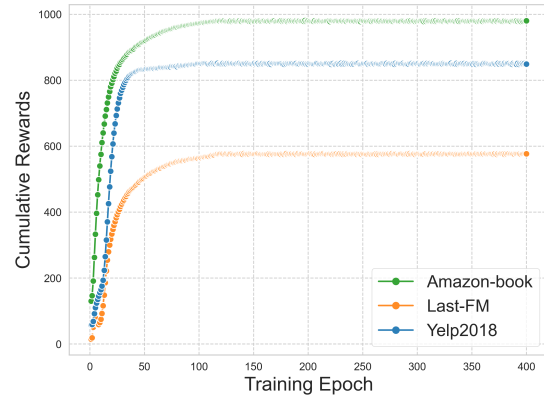


Fig. 5: Learning curves of cumulative rewards w.r.t. training epoch.

our model, which contain negative user preference signals to assist the better decision-making for inactive users. Besides, we explore higher-order connectivity among the CKG as side information for our model training, thus enhancing model robustness facing sparse datasets.

## VIII. CONCLUSION

In this work, we propose CERec, a reinforcement learning-based counterfactual explainable recommendation framework over a CKG. Our CERec is capable of generating attribute-based counterfactual explanations meanwhile provide precise recommendations. In particular, we design a counterfactual explanation model as a reinforcement learning agent to discover high-quality counterfactual items. The counterfactual explanation model takes paths sampled from our counterfactual path sampler as actions to optimize an explanation policy. By maximizing the counterfactual rewards of the deployed actions, the explanation policy is learnt to generate high-quality counterfactual items. In addition, we reduce the vast action space by utilizing attention mechanisms in our path sampler to yield effective paths from the CKG. Finally, the learnt explanation policy generates attribute-based counterfactual explanations for recommendations. We also deploy the explanation policy to a recommendation model to enhance the recommendation. Extensive explainability and recommendation evaluations on three large-scale datasets demonstrate CERec’s abilities to improve the recommendation and provide counterfactual explanations consistent with user preferences.

## ACKNOWLEDGMENT

This work is supported by the Australian Research Council (ARC) under Grant No. DP200101374, LP170100891, DP220103717 and LE220100078.

## REFERENCES

- [1] Q. Li, X. Wang, Z. Wang, and G. Xu, “Be causal: De-biasing social network confounding in recommendation,” *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 2022.

- [2] B. Hu, C. Shi, W. X. Zhao, and P. S. Yu, "Leveraging meta-path based context for top-n recommendation with a neural co-attention model," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 1531–1540.
- [3] Y. Zhang, X. Chen *et al.*, "Explainable recommendation: A survey and new perspectives," *Foundations and Trends® in Information Retrieval*, vol. 14, no. 1, pp. 1–101, 2020.
- [4] X. Wang, D. Wang, C. Xu, X. He, Y. Cao, and T.-S. Chua, "Explainable reasoning over knowledge graphs for recommendation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 5329–5336.
- [5] L. Xie, Z. Hu, X. Cai, W. Zhang, and J. Chen, "Explainable recommendation based on knowledge graph and multi-objective optimization," *Complex & Intelligent Systems*, vol. 7, no. 3, pp. 1241–1252, 2021.
- [6] N. Frosst and G. E. Hinton, "Distilling a neural network into a soft decision tree," in *Proceedings of the First International Workshop on Comprehensibility and Explanation in AI and ML 2017 co-located with 16th International Conference of the Italian Association for Artificial Intelligence (AI\*IA 2017), Bari, Italy, November 16th and 17th, 2017*, ser. CEUR Workshop Proceedings, T. R. Besold and O. Kutz, Eds., vol. 2071. CEUR-WS.org, 2017. [Online]. Available: [http://ceur-ws.org/Vol-2071/CExAIIA\\_2017\\_paper\\_3.pdf](http://ceur-ws.org/Vol-2071/CExAIIA_2017_paper_3.pdf)
- [7] E. Shulman and L. Wolf, "Meta decision trees for explainable recommendation systems," in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2020, pp. 365–371.
- [8] W. Lin, S. A. Alvarez, and C. Ruiz, "Collaborative recommendation via adaptive association rule mining," *Data Mining and Knowledge Discovery*, vol. 6, no. 1, pp. 83–105, 2000.
- [9] Y. Zhang, G. Lai, M. Zhang, Y. Zhang, Y. Liu, and S. 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*, 2014, pp. 83–92.
- [10] A. Ghazimatin, R. Saha Roy, and G. Weikum, "Fairy: A framework for understanding relationships between users' actions and their social feeds," in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, 2020, pp. 240–248.
- [11] X. Wang, Y. Chen, J. Yang, L. Wu, Z. Wu, and X. Xie, "A reinforcement learning framework for explainable recommendation," in *2018 IEEE international conference on data mining (ICDM)*. IEEE, 2018, pp. 587–596.
- [12] G. Peake and J. Wang, "Explanation mining: Post hoc interpretability of latent factor models for recommendation systems," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2060–2069.
- [13] Q. Ai and L. Narayanan, R. "Model-agnostic vs. model-intrinsic interpretability for explainable product search," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 5–15.
- [14] J. Singh and A. Anand, "Posthoc interpretability of learning to rank models using secondary training data," *CoRR*, vol. abs/1806.11330, 2018. [Online]. Available: <http://arxiv.org/abs/1806.11330>
- [15] J. Tan, S. Xu, Y. Ge, Y. Li, X. Chen, and Y. Zhang, "Counterfactual explainable recommendation," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 1784–1793.
- [16] S. Verma, J. Dickerson, and K. Hines, "Counterfactual explanations for machine learning: A review," 2020.
- [17] J. Pearl, "Causal inference in statistics: An overview," *Statistics surveys*, vol. 3, pp. 96–146, 2009.
- [18] V. Kaffes, D. Sacharidis, and G. Giannopoulos, "Model-agnostic counterfactual explanations of recommendations," in *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization*, 2021, pp. 280–285.
- [19] K. Xiong, W. Ye, X. Chen, Y. Zhang, W. X. Zhao, B. Hu, Z. Zhang, and J. Zhou, "Counterfactual review-based recommendation," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 2231–2240.
- [20] W. Cheng, Y. Shen, L. Huang, and Y. Zhu, "Incorporating interpretability into latent factor models via fast influence analysis," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 885–893.
- [21] A. Ghazimatin, O. Balalau, R. Saha Roy, and G. Weikum, "Prince: Provider-side interpretability with counterfactual explanations in recommender systems," in *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 196–204.
- [22] K. H. Tran, A. Ghazimatin, and R. Saha Roy, "Counterfactual explanations for neural recommenders," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 1627–1631.
- [23] N. Tintarev and J. Masthoff, "Evaluating the effectiveness of explanations for recommender systems," *User Modeling and User-Adapted Interaction*, vol. 22, no. 4, pp. 399–439, 2012.
- [24] K. Balog and F. Radlinski, "Measuring recommendation explanation quality: The conflicting goals of explanations," in *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, 2020, pp. 329–338.
- [25] F. Gedikli, D. Jannach, and M. Ge, "How should i explain? a comparison of different explanation types for recommender systems," *International Journal of Human-Computer Studies*, vol. 72, no. 4, pp. 367–382, 2014.
- [26] J. Vig, S. Sen, and J. Riedl, "Tagplanations: explaining recommendations using tags," in *Proceedings of the 14th international conference on Intelligent user interfaces*, 2009, pp. 47–56.
- [27] J. L. Herlocker, J. A. Konstan, and J. Riedl, "Explaining collaborative filtering recommendations," in *Proceedings of the 2000 ACM conference on Computer supported cooperative work*, 2000, pp. 241–250.
- [28] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Item-based collaborative filtering recommendation algorithms," in *Proceedings of the 10th international conference on World Wide Web*, 2001, pp. 285–295.
- [29] A. Papadimitriou, P. Symeonidis, and Y. Manolopoulos, "A generalized taxonomy of explanations styles for traditional and social recommender systems," *Data Mining and Knowledge Discovery*, vol. 24, no. 3, pp. 555–583, 2012.
- [30] X. Wang, Q. Li, D. Yu, P. Cui, Z. Wang, and G. Xu, "Causal disentanglement for semantics-aware intent learning in recommendation," *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- [31] Q. Li, X. Wang, Z. Wang, and G. Xu, "Be causal: De-biasing social network confounding in recommendation," *ACM Trans. Knowl. Discov. Data*, apr 2022. [Online]. Available: <https://doi.org/10.1145/3533725>
- [32] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *Advances in neural information processing systems*, vol. 30, 2017.
- [33] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, "Heterogeneous graph attention network," in *The World Wide Web Conference*, 2019, pp. 2022–2032.
- [34] Z. Hu, Y. Dong, K. Wang, and Y. Sun, "Heterogeneous graph transformer," in *Proceedings of The Web Conference 2020*, 2020, pp. 2704–2710.
- [35] C. Chen, M. Zhang, Y. Liu, and S. Ma, "Neural attentional rating regression with review-level explanations," in *Proceedings of the 2018 World Wide Web Conference*, 2018, pp. 1583–1592.
- [36] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua, "Neural collaborative filtering," in *Proceedings of the 26th international conference on world wide web*, 2017, pp. 173–182.
- [37] C. Nóbrega and L. Marinho, "Towards explaining recommendations through local surrogate models," in *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing*, 2019, pp. 1671–1678.
- [38] C. Fernandez, F. J. Provost, and X. Han, "Explaining data-driven decisions made by AI systems: The counterfactual approach," *CoRR*, vol. abs/2001.07417, 2020. [Online]. Available: <https://arxiv.org/abs/2001.07417>
- [39] J. Dodge, Q. V. Liao, Y. Zhang, R. K. E. Bellamy, and C. Dugan, "Explaining models: An empirical study of how explanations impact fairness judgment," *CoRR*, vol. abs/1901.07694, 2019. [Online]. Available: <http://arxiv.org/abs/1901.07694>
- [40] J. Woodward, *Making Things Happen: A Theory of Causal Explanation*, ser. Oxford Studies in the Philosophy of Science. New York: Oxford University Press, 2004. [Online]. Available: <https://oxford.universitypressscholarship.com/10.1093/0195155270.001.0001/acprof-9780195155273>
- [41] T. D. Duong, Q. Li, and G. Xu, "Stochastic intervention for causal effect estimation," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.
- [42] Q. Li, Z. Wang, S. Liu, G. Li, and G. Xu, "Deep treatment-adaptive network for causal inference," *The VLDB Journal*, pp. 1–16, 2022.
- [43] R. M. Byrne, *The rational imagination: How people create alternatives to reality*. MIT press, 2007.

- [44] Q. Li, Z. Wang, S. Liu, G. Li, and G. Xu, "Causal optimal transport for treatment effect estimation," *IEEE transactions on neural networks and learning systems*, 2021.
- [45] Q. Li, T. D. Duong, Z. Wang, S. Liu, D. Wang, and G. Xu, "Causal-aware generative imputation for automated underwriting," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 3916–3924.
- [46] X. Wang, X. He, Y. Cao, M. Liu, and T.-S. Chua, "Kgat: Knowledge graph attention network for recommendation," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 950–958.
- [47] R. K. Mothilal, A. Sharma, and C. Tan, "Explaining machine learning classifiers through diverse counterfactual explanations," in *Proceedings of the 2020 conference on fairness, accountability, and transparency*, 2020, pp. 607–617.
- [48] H. Narasimhan, A. Cotter, M. Gupta, and S. Wang, "Pairwise fairness for ranking and regression," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5248–5255.
- [49] Y. Shi, A. Karatzoglou, L. Baltrunas, M. Larson, N. Oliver, and A. Hanjalic, "Climf: Learning to maximize reciprocal rank with collaborative less-is-more filtering," in *Proceedings of the Sixth ACM Conference on Recommender Systems*, ser. RecSys '12. New York, NY, USA: Association for Computing Machinery, 2012, p. 139–146. [Online]. Available: <https://doi.org/10.1145/2365952.2365981>
- [50] S. Zhang, Z. Han, Y.-K. Lai, M. Zwicker, and H. Zhang, "Stylistic scene enhancement gan: mixed stylistic enhancement generation for 3d indoor scenes," *The Visual Computer*, vol. 35, no. 6, pp. 1157–1169, 2019.
- [51] J. Xu, Z. Li, B. Du, M. Zhang, and J. Liu, "Reluplex made more practical: Leaky relu," in *2020 IEEE Symposium on Computers and Communications (ISCC)*. IEEE, 2020, pp. 1–7.
- [52] K. Xu, C. Li, Y. Tian, T. Sonobe, K.-i. Kawarabayashi, and S. Jegelka, "Representation learning on graphs with jumping knowledge networks," in *International Conference on Machine Learning*. PMLR, 2018, pp. 5453–5462.
- [53] B. Jiang, T. Lin, S. Ma, and S. Zhang, "Structured nonconvex and nonsmooth optimization: algorithms and iteration complexity analysis," *Computational Optimization and Applications*, vol. 72, no. 1, pp. 115–157, 2019.
- [54] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme, "Bpr: Bayesian personalized ranking from implicit feedback," in *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, 2009, pp. 452–461.
- [55] L. Bottou, "Stochastic gradient descent tricks," in *Neural networks: Tricks of the trade*. Springer, 2012, pp. 421–436.
- [56] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," *Advances in neural information processing systems*, vol. 12, 1999.
- [57] K. Bollacker, R. Cook, and P. Tufts, "Freebase: A shared database of structured general human knowledge," in *AAAI*, vol. 7, 2007, pp. 1962–1963.
- [58] R. F. Woolson, *Wilcoxon Signed-Rank Test*. John Wiley & Sons, Ltd, 2008, pp. 1–3. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/9780471462422.eoct979>
- [59] X. Wang, Y. Xu, X. He, Y. Cao, M. Wang, and T.-S. Chua, "Reinforced negative sampling over knowledge graph for recommendation," in *Proceedings of The Web Conference 2020*, 2020, pp. 99–109.
- [60] S. Zhou, X. Dai, H. Chen, W. Zhang, K. Ren, R. Tang, X. He, and Y. Yu, "Interactive recommender system via knowledge graph-enhanced reinforcement learning," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 179–188.
- [61] X. Han, C. Shi, S. Wang, S. Y. Philip, and L. Song, "Aspect-level deep collaborative filtering via heterogeneous information networks," in *IJCAI*, 2018, pp. 3393–3399.
- [62] W. Zhang, T. Chen, J. Wang, and Y. Yu, "Optimizing top-n collaborative filtering via dynamic negative item sampling," in *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, 2013, pp. 785–788.
- [63] X. Wang, Q. Li, D. Yu, and G. Xu, "Off-policy learning over heterogeneous information for recommendation," ser. WWW '22. New York, NY, USA: Association for Computing Machinery, 2022, p. 2348–2359. [Online]. Available: <https://doi.org/10.1145/3485447.3512072>