# Knowledge graph enhanced recommender system

Zepeng Huai[1,2], Jianhua Tao[1,2,3], Feihu Che[1,2], Guohua Yang[1], Dawei Zhang[1]

[1]National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences

[2]School of Artificial Intelligence, University of Chinese Academy of Sciences

[3]CAS Center for Excellence in Brain Science and Intelligence Technology

Beijing, China

huaizepeng2020@ia.ac.cn,jhtao@nlpr.ia.ac.cn,chefeihu2017@ia.ac.cn,dawei.zhang@nlpr.ia.ac.cn,guohua.yang@nlpr.ia.ac.cn

## ABSTRACT

(Since KDD requires the papers that have been submitted to arXiv at least one month prior to the deadline need a different title and abstract, here we give a brief version of title and abstract.) Knowledge Graphs (KGs) have shown great success in recommendation. This is attributed to the rich attribute information contained in KG to improve item and user representations as side information. However, existing knowledge-aware methods leverage attribute information at a coarse-grained level both in item and user side. In this paper, we proposed a novel *attentive knowledge graph attribute network*(AKGAN) to learn item attributes and user interests via attribute information in KG. Technically, AKGAN adopts a heterogeneous graph neural network framework, which has a different design between the first layer and the latter layer. With one attribute placed in the corresponding range of element-wise positions, AKGAN employs a novel interest-aware attention network, which releases the limitation that the sum of attention weight is 1, to model the complexity and personality of user interests towards attributes. Experimental results on three benchmark datasets show the effectiveness and explainability of AKGAN.

## CCS CONCEPTS

• **Information systems** → **Recommender systems**.

## KEYWORDS

Recommendation, Knowledge Graph, Graph Neural Network

## 1 INTRODUCTION

Recommender systems have shown great success in e-commerce, online advertisement, and social medial platforms. The core task of recommender systems is to solve the overload information problem [31, 32, 40] and suggest items that users are potentially interested in. Traditional methods to achieve information filtering are content-based and collaborative filtering (CF)-based recommender systems[9, 20, 34], which utilize items' content features and the

similarity of users or items from interaction data respectively[7]. However, both of them don't introduce much side information.

In recent years, introducing knowledge graphs (KGs) into recommender systems as side information has been effective for improving recommendation performance. KGs represent real-world entities and illustrate the relationship between them with graph data structure, which can reveal multiple attributes of items and explore the potential reason for user-item interactions. Generally, a KG contains item nodes and attribute nodes, and attribute nodes can not only describe items' attributes directly but also represent other attribute nodes' attributes. For example, a movie knowledge graph (as shown in figure 1) has six types of nodes, where movie node (i.e., $e_1$) is item node and the others (i.e., $e_{2-10}$) are attribute nodes. $e_{2-4}$ represent the actor attribute of movie $e_1$, which means actors $e_{2-4}$ stared in movie $e_1$. And $e_{8,9}$ represent the singer attribute of song $e_5$, which means singer $e_{8,9}$ sang the song $e_5$. To integrate attribute information into recommender system, earlier works focus on embedding-based methods [1, 3, 12, 45, 46] and path-based methods [19, 26, 41–43]. The former exploits KGs with knowledge graph embedding (KGE) algorithm (i.e., TransE [2] and TransH [38]) to learn entity embeddings and then feed them into a recommender framework. The latter usually predefines a path scheme (i.e., meta-path) and leverages path-level semantic similarities of entities to refine the representations of users and items. However, both of them don't capture high-order connectivities and fail to exploit both the rich semantics and topology of KGs. More recently, the propagation-based methods, a.k.a. graph neural network (GNN)-based methods such as KGAT[33], KGNN-LS[30], KNI[23], AKGE[25], KGIN[35], have attracted considerable interest of researchers. With the attribute information iteratively propagating in KGs, GNN-based methods integrate multi-hop neighbors into representations and have achieved the state-of-the-art recommendation results.

Despite the success of GNN in exploiting multi-hop attribute information, we argue that there are still three shortcomings: (1)**The pollution caused by weighted sum operation in merging different attribute information**. Different attributes are independent in terms of semantics and user preference. Take the actor attribute $e_2$ and the singer attribute $e_8$ of movie $e_1$ in figure 1 as an example, semantics independence means $e_2$ doesn't co-occur with $e_8$ in each movie, since an actor and a singer are invited to work for a movie respectively. Preference independence means whether a user prefers actor $e_2$ is independent with whether he likes singer $e_8$. However, existing GNN-based methods pool embeddings from different attribute nodes with weighted sum (i.e., sum, mean, attention) operation, which brings about semantic pollution
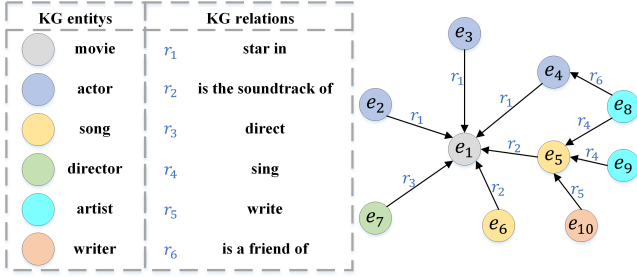
**Figure 1: An example of how KG contains multiple attributes information. Best viewed in color.**

and the difficulty to distill user interested attribute information. (2) **The nonlinearity between the distance and importance of attribute node relative to item node**. Usually, high-order neighbors are less relative to the center node in graph data, but this is not absolute. For example, singer attribute $e_8$ is more valued than director attribute $e_7$ by some movie viewers who like music, while $e_8$, a 2-hop neighbor, is further than $e_7$, a 1-hop neighbor. This problem is caused by inherent graph topology, which means some significant attribute nodes don't connect to item nodes directly, such that it requires multiple passes to integrate these high-order neighbors to the center node. However, existing GNN-based algorithms neglect this issue and decrease the weight of significant high-order neighbors coupled with the increase of propagation times. (3)**The complexity and personality of user interests towards different attributes**. User interests show the following pattern: a user just gets interested in a part of rather than all attributes of an item, and different users prefer different attributes even towards the same item. For example, both user $u_1$ and $u_2$ watch movie $e_1$ because $u_1$ likes $e_1$'s actor $e_2$ rather than director $e_7$ while $u_2$ prefers the theme song $e_5$ rather than actor $e_2$. Therefore, we should integrate actor $e_2$ rather than director $e_7$ attribute embedding into $u_1$'s representation and integrate song $e_5$ rather than actor $e_2$ attribute embedding into $u_2$'s representation. In other words, the item representation learned by GNN propagation contains noisy signals and we should distill part attributes interested by user personally. However, existing GNN-based algorithms recognize this pattern insufficiently, like [33] doesn't consider noisy attributes and [35] doesn't consider personal extraction.

To address the foregoing problems, we propose a novel *attentive knowledge graph attribute network* (AKGAN), which consists of two components: (1) **knowledge graph attribute network (KGAN)**. The core task of KGAN is to learn informative item representations without semantic pollution and weight decrease of significant attribute nodes. Technologically, KGAN has a different design between the first layer and the latter layer under GNN framework. The first layer, called attribute modeling layer, aims to generate initial item representations without semantic pollution. It regards each relation in KG as an attribute and has two key designs: embedding each entity in different attribute spaces and using concatenation operation to merge different attribute information. The latter layer, called attribute propagation layer, aims to remain semantics unpolluted and avoids weight decrease of significant neighbors

after GNN propagation. Finally, KGAN generates item representations where one attribute is represented within a specific range of element-wise positions independently. (2) **user interest-aware attention network**. With different attributes placed in corresponding element-wise positions, we design an interest-aware attention layer to distill user interested attributes. Specifically, for a particular user, we pool the representations of his interacted items and extract each attribute representation according to corresponding element-wise positions. Then each attribute representation will be fed into an attention module to calculate a personal interest score, which describes how much he prefers this attribute. We introduce a novel activation unit to release the limitation that the sum of attention weight is 1, which aims to reserve the intensity of user interests[50]. Finally all interest scores and item representations are further combined to infer user representations.

To summarize, the main contributions of this work are as follows:

- On the item side, we propose a novel knowledge graph attribute network, which employs a heterogeneous design in different layers under GNN framework, embeds each entity in different attribute spaces, and combines different attributes via concatenation operation, to avoid pollution caused by weighted sum operation and weight decrease of significant neighbors.
- On the user side, we develop an interest-aware attention network, which introduces a novel activation unit and releases the limitation that the sum of attention weight is 1 , to model user interests towards attributes personally.
- We conduct extensive experiments on three public benchmarks, and the results demonstrate the superior performance of AKGAN over state-of-the-art baselines. In-depth analyses are provided to illustrate the interpretability of AKGAN for user personal interests.

## 2 PROBLEM FORMALIZAITON

We begin by introducing some related notations and then defining the KG enhanced recommendation problem. Let $\mathcal{U} = \{u\}$ and $\mathcal{I} = \{i\}$ separately denote the user and item sets. A typical recommender system usually has historical user-item interactions, which is defined as $O^+ = \{(u, i)|u \in \mathcal{U}, i \in \mathcal{I}\}$. Each $(u, i)$ pair indicates user $u$ has interacted with item $i$ before, such as clicking, review, or purchasing. We also have a knowledge graph which stores the structured semantic information of real-world facts. We denote the entity and relation sets in KG as $\mathcal{R} = \{r\}$ and $\mathcal{E} = \{e\}$. KG is presented as $\mathcal{G} = \{(h, r, t)|h \in \mathcal{E}, r \in \mathcal{R}, t \in \mathcal{E}\}\}$, where each tripile means there is a relation $r$ from head entity $h$ to tail entity $t$. For example, $(James\ Cameron, Direct, Avader)$ describes the fact that James Cameron is the director of the movie Avatar. Note that $\mathcal{R}$ contains relations in both canonical direction (e.g., Direct) and inverse direction (e.g., DirectedBy). The bridge of KG and recommender system is that items are contained in entities. Specificlly, $\mathcal{E}$ consists of item node $\mathcal{I}(\mathcal{I} \subseteq \mathcal{E})$ as well as attribute node $\mathcal{E} \setminus \mathcal{I}$, which can futher refine user representation $e_u$ and item representation $e_i$.

We now formulate the KG enhanced recommendation task to be addressed in this paper:

- **Input**: a knowledge graph $\mathcal{G}$, which contains rich sturcture senmeantic information of item, and the user-item interaction data $O^+$

- **Output**: a scoring function that denotes the probability that user u will interact with item $i$.

## 3 METHODOLOGY

### 3.1 Model Overview

In this subsection, we introduce the framework of AKGAN. With a widely used two-tower structure in recommender model[8, 21, 47], AKGAN consists of two main modules: (1) knowledge graph attribute network, which is illustrated in Figure 2, and (2) user interest-aware attention module, which is illustrated in Figure 3.

The core task of KGAN is to learn KG enhanced item representations that contain multiple attribute information. KGAN adopts GNN framework and has a heterogeneous design in the first layer and the latter layer, which are named attribute modeling layer and attribute propagation layer respectively.

**Attribute modeling layer** (AML) is to construct initial entity representations by merging one-hop neiborhoods' attribute information in knowledge graph, as follows:

$$e_i^{(0)} = f_{\text{AML}}(e_i^{atr}, \mathcal{G}) \tag{1}$$

**Attribute propagation layer** (APL) recursively propagates attribute information to acquire more informative item representations as

$$e_i^{(l)} = f_{\text{APL}}(e_i^{(l-1)}, \mathcal{G}) \tag{2}$$

After performing $L$ layers, we obtain multiple item representations, namely $\{e_i^{(0)}, ..., e_i^{(L)}\}$. As the output of $l_{\text{th}}$ layer represents $l-$hop attribute information, we conduct sum opration to pool them and infer final item representations as

$$e_i^* = \sum_{l=0}^{L} e_i^{(l)} \tag{3}$$

When item representations have been learned, **interest-aware attention** (IAA) module is to learn user representations via interaction data with a novel relation-aware attention mechanism which represents user interests towards relations, a.k.a. attributes. Formally, user representations are obtained by IAA as

$$e_u^* = f_{\text{IAA}}(e_i, O^+) \tag{4}$$

When user representations and item representations are learned, a scoring function is used to predict their matching score and here we adopt inner product operation as

$$\hat{y}(u, i) = e_u^{*\top} e_i^* \tag{5}$$

### 3.2 Attribute Modeling Layer

In KGs, one entity has different types of relations with neighbors. For example, in figure 2, $e_8$ has different relations with $e_4$ and $e_5$. To model a pure initial representation without semantic pollution, we regard each relation as an attribute and embeds each entity in all attribute embedding spaces, which is similar to the design of FFM[14] that embeds each feature in different fields. Therefore each entity has sevecal embeddings and we denote embedding set in different attribute spaces as

$$\mathcal{E}^{atr} = \{(e_i^{r_1}, e_i^{r_2}, ..., e_i^{r_M}) | i \in \mathcal{E}, M = |\mathcal{R}|\} \tag{6}$$

where $e_i^{r_m} \in \mathbb{R}^{d^m}$ denotes the embedding of entity $i$ in attribute $r_m$ space and $d^m$ is the embedding dimension of attribute $r_m$ space.

Then we construct initial entity representations by aggregating attribute embeddings of one-hop neighbors. We can see that each neighbor has several embeddings and just one relation-aware embedding will be used. Taking figure 2 as an example, there are two links $(e_8, r_4, e_5)$ and $(e_8, r_6, e_4)$, we use embeddings $e_8^{r_4}$ to represents the singer attribute of song $e_5$ and embedding $e_8^{r_6}$ to represents the friend attribute of actor $e_4$, respectively. After we prepared the relation-aware embeddings, we adopt average operation to pool the same relation-aware neighbors to acquire the main semantics of this attribute as

$$e_{i'}^{r_j} = \frac{1}{|\mathcal{N}_i^{r_j}|} \sum_{j \in \mathcal{N}_i^{r_j}} e_j^{r_j} \tag{7}$$

where $j \in \mathcal{N}_i^{r_j}$ and $\mathcal{N}_i^{r_j} = \{j | j \in (j, r_j, i), (j, r_j, i) \in \mathcal{G}\}$ denotes the set of head entities that belong to the triplet where $i$ is tail entity and $r_j$ is relation. For example, if $e_{2,3}$ are both comedy actors and $e_4$ is an action actor, such that $mean(e_{2-4}^{r_1})$ represents movie $e_1$ is likely to show much funny performance. Note that not all attributes occur in one-hop range, like movie $e_1$ doesn't has friend attribute $r_6$. To obtain a fixed-length representation, we provide zero vector to the absent attribute. Finally, all attribute embeddings will be integrated into one representation by concatenation operation as

$$e_i^{(0)} = e_{i'}^{r_1} \| e_{i'}^{r_2} \| \cdots \| e_{i'}^{r_M} \tag{8}$$

Note that we adopt concatenation rather than weighted sum operation (i.e., sum, mean, attention), which aims to solve the problem of semantic pollution. More specfically, equation 8 shows that one attribute is represented within a specific range of element-wise positions. In addition, one type of attribute is placed in the same element-wise postions for all entities. For example, two representations $e_1^{(0)} = e_{1'}^{r_1} \| \cdots \| e_7^{r_3} \| \cdots$ and $e_{11}^{(0)} = e_{11'}^{r_1} \| \cdots \| e_7^{r_3} \| \cdots$ mean movie $e_1$ and movie $e_{11}$ have the same director and different actors. In a word, a concatenated representation avoids the interaction and preserves semantic independence of different attributes, which will be further advantageous to maintain the weight of siginificant high-order neighbors in subsection 3.3 and distill user interested attributes in subsection 3.4.

### 3.3 Attribute Propagation Layer

When we obtain initial entity representations via one-hop neighbors, an intuitive idea is to propagate them via GNN framework, so as to aggregate high-order neighbors into the center node and acquire more informative item representations.

Here we regard a KG as a heterogeneous graph and adopt a widely used two-step scheme[5, 36, 44] to aggregate the representations of neighbors: (1) same relation-aware neighbors aggregation; (2) relations combination.

The same relation-aware neighbors contains the similiar attribute types. For example, in figure 2, both $e_2$ and $e_3$ are actors and contain friend attribute while $e_5$ is a song and contains singer attribute. Therefore, we firstly aggregate representations of the same relation-aware neighbors. Secondly, to learn a more comprehensive representation, we need to fuse multiple attributes hold in different relation-aware neighbors. The pooling methods in the above two
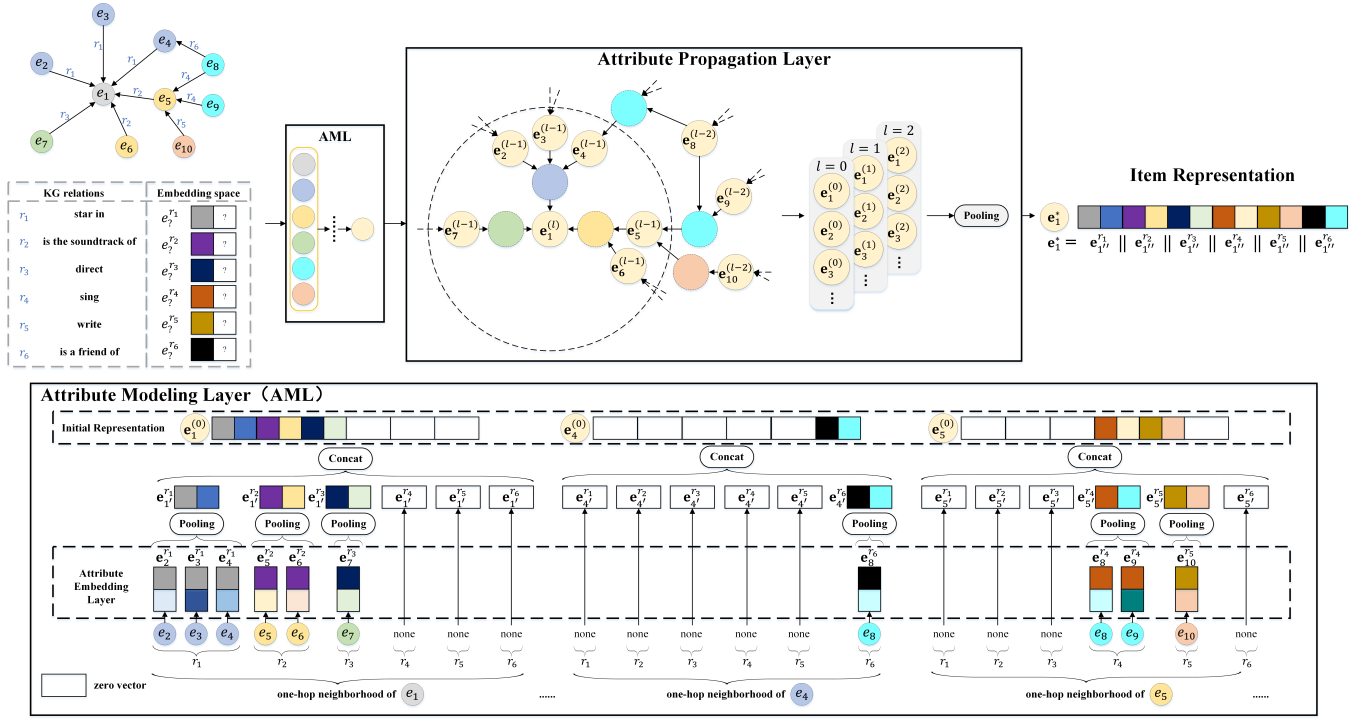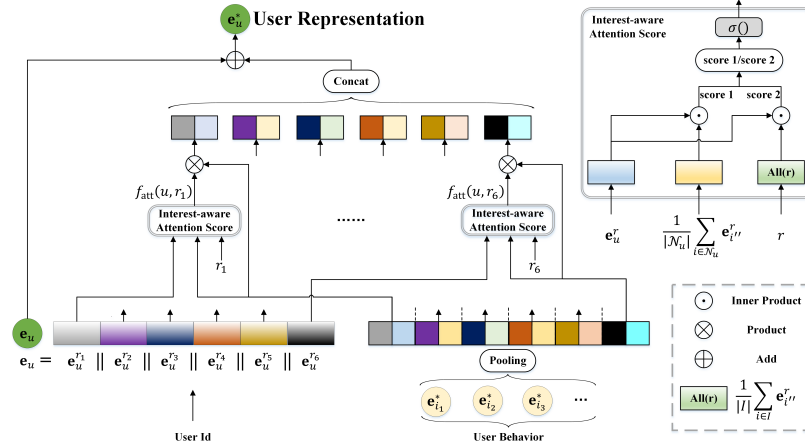
Figure 2: The structure of KGAN.



Figure 3: The structure of IAAN.

steps are average and sum operation, respectively. And we leave the further exploration of other pooling methods like attention as the future work. More formally, in the $l-$th layer, we recursively formulate the representation of an entity as:

$$\mathrm{e}_i^{(l)} = \sum_{r_j \in \mathcal{R}} \frac{1}{|\mathcal{N}_i^{r_j}|} \sum_{j \in \mathcal{N}_i^{r_j}} \mathrm{e}_j^{(l-1)} \qquad (9)$$

Now the final item representation $\mathrm{e}_i$ has been learned by equation 3, let's check how KGAN avoids pollution and weight decrease of significant high-order neighbors.

We re-examine $\mathrm{e}_i^*$ from the perspective of element-wise position. Without loss of generality, $\mathrm{e}_i$ is contructed as

$$\mathrm{e}_i^* = \mathrm{e}_{i''}^{r_1} \parallel \mathrm{e}_{i''}^{r_2} \parallel \cdots \parallel \mathrm{e}_{i''}^{r_M} \qquad (10)$$

where $\mathrm{e}_{i''}^{r_m} \in \mathbb{R}^{d^m}$ is a truncated vector in $\mathrm{e}_i^* \in \mathbb{R}^D$ and $D = \sum_{r_m \in \mathcal{R}} d^m$. The start and ending index of $\mathrm{e}_{i''}^{r_m}$ in $\mathrm{e}_i^*$ are $\sum_{p=1}^{m-1} d^p$ and $\sum_{p=1}^{m} d^p$, respectively. Reviewing the generation process of $\mathrm{e}_{i''}^{r_m}$, we can see that $\mathrm{e}_{i''}^{r_m}$ is learned by $i$'s neighbors' embeddings in attribute

$r_m$ space and doesn't contain any embeddings from other attribute spaces. For example, in figure 2, $e_{1''}^{r_1}$ is learned by $e_2^{r_1}$, $e_3^{r_1}$, and $e_4^{r_1}$. This means $e_{i''}^{r_m}$ maintains the independence of attribute $r_m$ and KGAN remains semantics unpolluted after multi-layer propagations. Furthermore, we check sepecific neighbors which are aggregated into $e_{i''}^{r_m}$. Take $e_{1''}^{r_5}$ as an exapmle, after 1-order propagation (one AML), $e_{1''}^{r_5}$ is a zero vector since attribute $r_5$ doesn't occur in one-hop neighbors of movie $e_1$. Then after 2-order propagation (one AML and one APL), $e_{1''}^{r_5} = e_{10}^{r_5}$, which seems as if there were a link $(e_{10}, r_5, e_1)$ and $e_1$ connected to $e_{10}$ directly. In general, when one relation, a.k.a attribute (i.e., $r_5$), firstly appear in the receptive field of center node (i.e., $e_1$) at $l-$hop (i.e., $2-$hop) position, the weight of its corresponding node (i.e., $e_{10}$) will not been decreased, which proves that KGAN maintains the weight of significant high-order neighbors.

## 3.4 Interest-aware Attention Layer

After item representations have been obtained by KGAN, a typical idea in recommender system is to enhance user representations by clicked items, like [16, 29, 39]. We take avearage pooling as an example as

$$e_u^* = e_u + \frac{1}{\mid \mathcal{N}_u \mid} \sum_{i \in \mathcal{N}_u} e_i^* \tag{11}$$

where $\mathcal{N}_u = \{i \mid (u, i) \in O\}$ and $e_u$ represents user embedding for collaborative filtering. However, avearage pooling doesn't consider user preferences personally, thereby recent works aims to model user interests via attention mechanism, such as [17, 22, 33, 35, 49].

Here we propose a novel attention module to model user interests towards different attributes. We assign an interest score $f_{\text{att}}(u, r_m)$ to each pair of attribute $r_m$ and user $u$, and personally generate user representation by combining interest scores and interacted items. We firstly introduce how to calculate $f_{\text{att}}(u, r_m)$ and then illustrate how to combine $f_{\text{att}}(u, r_m)$ with interacted items.

With different attributes placed in corresponding element-wise positions as shown in equation 10, we truncate user vector $e_u$ with the same strategy as

$$e_u = e_u^{r_1} \parallel e_u^{r_2} \parallel \cdots \parallel e_u^{r_M} \tag{12}$$

where $e_u^{r_m} \in \mathbb{R}^{d^m}$ is a truncated vector in $e_u \in \mathbb{R}^D$. The start and ending index of $e_u^{r_m}$ in $e_u$ are $\sum_{p=1}^{m-1} d^p$ and $\sum_{p=1}^{m} d^p$, respectively. Then the interest score of user $u$ towards attribute $r_m$ is calculated by

$$f_{\text{att}}(u, r_m) = \sigma \left( \tau \frac{\frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} e_u^{r_m \top} e_{i''}^{r_m}}{\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} e_u^{r_m \top} e_{i''}^{r_m}} \right) \tag{13}$$

where $\tau$ is temperature coefficient as a hyperparameter. In equation 13, $\frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} e_u^{r_m \top} e_{i''}^{r_m}$ denotes the interest-degree of user $u$ for his interacted items in attribute $r_m$, while $\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{N}_u} e_u^{r_m \top} e_{i''}^{r_m}$ denotes the interest-degree of user $u$ for all items in attribute $r_m$. If user $u$ attachs great importance to attribute $r_m$ when choosing items, the corresponding truncated representation $e_{i''}^{r_m}$ of his interacted items will be radically different from that of other uninterested items, such that $\frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} e_u^{r_m \top} e_{i''}^{r_m}$ will be much greater than $\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{N}_u} e_u^{r_m \top} e_{i''}^{r_m}$, and vice versa. Therefore, the ratio between

the above two expressions denotes the interest-degree of user $u$ for attribute $r_m$. Since a negative value of ratio is meaningless, we first feed this ratio to $\mathrm{Relu}$ and then select $\tanh$ as the nonlinear activation function, therefore $\sigma(x) = \tanh(\mathrm{Relu}(x))$.

After $f_{\text{att}}(u, r_m)$ has been prepared, we need to combine it with interacted items to learn user representations. We firstly re-express equation 11 as

$$e_u^* = e_u + \underset{r_m \in \mathcal{R}}{\parallel} \left( \frac{1}{\mid \mathcal{N}_u \mid} \sum_{i \in \mathcal{N}_u} e_{i''}^{r_m} \right) \tag{14}$$

where $\parallel$ is the concatenation operation. Therefore an intuitive idea is to use $f_{\text{att}}(u, r_m)$ to control how much attribute information, a.k.a. $\frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} e_{i''}^{r_m}$, will be passed to user. Consider the limit case, when $f_{\text{att}}(u, r_m) = 0$ which means user $u$ doesn't pay attention to attribute $r_m$ at all, any information of attribute $r_m$ should not be contained in user representaion. Formally, user representation is obtained by

$$e_u^* = e_u + \underset{r_m \in \mathcal{R}}{\parallel} \left( \frac{f_{\text{att}}(u, r_m)}{\mid \mathcal{N}_u \mid} \sum_{i \in \mathcal{N}_u} e_{i''}^{r_1} \right) \tag{15}$$

Note that here we relax the constraint that the sum of attention weights towards all attributes is 1, a.k.a. $\sum_{r_m \in \mathcal{R}} f_{\text{att}}(u, r_m) \neq 1$. The reason is as follows: when the number of members participating in the attention calculation is large and the limitation that the sum of attention weights is 1 is still reserved at this time, it will cause the attention weight of each member to be dispersed, making it difficult to learn the coefficients of important nodes. This problem has also appeared in Graphair [11], which shows estimating $O(\mid \mathcal{N} \mid^2)$ coefficients exposes the risk of overfitting. Therefore, we introduce the above novel activation unit to release the limitation that the sum of attention weight is 1, and the same idea is adopted by DIN [50].

## 3.5 Model Optimization

With the user representation $e_u^*$ and the item representation $e_i^*$ ready, equation 5 is adopted to calculate the prediction score of each pair of user and item. Then we employ the BPR loss [24] to encourage that the observed interactions should be assigned higher prediction values than unobserved ones. The objective function is formulated as

$$\mathcal{L} = \sum_{(u, i^+, i^-) \in O} - \ln \sigma \left( \hat{y}(u, i^+) - \hat{y}(u, i^-) \right) + \lambda \|\Theta\|_2^2 \tag{16}$$

where $O = \{(u, i^+, i^-) \mid (u, i^+) \in O^+, (u, i^-) \in O^-\}$ denotes the training set, which contains the observed interactions $O^+$ and the unobserved interactions $O^-$; $\sigma(\cdot)$ is the sigmoid function. $\Theta = \{e_i^{r_m}, e_u \mid i \in \mathcal{I}, u \in \mathcal{U}, r_m \in \mathcal{R}\}$ is the model parameter set. $L_2$ regularization parameterized by $\lambda$ on $\Theta$ is conducted to prevent overfitting. We employ the Adam [15] optimizer and use it in a mini-batch manner.

## 4 EXPIREMENT

We evaluate our proposed AKGAN method on three benchmark datasets to answer the following research questions:

**Table 1: Statistics of the datasets. 'Int./user' indicates the average number of interactions per user.**

| Datasets | Amazon-book | Last-FM | Alibaba-iFashion |
|---|---|---|---|
| # users | 70,679 | 23,566 | 114,737 |
| # items | 24,915 | 48,123 | 30,040 |
| # interactions | 847,733 | 3,034,796 | 1,781,093 |
| # Int./user | 11.99 | 128.78 | 15.52 |
| # entities | 88,572 | 58,266 | 59,156 |
| # relations | 39 | 9 | 51 |
| # triples | 2,557,746 | 464,567 | 279,155 |

- **RQ1**: Does our proposed AKGAN outperform the state-of-the-art recommendation methods?
- **RQ2**: How do different components (i.e., attribute modeling layer, attribute propagation layer, and interest-aware attention layer) affect AKGAN?
- **RQ3**: Can AKGAN provide potential explanations about user preferences towards attributes?

## 4.1 Experimental Settings

**Dataset Description**. We choose three benchmark datasets to evaluate our method: Amazon-Book[1], Last-FM[2], and Alibaba-iFashion[3]. The former two datasets are released in [33] and the last one is released in [35], and all of them are publicly available. Each dataset consists of two parts: user-item interactions and a corresponding knowledge graph. The basic statistics of the three datasets are presented in Table 1. We follow the same data partition used in [33, 37] to split the datasets into training and testing sets. For each observed user-item interaction, we randomly sample one negative item that the user has not interacted with before, and pair it with the user as a negative instance.

**Evaluation Metrics**. We evaluate our method in the task of top-K recommendation. For each user, we treat all the items that the user has not interacted with as negative and the observed items in the testing set as positive. Then we rank all these items and adopt two widely-used evaluation protocols: Recall@$K$ and NDCG@$K$, where $K$ is set as 20 by default. We report the average metrics for all users in the testing set.

**Baselines**. To demonstrate the effectiveness, we compare AKGAN with KG-free (MF), embedding-based (CKE), path-based (RippleNet), and GNN-based (KGAT, KGNN-LS, KGIN) methods:

- **MF** [24]: This is matrix factorization optimized by the Bayesian personalized ranking (BPR) loss, which only considers the user-item interactions and leaves KG untouched.
- **CKE** [45]: This method uses TransR [18], a typical knowledge graph embedding algorithm, to regularize the representations of items, which are fed into MF framework for recommendation.

[1] http://jmcauley.ucsd.edu/data/amazon
[2] https://grouplens.org/datasets/hetrec-2011/
[3] https://drive.google.com/drive/folders/1xFdx5xuNXHGsUVG2VIohFTXf9S7G5veq

- **RippleNet** [28]: This model combines embedding-based methods and path-based methods to propagate users' preferences on the KG for recommendation. RippleNet first assigns entities in the KG with initial embeddings using TransE and represents a user via entities related to his historically clicking items.
- **KGAT** [33]: This method encodes user behaviors and item knowledge as an unified knowledge graph to exploit high-order connectivity. KGAT applies an attentive neighborhood aggregation mechanism on a holistic graph and introduces TransR to regularize the representations.
- **KGNN-LS** [31]: It uses a user-specific relation scoring function to transform a heterogeneous KG into a user-personalized weighted graph and employs label smoothness regularization to avoid overfitting of edge weights.
- **KGIN** [35]: KGIN is the state-of-the-art GNN-based recommender. It models each intent as an attentive combination of KG relations to explore intents behind user-item interactions and adopts a novel relational path-aware aggregation scheme.

**Parameter Settings**. We implement our AGKAN model in Pytorch and Deep Graph Library (DGL)[4], which is a Python package for deep learning on graphs. We released all implementations (code, datasets, parameter settings, and training logs) to facilitate reproducibility. The embedding size of one attribute space in AGKAN varies from 4 to 64, and the detailed design will be introduced in Appendix A.1. For a fair comparison, we fix the size of ID embeddings as 64, which equals the max embedding size of AKGAN, for all baselines, except RippleNet 32 due to its high computational cost. We adopt Adam [15] as the optimizer and the batch size is fixed at 1024 for all methods. We use the Xavier initializer [6] to initialize model parameters. We apply a grid search for hyper-parameters: the learning rate is searched in $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$, the coefficient of L2 normalization is tuned in $\{10^{-5}, 10^{-4}, \cdots, 10^{-1}\}$, the number of GNN layers $L$ is searched in $\{1, 2, 3\}$ for GNN-based methods, and the dropout ratio is tuned in $\{0.0, 0.1, \cdots, 0.9\}$. Besides, we use the node dropout technique for KGAT, KGIN, and AKGAN, where the ratio is searched in $\{0.0, 0.1, \cdots, 0.9\}$. For RippleNet, we set the number of hops as 2, and the memory size as 5, 15, 8 for Alibaba-iFashion, Last-FM, and Amazon-book respectively to obtain the best performance. Since RippleNet is a model for CTR prediction, we generate top-K items with top-K scores in all items, which are compared with the test set to compute Recall@$K$ and NDCG@$K$. For KGAT, we use the pre-trained ID embeddings of MF as the initialization, which is also adopted by KGIN. Moreover, early stopping strategy is performed, i.e., premature stopping if Recall@20 on the test set does not increase for 10 successive epochs.

## 4.2 Performance Comparison(RQ1)

We report the empirical results in Table 2 where we highlight the results of the best baselines (starred) and our AKGAN (boldfaced). And we also use %Imp. to denote the percentage of relative improvement on each metric. The observations are as followed:

- AKGAN consistently achieves the best performance on three datasets in terms of all measures. Specifically, it achieves significant improvements over the strongest baselines w.r.t. NDCG@20

[4] https://github.com/dmlc/dgl

**Table 2: Overall Performance Comparison. The best performance is boldfaced; the runner up is labeled with '*'. '%Imp.' indicates the improvements.**

| Dataset | Metrics | KG-free | embedding-based | path-based | GNN-based | | | | Imp. |
|---|---|---|---|---|---|---|---|---|---|
| | | MF | CKE | RippleNet | KGAT | KGNN-LS | KGIN | AKGAN | |
| Amazon-Book | Recall | 0.1241 | 0.1287 | 0.1355 | 0.1473 | 0.1389 | 0.1687* | **0.1783** | 5.69% |
| | NDCG | 0.0650 | 0.0674 | 0.0763 | 0.0782 | 0.0614 | 0.0915* | **0.0994** | 8.63% |
| Last-FM | Recall | 0.0774 | 0.0780 | 0.0842 | 0.0876 | 0.0877 | 0.0978* | **0.1209** | 23.62% |
| | NDCG | 0.0669 | 0.0659 | 0.0766 | 0.0745 | 0.0653 | 0.0848* | **0.1066** | 25.71% |
| Alibaba-iFashion | Recall | 0.0921 | 0.1068 | 0.1121 | 0.1015 | 0.1046 | 0.1147* | **0.1253** | 9.24% |
| | NDCG | 0.0562 | 0.0633 | 0.0695 | 0.0616 | 0.0582 | 0.0716* | **0.0801** | 11.87% |

by 8.63%, 25.71%, and 11.87% in Amazon-Book, Last-FM, and Alibaba-iFashion, respectively. These improvements are attributed to the following reasons: (1) By combining all attributes using concatenation operation, AKGAN avoids semantic pollution caused by weighted sum operation and learns more high-quality item representations for recommendation. (2) Compared to GNN-based baselines (i.e., KGAT, KGNN-LS, KGIN), AKGAN maintains the weight of significant high-order neighbors by placing different attributes in corresponding element-wise positions. (3) Benefiting from our novel interest-aware attention module that assigns an interest score to each pair of user and attribute, AKGAN can recognize the pattern of user interest at a fine-grained level to conduct a better personal recommendation.

- Jointly analyzing AKGAN across the three datasets, we find that the improvement on Last-FM is more significant than that on Alibaba-iFashion and Amazon-Book. The main reason is that the interaction number per user of Last-FM (128.78) is much larger than that of the other two datasets (11.99, 15.52). Therefore, there exists richer interaction information on the Last-FM dataset for AKGAN to refine user and item representation by collaborative signals. This indicates that AKGAN will fully realize its potential in recommendation scenarios with dense interaction data.
- KG-free method (MF) underperforms knowledge-aware methods (i.e., CKE, RippleNet, AKGAN). A clear reason is MF doesn't leverage the rich attribute information in KG.
- GNN-based methods (i.e., KGAT, KGNN-LS, KGIN, AKGAN) achieve better performance than path-based (RippleNet) and embedding-based (CKE) methods. A possible reason is that these three kinds of recommenders adopt different usages of attributes. GNN-based methods aggregate neighbors' attribute information into item nodes to learn more informative representations. The other two methods have a common limitation that both of them don't break loose from employing knowledge graph embedding algorithms to model attribute information and regularize node representations.
- In four GNN-based methods, AKGAN performs best, KGIN is the second-best, while KGAT and KGNN-LS are at the same level and achieve the worst results. The decreasing performance is because that the level of how a recommender learns item attributes and user interests is in descending order, from fine-grained to coarse-grained. AKGAN uses concatenation operation to avoid attribute interaction while other models adopt weighted sum (attentive combination) operation. AKGAN maintains the weight of significant high-order neighbors while other models

**Table 3: GNN-based models Comparison: (1) CM - Combination methods of different attributes (Con - concatenation, WS - weighted sum); (2) WD - whether to avoid weight decrease of significant high-order neighbors; (3) GF - the GNN framework that a recommender adopts (Het - heterogeneous design in different layers, Hom - identical network in each layer; (4) IA - whether to learn user interests towards attributes; (5) level - the level of how a recommender learns item attributes and user interests (FG - fine-grained, CG - coarse-grained).**

| Model | CM | | WD | GF | | IA | level |
|---|---|---|---|---|---|---|---|
| | Con | WS | | Het | Hom | | |
| AKGAN | √ | × | √ | √ | × | √ | FG ↓ CG |
| KGIN | × | √ | × | × | √ | √ | |
| KGNN-LS | × | √ | × | × | √ | √ | |
| KGAT | × | √ | × | × | √ | × | |

neglect this issue. To achieve the above two advantages, AKGAN adopts a heterogeneous design in different layers while the others use a homogeneous GNN framework. AKGAN, KGIN, and KGNN-LS all learn user interests towards attributes explicitly while KGAT doesn't. The above comparison is listed in Table 3.

## 4.3 Study of AKGAN(RQ2)

In this section, we first conduct an ablation study to investigate the effect of cancatenation operation and interest-aware attention layer. Towards the further analysis, we study the influence of layer numbers. In what follows, we explore how the hyperparameter, a.k.a. temperature coefficient, affects the performance.

**Impact of cancatenation operation & interest score**. To demonstrate the necessity of cancatenation oepration and interest score, we compare the performance of AKGAN with the following three variants: (1) combing different attributes with average operation and discarding interest score, termed AKGAN-mean, (2) combing different attributes with sum operation and discarding interest score, termed AKGAN-sum, (3) only discarding interest score, termed AKGAN-noatt. Discarding interest score means we use equation 11 to learn user representations. Note that we don't adopt weighted sum operation and interest-aware attention layer simultaneously. The reason is that cancatenation operation is the prerequisite of interest-aware attention layer, therefore we have to discard interest

**Table 4: Impact of cancatenation operation and interest score.**

| | Amazon-Book | | Last-FM | | Alibaba-iFashion | |
|---|---|---|---|---|---|---|
| | Recall | NDCG | Recall | NDCG | Recall | NDCG |
| AKGAN-mean | 0.1504 | 0.0799 | 0.0849 | 0.0713 | 0.1064 | 0.0660 |
| AKGAN-sum | 0.1489 | 0.0784 | 0.0886 | 0.0736 | 0.1045 | 0.0643 |
| AKGAN-noatt | 0.1768 | 0.0969 | 0.1141 | 0.1012 | 0.1173 | 0.0739 |

**Table 5: Impact of the number of layers $L$.**

| | Amazon-Book | | Last-FM | | Alibaba-iFashion | |
|---|---|---|---|---|---|---|
| | Recall | NDCG | Recall | NDCG | Recall | NDCG |
| AKGAN-1 | 0.1737 | 0.0974 | 0.1192 | 0.1060 | 0.1245 | 0.0791 |
| AKGAN-2 | 0.1752 | 0.0977 | 0.1205 | 0.1069 | 0.1253 | 0.0801 |
| AKGAN-3 | 0.1783 | 0.0994 | 0.1209 | 0.1066 | 0.1221 | 0.0776 |

score when testing weighted sum operation. The results are shown in Table 4 and we summarize the major findings as below:

- Comparing AKGAN-mean/AKGAN-sum with AKGAN-noatt, we can clearly see that replacing concatenation operation with weighted sum operation dramatically degrades performance of recommendation, which indicates the superiority of cancatenation operation.
- The improvement from AKGAN-noatt to AKGAN verifies the necessity of interest score.
- Jointly comparing AKGAN-noatt and AKGAN across the three datasets, we find that the improvement on Last-FM is more significant than that on Alibaba-iFashion and Amazon-Book, which is consistent with the aforesaid conclusion in subsection 4.2. The main reason is that Last-FM has more dense interaction data than the other two datasets, such that AKGAN can better learn user interest in Last-FM.

**Impact of model depth**. We investigate the influence of depth of receptive field in AKGAN by searching $L$ in the range of $\{1, 2, 3\}$. Particularly, $L = 1$ means AKGAN has only one attribute modeling layer. The results are reported in Table 5. In Amazon-Book and Last-FM, we can see that increasing propagation times of attributes can boost the performance, because more attribute information is aggregated into the center node to learn more informative user and item representations. While in Alibaba-iFashion, AKGAN-2 is the best and AKGAN-3 is the worst, which is caused by its inherent topology. Alibaba-iFashion just has two kinds of triplets: the first-order connectivity of an item is its components, a.k.a. *(fashion outfit, including, fashion staff)*, and the second-order connectivity is the category of staff, a.k.a. *(staff, having-category-?, ?)*. Therefore, all significant attribute information has been captured in the 2-hop range, leading to the best performance of AKGAN-2.

**Impact of temperature coefficient**. We vary the temperature coefficient in the range of $\tau = \{0.01, 0.1, 0.2, ..., 0.9, 1\}$ to study its influence on the performance of AKGAN. The experiment is ongoing.

## 4.4 Case Study(RQ3)

In this section, we visualize the interest score to show how AKGAN learns user preferences towards attributes. We choose Alibaba-iFashion as the example dataset since it helps to provide an intuitive explanation. We begin with introducing this dataset briefly for further better understanding. Alibaba-iFashion is an E-commerce dataset, which contains user-outfit click history for recommendation. Each outfit consists of several fashion staffs (e.g., tops, bottoms, shoes) and each staff is assigned with different categories. For example, trench coat, T-shirt and sweater are three kinds of tops, pants and long skirt are two kinds of bottoms. Alibaba-iFashion regards these categories as relations in KG and there is one more relation between outfit and staff, called *including*. Therefore, Alibaba-iFashion KG has two kinds of triples: *(outfit, including, staff)* and *(staff, having-category-?, ?)*, where ? is the specific category like T-shirt. According to Table 1, Alibaba-iFashion KG has 51 relations, where we number relation *including* as 0 and number 50 staff categories from 1 to 50. In figure 4, we first exhibit the training set and testing set of a specific user (user0). Then we visualize the interest score learned from the training set by AKGAN and use the testing set to validate the effectiveness of these scores. From figure 4 we can see that:

- The interest score accurately models the user preferences towards attributes. Among 11 outfits clicked by user0, the tops contain four main kinds: T-shirt, shirt, trench coat, and sweater, which means user0 prefers the above four types of tops. The corresponding interest score are 0.3254, 0.4379, 0.6793, 0.7256, respectively. The most clicked bottoms are women's pants which are assigned with a high score 0.8223, while the jeans and skirt occur less frequently and their scores, 0.0601 and 0.0671, are particularly small. As for shoes, user0 prefers women's shoes and boots only appear once, their scores, 0.7623 and 0.0431, also reflect this difference. The earrings and hat are two main accessories and both of them are assigned with high scores, 0.5506 and 0.3506.
- All outfits clicked by user0 have only one kind of handbag and the interest score $f_{\text{att}}(u_0, r_2) = 0.1506$, which is a small value. A possible reason is that outfits are manually created by Taobao's fashion experts who prefer to include handbags for the completion of the outfit [4]. Therefore user0 judges whether to click an outfit based on other staffs rather than the handbag. Another explanation is that handbag is not as frequently changed as other staffs like tops in daily dressing, such that user0 pays less attention to the handbag when browsing through outfits.
- Testing set proves the effectiveness of interest score. For example, $f_{\text{att}}(u_0, r_{40}) = 0.6793$ is proved to be reasonable since outfit12 which contains a trench coat appears in the testing set. And $f_{\text{att}}(u_0, r_9) = 0.8233$ is validated by outfit13 and outfit14 with same reason.
- In one staff, the interest degree varies according to category. Take accessories as an example, user0 prefers earrings than hat and $f_{\text{att}}(u_0, r_{23}) = 0.5506 > f_{\text{att}}(u_0, r_{13}) = 0.3506$. This result is also proved by the testing set, where two outfits contain earrings while only one outfit has a hat.

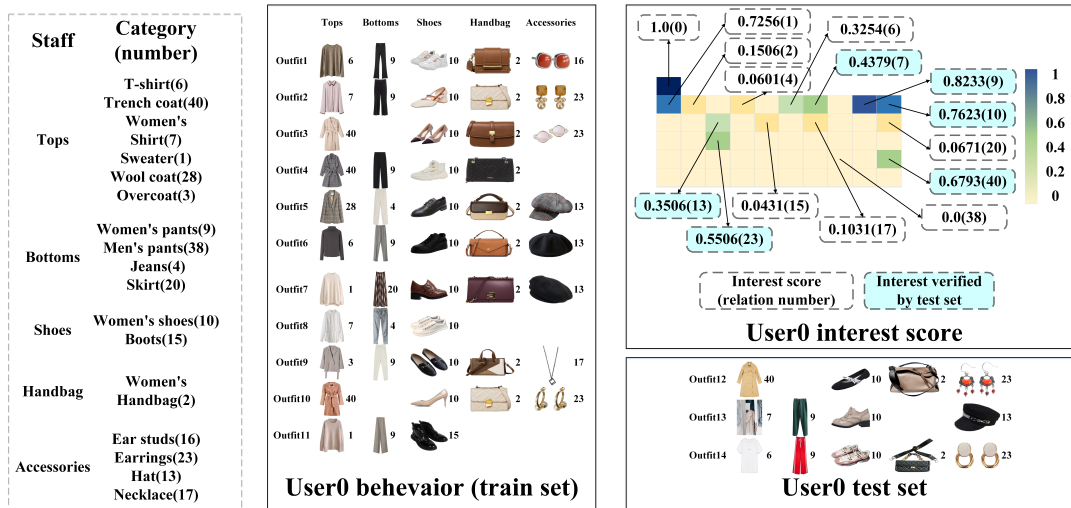**Figure 4: Case study.**

- Unconcerned attribute is assigned with a zero score, like $f_{\text{att}}(u_0, r_{38}) = 0$. The reason is that user0 perhaps is female since her click histories are all women's clothing and she isn't interested in men's pants.
- Relation *including* has the highest score 1. The reason is that this relation *(outfit, including, staff)* is the necessary bridge to acquire what staffs an outfit consists of. And the same reason has been reflected in the impact of model depth.

## 5 RELATED WORK

Our work is highly related with the knowledge-aware recommendation, which can be grouped into three categories.

**Embedding-based Methods** [1, 3, 12, 29, 45, 46] hire KG embedding algorithm (i.e., TransE [2] and TransH [38]) to model prior representations of item, which are used to guide the recommender model. For example, DKN [29] learns knowledge-level embedding of entities in news content via TransD [13] for news recommendation. KSR [12] utilizes knowledge base information learned with TransE as attribute-level preference to enhance the sequential recommendation.

**Path-based Methods** [19, 26, 41–43] usually predefines a path scheme (i.e., meta-path) and leverages path-level semantic similarities of entities to refine the representations of users and items. For example, MCRec [10] learns the explicit representations of meta-paths to depict the interaction context of user-item pairs. RKGE [27] mines the path relation between user and item automatically and encodes the entire path using a recurrent network to predict user preference towards this item.

**GNN-based Methods** [23, 25, 30, 33, 35, 48] utilizes the message-passing mechanism in graph to aggregate high-order attribute informations into item representation for enhanced and explainable recommendation. KGAT[33] regards user-item interaction as a new

relation added to KG and then employees attentive mechanism to propagate attribute information. IntentGC [48] reconstructs user-to-user relationships and item-to-item relationships based on KG and proposes a novel graph convolutional network to aggregate the attribute information from neighbors. KGIN [35] learn user interest via an attentive combination of attributes and integrates relational information from multi-hop paths to refine the representations.

## 6 CONCLUSION

In this paper, we study the attribute information in knowledge graphs intending to improve the recommendation performance. On the item side, the proposed AKGAN can learn more high-quality item representation by remaining the independency of attributes and maintaining the weight of high-order significant attribute nodes. On the user side, AKGAN mines user interests towards attributes and provides a personal recommendation. Extensive experiments demonstrate the effectiveness and explainability of AKGAN.

For future work, we plan to investigate how to model the evolving process of user interests towards attributes based on the long-term behavior sequence. Another direction is to explore whether attribute interaction will benefit recommendation since feature interaction has shown great success in click-through rate prediction.

## REFERENCES

[1] Qingyao Ai, Vahid Azizi, Xu Chen, and Yongfeng Zhang. 2018. Learning heterogeneous knowledge base embeddings for explainable recommendation. *Algorithms* 11, 9 (2018), 137.

[2] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems* 26 (2013).

[3] Yixin Cao, Xiang Wang, Xiangnan He, Zikun Hu, and Tat-Seng Chua. 2019. Unifying knowledge graph learning and recommendation: Towards a better understanding of user preferences. In *The world wide web conference*. 151–161.

[4] Wen Chen, Pipei Huang, Jiaming Xu, Xin Guo, Cheng Guo, Fei Sun, Chao Li, Andreas Pfadler, Huan Zhao, and Binqiang Zhao. 2019. POG: personalized outfit generation for fashion recommendation at Alibaba iFashion. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 2662–2670.

[5] Xinyu Fu, Jiani Zhang, Ziqiao Meng, and Irwin King. 2020. Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding. In *Proceedings of The Web Conference 2020*. 2331–2341.

[6] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 249–256.

[7] Qingyu Guo, Fuzhen Zhuang, Chuan Qin, Hengshu Zhu, Xing Xie, Hui Xiong, and Qing He. 2020. A survey on knowledge graph-based recommender systems. *IEEE Transactions on Knowledge and Data Engineering* (2020).

[8] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.

[9] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.

[10] Binbin Hu, Chuan Shi, Wayne Xin Zhao, and Philip S Yu. 2018. Leveraging metapath 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*. 1531–1540.

[11] Fenyu Hu, Yanqiao Zhu, Shu Wu, Weiran Huang, Liang Wang, and Tieniu Tan. 2021. Graphair: Graph representation learning with neighborhood aggregation and interaction. *Pattern Recognition* 112 (2021), 107745.

[12] Jin Huang, Wayne Xin Zhao, Hongjian Dou, Ji-Rong Wen, and Edward Y Chang. 2018. Improving sequential recommendation with knowledge-enhanced memory networks. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 505–514.

[13] Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. 2015. Knowledge graph embedding via dynamic mapping matrix. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 687–696.

[14] Yuchin Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. 2016. Field-aware factorization machines for CTR prediction. In *Proceedings of the 10th ACM conference on recommender systems*. 43–50.

[15] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[16] Yehuda Koren. 2008. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. 426–434.

[17] Chao Li, Zhiyuan Liu, Mengmeng Wu, Yuchi Xu, Huan Zhao, Pipei Huang, Guoliang Kang, Qiwei Chen, Wei Li, and Dik Lun Lee. 2019. Multi-interest network with dynamic routing for recommendation at Tmall. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 2615–2623.

[18] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning entity and relation embeddings for knowledge graph completion. In *Twenty-ninth AAAI conference on artificial intelligence*.

[19] Chen Luo, Wei Pang, Zhe Wang, and Chenghua Lin. 2014. Hete-cf: Social-based collaborative filtering recommendation using heterogeneous relations. In *2014 IEEE International Conference on Data Mining*. IEEE, 917–922.

[20] Hao Ma, Haixuan Yang, Michael R Lyu, and Irwin King. 2008. Sorec: social recommendation using probabilistic matrix factorization. In *Proceedings of the 17th ACM conference on Information and knowledge management*. 931–940.

[21] Xichuan Niu, Bofang Li, Chenliang Li, Rong Xiao, Haochuan Sun, Hongbo Deng, and Zhenzhong Chen. 2020. A dual heterogeneous graph attention network to improve long-tail performance for shop search in e-commerce. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3405–3415.

[22] Jiarui Qin, Weinan Zhang, Xin Wu, Jiarui Jin, Yuchen Fang, and Yong Yu. 2020. User behavior retrieval for click-through rate prediction. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2347–2356.

[23] Yanru Qu, Ting Bai, Weinan Zhang, Jianyun Nie, and Jian Tang. 2019. An end-to-end neighborhood-based interaction model for knowledge-enhanced recommendation. In *Proceedings of the 1st International Workshop on Deep Learning Practice for High-Dimensional Sparse Data*. 1–9.

[24] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).

[25] Xiao Sha, Zhu Sun, and Jie Zhang. 2019. Attentive knowledge graph embedding for personalized recommendation. *arXiv preprint arXiv:1910.08288* (2019).

[26] Chuan Shi, Zhiqiang Zhang, Ping Luo, Philip S Yu, Yading Yue, and Bin Wu. 2015. Semantic path based personalized recommendation on weighted heterogeneous information networks. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. 453–462.

[27] Zhu Sun, Jie Yang, Jie Zhang, Alessandro Bozzon, Long-Kai Huang, and Chi Xu. 2018. Recurrent knowledge graph embedding for effective recommendation. In *Proceedings of the 12th ACM Conference on Recommender Systems*. 297–305.

[28] Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. 2018. Ripplenet: Propagating user preferences on the knowledge graph for recommender systems. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. 417–426.

[29] Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. 2018. DKN: Deep knowledge-aware network for news recommendation. In *Proceedings of the 2018 world wide web conference*. 1835–1844.

[30] Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, and Zhongyuan Wang. 2019. Knowledge-aware graph neural networks with label smoothness regularization for recommender systems. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 968–977.

[31] Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, and Zhongyuan Wang. 2019. Knowledge Graph Convolutional Networks for Recommender Systems with Label Smoothness Regularization. CoRR abs/1905.04413 (2019). *arXiv preprint arXiv:1905.04413* (2019).

[32] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. 2018. Billion-scale commodity embedding for e-commerce recommendation in alibaba. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 839–848.

[33] Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and Tat-Seng Chua. 2019. Kgat: Knowledge graph attention network for recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 950–958.

[34] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*. 165–174.

[35] Xiang Wang, Tinglin Huang, Dingxian Wang, Yancheng Yuan, Zhenguang Liu, Xiangnan He, and Tat-Seng Chua. 2021. Learning Intents behind Interactions with Knowledge Graph for Recommendation. In *Proceedings of the Web Conference 2021*. 878–887.

[36] Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S Yu. 2019. Heterogeneous graph attention network. In *The World Wide Web Conference*. 2022–2032.

[37] Xiang Wang, Yaokun Xu, Xiangnan He, Yixin Cao, Meng Wang, and Tat-Seng Chua. 2020. Reinforced negative sampling over knowledge graph for recommendation. In *Proceedings of The Web Conference 2020*. 99–109.

[38] Zhigang Wang, Juanzi Li, Zhiyuan Liu, and Jie Tang. 2016. Text-enhanced representation learning for knowledge graph. In *Proceedings of International Joint Conference on Artificial Intelligent (IJCAI)*. 4–17.

[39] Le Wu, Peijie Sun, Yanjie Fu, Richang Hong, Xiting Wang, and Meng Wang. 2019. A neural influence diffusion model for social recommendation. In *Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval*. 235–244.

[40] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 974–983.

[41] Xiao Yu, Xiang Ren, Quanquan Gu, Yizhou Sun, and Jiawei Han. 2013. Collaborative filtering with entity similarity regularization in heterogeneous information networks. *IJCAI HINA* 27 (2013).

[42] Xiao Yu, Xiang Ren, Yizhou Sun, Quanquan Gu, Bradley Sturt, Urvashi Khandelwal, Brandon Norick, and Jiawei Han. 2014. Personalized entity recommendation: A heterogeneous information network approach. In *Proceedings of the 7th ACM international conference on Web search and data mining*. 283–292.

[43] Xiao Yu, Xiang Ren, Yizhou Sun, Bradley Sturt, Urvashi Khandelwal, Quanquan Gu, Brandon Norick, and Jiawei Han. 2013. Recommendation in heterogeneous information networks with implicit user feedback. In *Proceedings of the 7th ACM conference on Recommender systems*. 347–350.

[44] Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V Chawla. 2019. Heterogeneous graph neural network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 793–803.

[45] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. 2016. Collaborative knowledge base embedding for recommender systems. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 353–362.

[46] Yongfeng Zhang, Qingyao Ai, Xu Chen, and Pengfei Wang. 2018. Learning over knowledge-base embeddings for recommendation. *arXiv preprint arXiv:1803.06540* (2018).

[47] Yin Zhang, Derek Zhiyuan Cheng, Tiansheng Yao, Xinyang Yi, Lichan Hong, and Ed H Chi. 2021. A Model of Two Tales: Dual Transfer Learning Framework for Improved Long-tail Item Recommendation. In *Proceedings of the Web Conference 2021*. 2220–2231.

[48] Jun Zhao, Zhou Zhou, Ziyu Guan, Wei Zhao, Wei Ning, Guang Qiu, and Xiaofei He. 2019. Intentgc: a scalable graph convolution framework fusing heterogeneous information for recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2347–2357.

[49] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 5941–5948.

[50] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1059–1068.

## A APPENDIX

### A.1 Reproducibility Settings

In KG, the number of edges varies acoording to relation type. We can set $d^m$ as a fixed constant value which is the same as that all latent vectors of different fields share the same dimension in FFM[14]. However, KGs contain hundreds of relation types, which requires huge computational resources both in memory space and time. Table 6 shows an example on the Alibaba-iFashion dataset. Therefore, we design variable dimensions of different attribute spaces. The principle is that more edges belonging to one specific relation type bring in a larger dimension of this attribute space. Formally, we make the dimension increases linearly with the edge number as below:

$$d^m = \begin{cases} d_{min} + \dfrac{d_{max} - d_{min}}{c}|r_m|, & |r_m| \leq \dfrac{c}{\dfrac{d_{max} - d_{min}}{d_{max} - d_{min}}} \\ d_{max}, & |r_m| > \dfrac{c}{d_{max} - d_{min}} \end{cases} \quad (17)$$

where $d_{max}, d_{min}, c, c$ are hyperpatameters and $|r_m|$ is the number of edges belonging to relation $r_m$, like $|r_0| = 260477$ in Alibaba-iFashion.

**Table 6: Edge number of Alibaba-iFashion dataset**

| # $r_1$ - $r_{10}$ | 260,477 469 2518 482 1177 343 429 618 143 865 |
|---|---|
| # $r_{11}$ - $r_{20}$ | 2443 1572 155 939 221 934 244 187 162 93 |
| # $r_{21}$ - $r_{30}$ | 833 161 129 736 22 61 318 142 435 36 |
| # $r_{31}$ - $r_{40}$ | 296 459 19 102 73 23 53 321 63 45 |
| # $r_{41}$ - $r_{51}$ | 132 7 10 40 42 15 78 11 11 5 6 |

We list the parameter settings of AKGAN on three datasets in Table 7, where the hyperparameters include the learning rate $l_r$, the coefficient $\lambda$ of $L_2$ regularization, the temperature coefficient $\tau$, and three coefficients about dimension $d_{max}, d_{min}, c$. We have released our codes, datasets, model parameters, and training logs at https://github.com/huaizepeng2020/AKGAN.

**Table 7: Hyperparameter settings of AKGAN**

| | $l_r$ | $\lambda$ | $\tau$ | $d_{max}$ | $d_{min}$ | $c$ |
|---|---|---|---|---|---|---|
| Amazon-Book | $10^{-4}$ | $10^{-5}$ | 0.25 | 32 | 4 | 5000 |
| Last-FM | $10^{-4}$ | $10^{-5}$ | 0.5 | 64 | 16 | 5000 |
| Alibaba-iFashion | $10^{-4}$ | $10^{-5}$ | 0.1 | 64 | 4 | 5000 |