

Set2setRank: Collaborative Set to Set Ranking for Implicit Feedback based Recommendation

Lei Chen^{1,2}, Le Wu^{1,2,3,*}, Kun Zhang^{1,2}, Richang Hong^{1,2}, Meng Wang^{1,2,3}

¹ Key Laboratory of Knowledge Engineering with Big Data, Hefei University of Technology, China

² School of Computer Science and Information Engineering, Hefei University of Technology, China

³ Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, China

{chenlei.hfut, lewu.ustc, zhang1028kun, hongrc.hfut, eric.mengwang}@gmail.com

ABSTRACT

As users often express their preferences with binary behavior data (implicit feedback), such as clicking items or buying products, implicit feedback based Collaborative Filtering (CF) models predict the top ranked items a user might like by leveraging implicit user-item interaction data. For each user, the implicit feedback is divided into two sets: an observed item set with limited observed behaviors, and a large unobserved item set that is mixed with negative item behaviors and unknown behaviors. Given any user preference prediction model, researchers either designed ranking based optimization goals or relied on negative item mining techniques for better optimization. Despite the performance gain of these implicit feedback based models, the recommendation results are still far from satisfactory due to the sparsity of the observed item set for each user. To this end, in this paper, we explore the unique characteristics of the implicit feedback and propose Set2setRank framework for recommendation. The optimization criteria of Set2setRank are two folds: First, we design an item to an item set comparison that encourages each observed item from the sampled observed set is ranked higher than any unobserved item from the sampled unobserved set. Second, we model set level comparison that encourages a margin between the distance summarized from the observed item set and the most “hard” unobserved item from the sampled negative set. Further, an adaptive sampling technique is designed to implement these two goals. We have to note that our proposed framework is model-agnostic and can be easily applied to most recommendation prediction approaches, and is time efficient in practice. Finally, extensive experiments on three real-world datasets demonstrate the superiority of our proposed approach.

CCS CONCEPTS

• **Information systems** → **Recommender systems; Learning to rank.**

Le Wu is the corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '21, July 11–15, 2021, Virtual Event, Canada

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8037-9/21/07...\$15.00

<https://doi.org/10.1145/3404835.3462886>

KEYWORDS

recommendation; implicit feedback; collaborative ranking; adaptive sampling; self-supervised ranking

ACM Reference Format:

Lei Chen^{1,2}, Le Wu^{1,2,3,*}, Kun Zhang^{1,2}, Richang Hong^{1,2}, Meng Wang^{1,2,3}. 2021. Set2setRank: Collaborative Set to Set Ranking for Implicit Feedback based Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21), July 11–15, 2021, Virtual Event, Canada*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3404835.3462886>

1 INTRODUCTION

Collaborative Filtering (CF) provides personalized ranking list for each user by leveraging the collaborative signals from user-item interaction data, and are popular in most recommender systems with easy to collect data and relatively high performance [3, 40, 43, 46]. Earlier approaches worked on the explicit user feedback data (e.g., 1-5 rating scale), and relied on pointwise optimization functions by comparing the error of each predicted rating and the real rating value [19, 29, 53]. Nevertheless, in real-world applications, most feedback from users are presented in the implicit manner, such as clicking a product, purchasing an item, and visiting a restaurant. The implicit feedback is more common and much easier to collect than the explicit feedback. Therefore, researchers focus on the implicit feedback based recommendation.

Different from the explicit rating values, in implicit feedback recommendation scenario, we only have users' observed behaviors while the large majority of users' preferences are unknown [25, 27]. As the unknown records are much larger than the observed feedbacks, some researchers proposed to use weighted regression optimization, e.g., weighted matrix factorization to assign lower weights for unobserved behaviors than the observed behaviors [17]. Since the task of recommendation is to provide a personalized ranking to users, a natural idea is to use ranking based optimization. The current ranking based optimization models can be mainly categorized into pairwise and listwise models. Pairwise models are based on the pairwise comparison between a relevant item and an irrelevant item at each time [27, 31, 51]. E.g., the most widely used pairwise method of Bayesian Personalized Ranking, assumes that a user prefers an observed item than a randomly selected unobserved item [31, 44]. Most listwise models optimize the top-N ranking oriented measure or maximize the permutation probability of the most likely permutation of the defined list, such as the overall item set [14, 47], or the list that is composed of an observed item and unobserved items [45]. As computing top-N probability grows exponentially with N, nearly all listwise models simplify N as 1

for efficiency consideration. These ranking based models treat the recommendation as a ranking task, and show better performance than pointwise metrics for recommendation.

In these ranking based methods, as the size of the unobserved items are much larger than the size of the observed items, researchers proposed to randomly sample a portion of items from the unobserved set as negative items at each iteration, in order to reduce time complexity without much performance loss. Instead of the random selection, researchers proposed to select negative items based on heuristics, such as sampling items based on the popularity [10, 36, 43]. Another kind of models argued that different unobserved items have different importance in the ranking process, instead of random selection or heuristic selection models, researchers designed models to dynamically choose negative training samples from the ranked list produced in the current prediction model [4, 11, 54, 55]. E.g., researchers proposed a non-uniform sampling distribution that adapts to both the context and the current state of learning [30]. The negative item sampling techniques further improve collaborative ranking performance by mining the valuable samples from the large set of unobserved item behaviors.

Despite the advances of collaborative ranking based approaches for implicit feedback, we argue that performance of current ranking models for CF are still far from satisfactory. For each user, the implicit feedback are divided into two sets: an observed item set with limited observed behaviors, and a large unobserved item set that is mixed with negative and unknown behaviors. Each user's observed item size is very small, and usually far less than the size of unobserved items. E.g., the density of most user-item interaction matrix is less than 1%. How to well exploit the structure information hidden in two sets, especially the observed set is the key challenge. For most ranking based CF models, the limited observed set is only used at element wise level independently, while the relationships of observed items have not been well exploited. Specifically, pairwise approaches take each observed-unobserved item pair with independence assumption of each selected pair. Listwise approaches resort to top-1 ranking for efficiency consideration and could not model observed set from a global perspective. Given the sparse observed set and the large unobserved set, can we better learn patterns between two sets, such that more self-supervised information hidden in users' behaviors can be learned for recommendation?

In this paper, we design Set2setRank for collaborative ranking with implicit feedback. The key idea of our proposed ranking framework is that: instead of treating each element in the observed set independently, we explore how to better exploit the structure hidden in each set, and the structure information between two sets for implicit feedback. To this end, we design a novel Set2setRank approach to better exploit the structure information hidden in the observed set and unobserved set to guide implicit feedback based ranking. At each iteration, we sample an observed set and an unobserved set from implicit feedback. We design a newly two-level comparison to learn more self-supervised information hidden in users' behaviors. The first level is an item to an item set comparison that encourages each sampled unobserved item is ranked lower than any sampled observed items. The second level is a set to set comparison to encourage a margin between distance of observed items summarized from the observed set, and distance of the most "hard" item selected from the first level comparison. Moreover, we

design an adaptive sampling algorithm to implement our proposed two-level comparison, where the size of the two sampled sets vary at each iteration. We have to note that our proposed framework is model-agnostic and can be easily applied to most recommendation prediction approaches, and is time efficient in practice. Finally, extensive experiments on three real-world datasets demonstrate the superiority of our proposed approach.

2 RELATED WORK

CF exploits the collaborative signals from user-item interaction behaviors for recommendation [3, 40, 43]. Among all models for CF, learning user and item embeddings have been popular for modern recommender systems, as embedding learning shows flexibility and relative high performance [20, 40]. After that, the preference of a user to an item is predicted by the inner product of the user and item embedding. Most research works have been focused on designing user and item embedding architecture. Matrix factorization based approaches learn user and item embedding based on low rank decomposition of the user-item interaction matrix [16, 48]. Recently, due to the huge success of graph neural networks, researchers have treated the user-item behaviors as a user-item bipartite graph, and designed neural graph models [5, 42]. E.g., NGCF is one of the first few attempts that designed embedding propagation to inject the node centric graph structure for user and item embedding learning [42]. LightGCN [12] and LR-GCCF [5] use simple graph convolutions, which are easier to tune and show better performance.

Another research line of CF designed optimization and ranking goals. Earlier works modeled explicit rating values, and most of them directly use pointwise loss of root mean squared error [15, 29, 53]. As implicit feedback is more common in recommender systems, many researchers worked on implicit feedback based recommendation [22, 52]. In implicit feedback, there is a limited number of observed behaviors with a rating value of 1, and the remaining are unobserved behaviors. Some approaches extended pointwise approaches, and treated all unobserved items as negative items by assigning them with smaller confidence values [49].

As recommendation is a ranking problem that provides top ranked items, researchers proposed ranking based optimization approaches for recommendation. These ranking approaches can be divided into pairwise models [27, 31] and listwise models [2, 14]. Bayesian Personalized Ranking (BPR) is a popular pairwise approach that is designed for implicit feedback based recommendation [31]. Given a user, by selecting an observed-unobserved item pair, BPR assumes a user prefers an observed item compared to the unobserved item. BPR is easy to implement, and the time complexity is linear with the observed ratings. However, pairwise approaches treat each item pair independently, and ignore the correlation of multiple observed items and multiple unobserved items. To explore the correlation of more data, GBPR [27] and ABPR [24] built the correlation among users. Specifically, users with the same interests of items are grouped together, then these two models defined user group preference to build a relationship among users. Another direction involves modeling the correlation among multiple items [13, 26]. For example, by computing the mean value of each item in a set of items to form a set preference, CoFiSet is proposed

to introduce the correlation between multiple items for enhancing ranking performance [26]. However, most of these methods only focus on a single aspect of correlation. How to better exploit the structure information of implicit feedback, especially with the limited observed behaviors is still under explored.

The listwise approaches formulate a listwise based ranking loss to measure the distance between predicted list and true list [47], such as ListNet [2], ListRank-MF [35] and ListCF [14]. Some listwise approaches were applied for memory based CF, and needed explicit ratings for nearest neighbor ranking [14, 41]. SQL-Rank is proposed to cast listwise collaborative ranking as maximum likelihood under a permutation model which applies probabilities to permutations of the item set [45]. As the observed items are all treated as rating 1, SetRank is proposed to use the permutation probability to encourage one observed item ranks in front of multiple unobserved items in each list [38]. SetRank provides a new research perspective for listwise learning for implicit feedback, and achieves state-of-the-art ranking results. Another direction is directly maximizing the ranking metrics [1, 21, 33, 50]. Since most of the ranking metrics are not differentiable, existing models approximate the ranking metrics and optimize a smoothed version, such as deriving the lower bound. CLiMF optimizes a lower bound of Mean Reciprocal Rank (MRR) and becomes one of the popular listwise models for implicit feedback [34]. Despite the success of existing listwise models, most of them simplified the ranking problem for top-1 recommendation with surrogate optimization goals [14, 35, 38, 41].

In implicit feedback, the data is composed of a small portion of observed behaviors and a large portion of unobserved interactions. Some algorithms sampled negative items from a large portion of unobserved interactions with equal probabilities for each user [9, 25, 31, 39]. Other researchers developed techniques to pick valuable negative samples [4, 11, 54, 55]. By combining the context like a set of users or additional variables, a non-uniform item sampler is proposed for negative sampling [30]. Instead of explicitly picking unobserved items as negative items, NCE-PLRec generates negative instances by sampling from the popularity distribution [43]. These negative sampling methods further improve collaborative ranking performance. We also borrow the idea of negative sampling, and we explore how to better mine the relationships of positive samples from the positive set for “hard” negative item mining.

3 PRELIMINARY

Now, we define a general recommendation problem. Let user set \mathbf{U} and item set \mathbf{V} denote two kinds of entities of a recommender system, with $|U|$ and $|V|$ are the size of \mathbf{U} and \mathbf{V} . Let $\mathbf{R} = [r_{uv}]_{|U| \times |V|} \in \{0, 1\}$ denote the user-item interaction matrix, with each element r_{uv} in \mathbf{R} represents the implicit preference of user u to item v . $r_{uv} = 1$ indicates the observed feedback (i.e., user u likes the item v), $r_{uv} = 0$ indicates the unobserved feedback (i.e., the interaction between u and v is unobserved). For each user $u \in \mathbf{U}$, we use R_u^+ to indicate observed feedback, by contrast, $R_u^- = \mathbf{R} \setminus R_u^+$ to indicate non-interactive feedback. We have $R_u^+ \cap R_u^- = \emptyset$ and $R_u^+ \cup R_u^- = \mathbf{V}$. Given the user interaction matrix \mathbf{R} , the goal of CF ranking is to identify user’s preferences and predict the preference score matrix $\hat{\mathbf{R}} = [\hat{r}_{uv}]_{|U| \times |V|}$ with ranking based optimization functions.

We focus on design a ranking approach to predict $\hat{\mathbf{R}}$. As embedding based models have been the default choice of recommendation architecture, we can use any embedding based CF model to learn users’ and items’ embeddings $\mathbf{E} = \{\mathbf{E}_U, \mathbf{E}_V\}$. Specifically, $\mathbf{E}_U = \{e_1, \dots, e_u, \dots, e_{|U|}\}$ is the user embeddings, and $\mathbf{E}_V = \{e_{|U|+1}, \dots, e_v, \dots, e_{|U|+|V|}\}$ is the item embeddings. After that, we can predict each user u ’s preference to item v as $\hat{r}_{uv} = e_u^T e_v$. Please note that as we focus on ranking based optimization design, we assume the architecture of embedding learning is available, such as matrix factorization models [31] or neural graph models [5].

3.1 Pairwise and Listwise Learning

In order to predict users’ preferences, we need to use information of implicit feedback. In common implicit feedback, the mainly self-supervised information is a kind of preference assumption between observed and unobserved feedback, such as observed feedback $>_u$ unobserved feedback where $>_u$ represents the preference structure of user u . Pairwise learning is one of the most widely techniques to define a pairwise preference assumption. For each user u , pairwise learning assumes and uses self-information that the connected observed item i is more relevant than sampling unobserved item j , which can be formulated as follows:

$$i >_u j, \quad \forall i \in R_u^+, \quad \forall j \in R_u^-. \quad (1)$$

To this end, the corresponding ratings can be formulated as follows:

$$\hat{r}_{ui} > \hat{r}_{uj}, \quad (2)$$

where $i \in R_u^+$ and $j \in R_u^-$. In this situation, each pairwise item preference data is modeled independently without considering the complex correlations of multiple items [38]. For example, a user u likes item i_1 and i_2 , and has an unobserved item set of $[j_1, j_2, j_3]$. By independently sampling two pairwise relations: $i_1 >_u j_1$ and $i_2 >_u j_2$, the ground truth of $i_2 >_u j_1$ is not considered due to the independence assumption. Thus, the self-supervised information of implicit feedback is not fully utilized.

Listwise learning directly looks at all items [45] or entire observed items [34] at one time, and assumes that the correct sorting is better than the incorrect sorting. Let π be a ranking permutation (or sorting) of a list of items for a user. With the predicted ratings of each ranked user-item pair as \hat{r}_{π_j} , the permutation probability can be calculated as: $P(\pi) = \prod_{j=1}^{|\pi|} \frac{\phi(\hat{r}_{\pi_j})}{\sum_{l=j}^{|\pi|} \phi(\hat{r}_{\pi_l})}$, where $\phi(\cdot)$

is a monotonically increasing function that transforms the predicted rating into probability. Then, the goal of the listwise ranking is to find the best permutation probability of all users given the ground-truth rankings from the implicit feedback. In fact, this permutation probability needs to be calculated every time when the list has changed its order. Due to the high time cost of permutation probability calculation, nearly all listwise models consider top-1 for simplicity. However, this simplification ignores most of the information hidden in the entire observed item set, which leads to insufficient usage of implicit feedback information.

It is important to understand that these two ways (pairwise learning and listwise learning) do not make full use of implicit feedback. Given the uniqueness of the implicit feedback based CF, the following question should be considered: since observed items

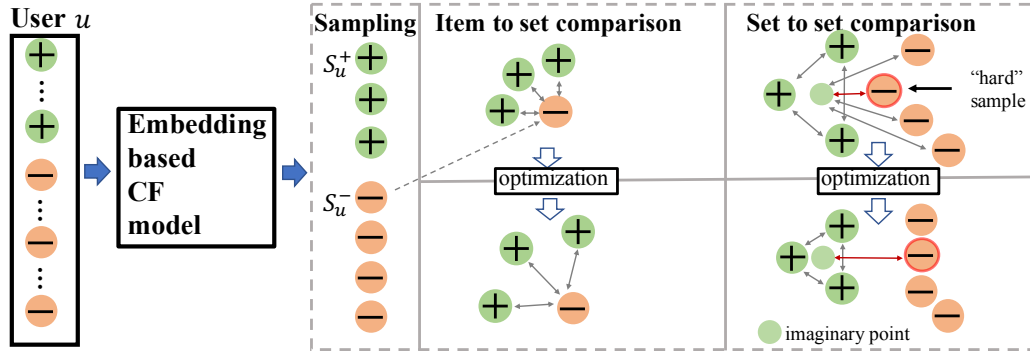


Figure 1: The overall structure of our proposed framework. After sampling observed set S_u^+ and unobserved set S_u^- , we build two-level comparisons that encourage the ranking positions of observed set before unobserved set. The item to set comparison pushes each unobserved sample far from observed set. The set to set comparison assumes that the summarized distance among observed set should be smaller than the distance between observed samples and the most “hard” unobserved sample.

and unobserved items are disordered inside each set, sampling entire items and sorting them is time consuming and unnecessary. On the other hand, considering only one positive item is insufficient for the utilization of limited observed items. Therefore, how to explore the self-supervised relationship for implicit feedback based ranking is the main challenge that we need to tackle.

4 THE PROPOSED SET2SET FRAMEWORK

In this section, we propose Set2set framework. In order to make full use of the implicit feedback, we propose to analyze the interactions between items and users at set level. Specifically, we first sample an observed set and an unobserved set. Then, we introduce the optimization and two-level comparison of our proposed Set2setRank framework. For easier understanding, Figure. 1 illustrates the overall architecture of our proposed Set2setRank framework with the two-level comparison. After that, we design an adaptive sampling technique to further improve the ability of Set2setRank framework. Finally, we make a discussion about the properties of our proposed Set2setRank framework.

4.1 Construction of Two Sets

As mentioned before, pairwise method suffers from independent assumption. Meanwhile, comparing all items in listwise methods leads to low efficiency problem. To this end, we try to make a compromise between these two methods. Specifically, we focus on an intermediate sampling to flexible use of observed and unobserved feedback. We intend to sample L observed items from the interactive item set R_u^+ , and utilize the corresponding predictive score to construct the observed set S_u^+ during training. Meanwhile, the same operation is applied to construct the unobserved set S_u^- . This process can be formulated as follows:

$$\begin{aligned} S_u^+ &= \{\hat{r}_{ui}^+ | i \in R_u^+\}, & |S_u^+| &= L, \\ S_u^- &= \{\hat{r}_{uj}^-\} | j \in R_u^-\}, & |S_u^-| &= K, \end{aligned} \tag{3}$$

where i and j indicate the i^{th} observed sample and the j^{th} unobserved sample for user u . For user u , \hat{r}_{ui}^+ and \hat{r}_{uj}^- denote the predicted ratings for the sampled observed item i and sampled unobserved item j . K and L are the sample size of S_u^+ and S_u^- . Since the size of observed items is much smaller than unobserved items, comparing

them at set level can not only fully exploit the structure information hidden in each set, but also effectively analyze the noise unobserved information with the guidance of observed information.

4.2 Optimization Criteria of Set2setRank

Since we have sampled two sets of each user, we aim at maximizing the following preference:

$$S_u^+ >_u S_u^- \tag{4}$$

where $>_u$ denotes the preference structure of user u . In order to model the interactions of users and items based on the observed set and unobserved set in a comprehensive way, we design two-level comparisons (i.e. *item to set comparison* and *set to set comparison*) to achieve this constraint (Eq.(4)) concretely. In the following part, we first introduce the preference gain. Then, we give a detailed analysis about our newly designed two-level comparison.

4.2.1 Preference Gain Function and Comparison Function. Eq.(4) describes the preference of user u , in which the observed items should be rated higher than the unobserved items. One step further, we should rank the observed items in front of the unobserved items. However, Eq.(4) is not differentiable. We need a replacement to satisfy the requirements above and also be differentiable. To this end, we define a preference gain $D(x, y)$, where x and y are sampled from predictions of the observed set S_u^+ and predictions of unobserved set S_u^- , separately. The output should be large, indicating the predicted ratings of the two samples are far away. Since the pairwise loss function in BPR is an intuitive and widely used preference function in implicit feedback, we adopt the same preference metric as BPR [31], which is formulated as follows:

$$D(x, y) = \sigma(x - y), \tag{5}$$

where $\sigma(\cdot)$ is the activation function. $\{x, y\}$ is the input of this function, and is two scalar data like two samples from S_u^+ or S_u^- .

For set comparison, we intend to fully exploit the structure information hidden in the implicit feedback and make full use of observed set to model user preference. Therefore, we try to build multiple comparisons/views between observed set and unobserved set. More specifically, we need each sample in one set (e.g., S_u^+) to be far away from the entire another set (e.g., S_u^-). To this end, we define a comparison function for our proposed set comparison on

the set level as follows:

$$F(X, y_j) = \sum_{i=1}^{|X|} [D(x_i, y_j)], \quad (6)$$

where y_j denotes the j^{th} elements in the set Y (e.g., unobserved set S_u^-), and X represents another set (e.g., observed set S_u^+). $|X|$ is the number of elements in the set X . Next, we introduce the technical details of our proposed two-level comparison on the sets.

4.2.2 Item to Set Comparison. Different from the comparison of a pair of items at each time of pairwise learning models, the comparison between observed set and unobserved set has to deal with multiple samples. Therefore, how to effectively and fully use set data to ensure the personalized recommendation performance is our main focus. There is explicit set-level preference information. For one user, we have confidence that most unobserved samples have lower predicted ratings than ratings of the observed set. In other words, the predicted rating of each unobserved item is encouraged to be smaller than all the predicted ratings of observed items. This constraint can be implemented as follows:

$$L_2(u) = \mathbb{E}_{S_u^+} \mathbb{E}_{S_u^-} \left[\sum_{j=1}^K \ln F(S_u^+, \hat{r}_{uj}^-) \right], \quad (7)$$

where $\{\hat{r}_{uj}^- | j = 1, \dots, K\}$ is all the samples in S_u^- .

By utilizing this loss function, the recommendation method can push the distance between observed samples(set) and unobserved set(samples) larger. Please note that, this formulation differs from pairwise models, as pairwise models can only push two items far away at one time. In contrast, the item to set level comparison pushes each unobserved item far away from entire items in the observed set for optimization.

4.2.3 Set to Set Comparison. The item to set comparison has exploited the structure information hidden in implicit feedback by optimizing $L + 1$ items (one negative item and L observed items from the sampled observed set) at one optimization step. Since the size of observed items is much smaller than unobserved items in the implicit feedback, item to set comparison still is in short of fully utilizing the structure information in observed set. In view of this drawback, we propose a set to set comparison to fully explore observed set and exploit useful information from the noise unobserved set. We assume that the summarized distance among observed set should be smaller than the distance between observed samples and the most “hard” unobserved sample.

Specifically, for each element \hat{r}_{ui}^+ in the observed set, the summary of an observed set S_u^+ can be implemented by directly measuring the distance between observed samples:

$$f_{pos}(S_u^+, S_u^+) = 1/L \sum_{i=1}^L \ln F(S_u^+, \hat{r}_{ui}^+). \quad (8)$$

For each item from the unobserved set S_u^- , since we have already calculated the distance between the observed set and each unobserved item in Eq.(7), we pick the “hard” unobserved item as:

$$\begin{aligned} f_{neg}(S_u^+, S_u^-) &= \min([\ln F(S_u^+, \hat{r}_{u1}^-), \ln F(S_u^+, \hat{r}_{u2}^-), \dots, \ln F(S_u^+, \hat{r}_{uK}^-)]) \\ &= \min([\ln F(S_u^+, \hat{r}_{uj}^-) | j = 1, \dots, K]), \end{aligned} \quad (9)$$

where $S_u^- = \{\hat{r}_{u1}^-, \dots, \hat{r}_{uK}^-\}$, \min represents picking the most “hard” unobserved item that is closely to the observed set. It combines the idea of negative sampling and has the advantage of low computational complexity.

As the distance of observed items (Eq.(8)) is smaller than the distance of the most “hard” unobserved items, we enlarge the margin:

$$L_3(u) = \mathbb{E}_{S_u^+} \mathbb{E}_{S_u^-} \ln D(f_{neg}(S_u^+, S_u^-), \beta * f_{pos}(S_u^+, S_u^+)), \quad (10)$$

where β is a margin parameter to adjust the distance between sets.

Another intuitive idea is that, similar as selecting the “hard” unobserved sample to keep the set-level distance, we can also select the “easy” observed sample for set to set comparison. Then, we push the “easy” observed sample far from the “hard” unobserved sample. This process can be achieved as:

$$\begin{aligned} g_{pos}(S_u^+) &= \max([\ln F(S_u^+, \hat{r}_{ui}^+) | j = 1, 2, \dots, L]), \\ L_4(u) &= \mathbb{E}_{S_u^+} \mathbb{E}_{S_u^-} \ln D(f_{neg}(S_u^+, S_u^-), \beta * g_{pos}(S_u^+)). \end{aligned} \quad (11)$$

In practice, we find that $g_{pos}(S_u^+)$ performs a little worse than the performance as $f_{pos}(S_u^+, S_u^+)$. A possible reason is that, the “easy” observed positive item has already been calculated in the item to set level comparison. By selecting the “easy” observed positive item, we only increase the weight of a single observed item. In the following, we select $f_{pos}(S_u^+, S_u^+)$ to achieve the set-level constraint.

4.2.4 Objective Function. We integrate two set comparison losses to obtain two correlated views of the same user. The overall loss can be formulated as maximizing the following loss:

$$Loss = \sum_{u \in U} Loss(u) = \sum_{u \in U} [L_2(u) + \lambda * L_3(u)], \quad (12)$$

where λ is a weighting parameter. By taking the consideration of two-level comparison, $Loss(u)$ can force Set2setRank to make observed items closer to the user and unobserved items far away from the user, which is in favor of providing better results.

4.3 Adaptive Set2setRank

In the previous part, we require each observed set have L items. In fact, the number of observed items varies greatly among different users, with some users have very large number of interaction records, while others have only several ratings [6, 32]. Therefore, instead of a fixed size L of the sampled observed set, the size of the observed set can be flexibly designed. To further enhance the ability of our proposed framework, we design a simple and effective mask method, which can flexibly select the set size.

Word masking randomly picks and replaces some words, which is widely applied in representation learning in many natural language processing scenarios [7, 8]. Inspired by word masking, we enhance Set2setRank in a similar way. Specifically, in each observed set S_u^+ , we use a random mask to remove some observed samples and keep at least two observed samples. The mask $Mask = \{Mask_1, \dots, Mask_1, \dots, Mask_L\}$ randomly selects some observed samples from observed set in each update. In order to integrate the $Mask$

into our framework, we modify Eq.(6) with the *Mask* as follows:

$$Mask_l = \begin{cases} 1, & \text{selected} \\ 0, & \text{else} \end{cases} \text{ and } \sum_{l=1}^L Mask_l \geq 2, \quad (13)$$

$$\tilde{F}(X, y_j) = \sum_{i=1}^{|X|} [D(x_i, y_j) * Mask_i], \quad (14)$$

please note that, *Mask* is generated by a random probability.

We have to note that our proposed *Mask* is compositional, which can be flexibly formed with different combinations of set lengths during training. This is equivalent to using a variety of sampling methods to form the observed set, thus observed set is actually changeable and dynamic to suit different lengths. Therefore, using the mask method can adaptively deal with the different numbers of observed samples. In this way, the item to set comparison L_2 in Eq.(7) will be modified as follows:

$$L_2(u) = \mathbb{E}_{S_u^+, S_u^-} \left[\sum_{j=1}^K \ln \tilde{F}(S_u^+, \hat{r}_{u_j}^-) \right]. \quad (15)$$

Meanwhile, Eq.(8) and Eq.(9) will be changed into:

$$\tilde{f}_{pos}(S_u^+, S_u^+) = 1/L \sum_{i=1}^L \ln \tilde{F}(S_u^+, \hat{r}_{u_i}^+), \quad (16)$$

$$\tilde{f}_{neg}(S_u^+, S_u^-) = \min([\ln \tilde{F}(S_u^+, \hat{r}_{u_j}^-) | j = 1, 2, \dots, K]). \quad (17)$$

By applying mask technology into Set2setRank, our proposed framework can deal with different numbers of observed samples for each user. Along this line, it can flexibly accommodate different combinations of observed samples. Thus, it is able to cope with different interactions and improves its generalization ability.

4.4 Model Discussion

Connections with Previous Works. In order to make full use of implicit feedback, we propose to compare the observed interactions and unobserved interactions at set level. This is the main difference of Set2setRank compared with pairwise and listwise methods. By introducing the item to set level comparison, we break the independence assumption of each item pair in pairwise based methods, and can jointly compare multiple unobserved items at the same time. Besides, the set to set level comparison pushes away the distance of positive items, and distance of the ‘‘hard’’ negative items picked from the item to set level calculation. The set to set level comparison jointly considers observed items, unobserved items together by modeling self-supervised information of the limited observed behaviors. However, listwise models either considered top-1 recommendation or surrogate optimization function with efficiency consideration, and could not well exploit the limited observed behavior for better ranking.

In practice, if we simply set $\lambda = 0$ and $|S_u^+| = |S_u^-| = 1$ in Eq.(12), $Loss(u)$ degenerates to pairwise loss, which is a typical target of pairwise learning. Furthermore, if we enlarge $|S_u^+|$ and $|S_u^-|$ to maximum lengths, Set2setRank is modified to a special listwise method, where a list of all observed items and all unobserved items are considered for ranking.

Complexity Analysis. We perform complexity analysis in this part. First of all, we introduce some notations. Let M_{avg} denote

the average number of user’s observed items, M_{max} denote the maximum number of users’ observed items. As each user has very limited observed behaviors, e.g., the density of most CF data is less than 1%, the average number of unobserved items can be regarded as the item set size. Therefore, almost all methods use sampling techniques to pick some unobserved items, and sample K unobserved items at a time. We analyze our proposed framework with four typical ranking methods, including pairwise ranking model of BPR [31], listwise ranking models of CLiMF [34] and SetRank [28], and a negative sampling based approach of AoBPR [30]. As all of these models are designed for ranking based optimization, for fair comparison, we omit the time complexity of predicting users’ preferences, and only compare the time complexity of the loss optimization part.

- BPR [31]: BPR compares one observed item and one unobserved item at each time, and the total number of unobserved items is $|U|M_{avg}K$. Thus, its time complexity is $O(|U|M_{avg}K)$.
- CLiMF [34]: It compares the relevance scores between all the observed items to optimize a smoothed version of the MRR. As the number of observed items vary largely of different users, the time complexity depends on the users that have largest observed set size. Therefore, its time complexity can be calculated as $O(|U|M_{max}^2)$.
- SetRank [38]: SetRank compares each observed item and K unobserved items to form a list of $K + 1$ items with permutation probability calculation. For each user, SetRank needs to formulate M_{max} sets to ensure each observed item is considered. Its time complexity is $O(|U|M_{max}(K + 1))$.
- AoBPR [30]: It is an adaptive sampling model based on BPR. The time complexity of adaptive sampling is $O(D + |V|\log|V|)$, where D is the dimension of the embedding vector. Therefore, the overall time complexity is the time complexity of BPR plus the adaptive sampling, which equals to $O(|U|M_{avg}K(D + |V|\log|V|))$.

For Set2setRank, the time cost mainly depends on the complexity of item to set comparison (Eq.(7)) and set to set comparison (Eq. (10)). For each user, Set2setRank samples L observed items and K unobserved items at each time. Since each user’s observed items need to be considered, all observed rating records ($|U|M_{avg}$) need to be sampled $|U|M_{avg}/L$ times. For item to set comparison (Eq.(7)), it involves KL comparisons of L observed items and K unobserved items. Thus, its time complexity can be calculated as $O((|U|M_{avg}/L)KL) = O(|U|M_{avg}K)$. The KL comparison in Eq.(7) can be reused in Eq.(9). The additional time complexity in set to set comparison (Eq.(10)) only comes from f_{pos} (Eq.(8)). The f_{pos} involves L^2 comparisons, its time complexity can be calculated as $O((|U|M_{avg}/L)L^2) = O(|U|M_{avg}L)$. Therefore, the overall time complexity is $O(|U|M_{avg}(K + L))$.

For most recommender systems, $M_{max} \gg M_{avg}$ always exists (e.g., $M_{max} = 6479$ and $M_{avg} = 101$ on MovieLens-20M dataset). Thus, two listwise approaches (CLiMF and SetRank) have the higher time complexity. For BPR and Set2setRank, their time complexity are very similar. In practice, for our proposed Set2setRank, the values of K and L are very small, and we find $L = K - 1$ achieves relatively high ranking performance. Thus, we have $O(|U|M_{avg}(K + L)) = O(|U|M_{avg}K)$. The above time complexity of Set2setRank is very similar to the time complexity of BPR.

Table 1: Overall performance of HR@N and NDCG@N on three datasets.

Dataset	Models	N=10		N=20		N=30		N=40		N=50	
		HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG
AmazonBooks	BPR	0.01851	0.01710	0.02853	0.02169	0.03821	0.02564	0.04737	0.02911	0.05556	0.03205
	Pop-Sampling	0.01945	0.01774	0.03028	0.02268	0.04072	0.02693	0.05033	0.03056	0.05914	0.03373
	AoBPR	0.01960	0.01787	0.03061	0.02288	0.04157	0.02735	0.05123	0.03101	0.06019	0.03423
	CLiMF	0.01546	0.01431	0.02351	0.01790	0.03079	0.02086	0.03761	0.02343	0.04392	0.02568
	Deep-SetRank	0.01950	0.01827	0.02913	0.02259	0.03879	0.02650	0.04751	0.02980	0.05592	0.03282
	Set2setRank_BPR	0.02084	0.01900	0.03221	0.02417	0.04276	0.02846	0.05288	0.03228	0.06189	0.03550
	Set2setRank(A)_BPR	0.02169	0.01985	0.03343	0.02520	0.04446	0.02968	0.05441	0.03345	0.06359	0.03675
	LR-GCCF	0.02209	0.02040	0.03407	0.02583	0.04532	0.03039	0.05532	0.03416	0.06498	0.03761
	Set2setRank_GCN	0.02247	0.02079	0.03454	0.02629	0.04567	0.03089	0.05648	0.03490	0.06597	0.03830
Set2setRank(A)_GCN	0.02264	0.02086	0.03453	0.02626	0.04598	0.03092	0.05658	0.03492	0.06620	0.03838	
Yelp	BPR	0.2380	0.2504	0.2512	0.2464	0.2767	0.2522	0.3012	0.2593	0.3240	0.2668
	Pop-Sampling	0.2545	0.2707	0.2667	0.2654	0.2935	0.2714	0.3195	0.2791	0.3428	0.2871
	AoBPR	0.2381	0.2523	0.2522	0.2487	0.2784	0.2547	0.3027	0.2618	0.3254	0.2683
	CLiMF	0.1130	0.1268	0.1113	0.1174	0.1185	0.1163	0.1279	0.1178	0.1373	0.1202
	Deep-SetRank	0.2759	0.2857	0.3010	0.2880	0.3314	0.2966	0.3587	0.3056	0.3841	0.3144
	Set2setRank_BPR	0.2899	0.3047	0.3046	0.3005	0.3333	0.3077	0.3623	0.3168	0.3886	0.3257
	Set2setRank(A)_BPR	0.3067	0.3197	0.3258	0.3176	0.3567	0.3255	0.3869	0.3355	0.4144	0.3445
	LR-GCCF	0.2672	0.2764	0.2832	0.2752	0.3125	0.2837	0.3426	0.2929	0.3715	0.3042
	Set2setRank_GCN	0.2689	0.2801	0.2832	0.2771	0.3114	0.2847	0.3403	0.2942	0.3675	0.3036
Set2setRank(A)_GCN	0.2920	0.3011	0.3135	0.3015	0.3470	0.3111	0.3805	0.3222	0.4115	0.3329	
MovieLens-20M	BPR	0.03708	0.02929	0.05952	0.03786	0.07788	0.04406	0.09523	0.04942	0.11064	0.0539
	Pop-Sampling	0.03580	0.02806	0.05655	0.03603	0.07527	0.04229	0.09116	0.04721	0.1047	0.05118
	AoBPR	0.03788	0.02951	0.06024	0.03808	0.08006	0.04469	0.09762	0.05015	0.11348	0.05478
	CLiMF	0.01752	0.01329	0.03043	0.01828	0.04241	0.02235	0.05341	0.02576	0.06408	0.02887
	Deep-SetRank	0.03801	0.02945	0.06242	0.03880	0.08353	0.04589	0.10219	0.05164	0.11886	0.05649
	Set2setRank_BPR	0.03855	0.03054	0.06120	0.03918	0.08161	0.04604	0.09939	0.05153	0.11532	0.05618
	Set2setRank(A)_BPR	0.03920	0.03104	0.06289	0.04010	0.08393	0.04716	0.10228	0.05283	0.11924	0.05778
	LR-GCCF	0.03919	0.03118	0.06210	0.03991	0.08224	0.04663	0.10078	0.05236	0.11755	0.05725
	Set2setRank_GCN	0.03990	0.03180	0.06238	0.04036	0.08300	0.04726	0.10092	0.05279	0.11710	0.05748
Set2setRank(A)_GCN	0.04058	0.03258	0.06421	0.04162	0.08509	0.04864	0.10418	0.05454	0.12070	0.05937	

5 EXPERIMENTS

5.1 Experimental Setup

Datasets and Evaluation Metrics. We use three publicly available datasets: Amazon Books¹, MovieLens-20M² and Yelp³. We follow the setting of Amazon Books in previous works [5, 42], Amazon Books provides 3 million ratings from 52 thousand users on 91 thousand items. Yelp dataset is adopted from the 2018 edition of the Yelp challenge, and we use the 10-core setting. Then Yelp provides 1 million ratings from 45 thousand users on 45 thousand items. MovieLens-20M dataset provides 20 million ratings from 138 thousand users on 27 thousand movies. Since our framework focuses on the top-N recommendations, we employed two widely ranking metrics for model evaluation: HR and NDCG [37]. For each user, we select all unrated items as candidates, and mix them with the records in the validation and test sets to select the top-N results.

Baseline. As mentioned before, our proposed framework is model-agnostic and can be applied to CF-based models. Thus, we adopt two typical models as the base CF model, including the classical matrix factorization model (BPR) [31] and a state-of-the-art graph-based recommendation model (LR-GCCF) [5]. For the sake of convenience, we use Set2setRank_BPR and Set2setRank_GCN to distinguish different base rating prediction models of our proposed framework. Besides, we use Set2setRank(A)_BPR and Set2setRank(A)_GCN to denote the adaptive Set2setRank ranking approach that is proposed in Section 4.3. In order to better verify the effectiveness of our proposed framework, we select the following baselines:

- BPR [31]: BPR is a widely used pairwise method.

- LR-GCCF [5]: LR-GCCF is a state-of-the-art graph-based recommendation model, which is easier to tune and shows better performance than BPR. We use the pairwise ranking loss as BPR as the optimization goal of LR-GCCF.
- CLiMF [34]: CLiMF is a typical listwise method for implicit feedback, which optimizes the lower bound of the MRR.
- Deep-SetRank [38]: Deep-SetRank is a state-of-the-art listwise learning model that achieves better performance than most permutation based ranking models. This method maximizes top-1 permutation probability to guarantee that each user prefers an observed item to multiple unobserved items.
- Pop-Sampling: Instead of the randomly sampling negative items for BPR, Pop-Sampling samples unobserved items based on the popularity of each item in the training data [30].
- AoBPR [30]: AoBPR proposed an adaptive and context-dependent sampling distribution to build a non-uniform item sampler.

Parameter Settings. In this part, we list the parameter settings of Set2setRank. The number K of sampled unobserved set size influences all ranking based models, for fair comparison, we set the number of unobserved item to 5 for all models, including our proposed model. We set $L = 2$ in Set2setRank_BPR and Set2setRank_GCN. To evaluate the model performance more comprehensively, we select the top-N values in the range of $N = \{10, 20, 30, 40, 50\}$. The parameter β is $\beta = 0.5$ for Set2setRank and $\beta = 0.2$ for Set2setRank(A). The parameter λ is $\lambda = 1$. We test parameters using the MindSpore [23] and other tools. We implement the proposed framework using the MindSpore [23] and other open-source Python libraries.

5.2 Experimental Results and Analyses

5.2.1 Overall Performance Comparison. Tables 1 reports the overall experimental results on different datasets with different evaluation metrics. We can obtain the following observations:

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

²<https://grouplens.org/datasets/movielens/20m>.

³<https://www.yelp.com/dataset>.

Table 2: Performance of using different observed set size K and unobserved set size L on Amazon Books.

Models	N=10		N=20		N=30		N=40		N=50	
	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG
BPR	0.01851	0.01710	0.02853	0.02169	0.03821	0.02564	0.04737	0.02911	0.05556	0.03205
Set2setRank_BPR(L=2,K=5)	0.02084	0.01900	0.03221	0.02417	0.04276	0.02846	0.05288	0.03228	0.06189	0.03550
Set2setRank_BPR(L=3,K=5)	0.02144	0.0197	0.03341	0.02517	0.04448	0.02968	0.05459	0.03350	0.06378	0.03679
Set2setRank_BPR(L=4,K=5)	0.02095	0.01904	0.03256	0.02432	0.04366	0.02883	0.05378	0.03265	0.06277	0.03588
Set2setRank_BPR(L=5,K=5)	0.02094	0.01893	0.03187	0.02390	0.04209	0.02805	0.05151	0.03160	0.06029	0.03474
Set2setRank_BPR(L=2,K=10)	0.02061	0.0190	0.03224	0.02429	0.04311	0.02873	0.05291	0.03244	0.06182	0.03564
Set2setRank_BPR(L=2,K=20)	0.02174	0.01998	0.03413	0.02558	0.04566	0.03026	0.05607	0.03420	0.06585	0.03771
Set2setRank_BPR(L=2,K=50)	0.02160	0.01987	0.03404	0.02554	0.04567	0.03027	0.05636	0.03431	0.06627	0.03785
Set2setRank_BPR(L=2,K=100)	0.02166	0.01983	0.03394	0.02541	0.04549	0.03009	0.05595	0.03405	0.06550	0.03745

Table 3: Performance of using different observed set size K and unobserved set size L on MovieLens-20M.

Models	N=10		N=20		N=30		N=40		N=50	
	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG
BPR	0.2380	0.2504	0.2512	0.2464	0.2767	0.2522	0.3012	0.2593	0.3240	0.2668
Set2setRank_BPR(L=2,K=5)	0.3072	0.3234	0.3245	0.3195	0.3539	0.3266	0.3838	0.3360	0.4112	0.3453
Set2setRank_BPR(L=3,K=5)	0.3071	0.3204	0.3235	0.3169	0.3550	0.3248	0.3857	0.3344	0.4141	0.3439
Set2setRank_BPR(L=4,K=5)	0.3052	0.3067	0.3261	0.3077	0.3584	0.3168	0.3899	0.3270	0.4188	0.3369
Set2setRank_BPR(L=5,K=5)	0.2902	0.2844	0.3175	0.2914	0.3515	0.3022	0.3835	0.3131	0.4130	0.3235
Set2setRank_BPR(L=2,K=10)	0.3101	0.3262	0.3239	0.3205	0.3550	0.3281	0.3860	0.3378	0.4142	0.3473
Set2setRank_BPR(L=2,K=20)	0.3110	0.3269	0.3280	0.3230	0.3585	0.3305	0.3882	0.3398	0.4151	0.3490
Set2setRank_BPR(L=2,K=50)	0.3321	0.34509	0.3559	0.3451	0.3901	0.3548	0.4222	0.3656	0.4517	0.3761
Set2setRank_BPR(L=2,K=100)	0.3258	0.3360	0.3531	0.3391	0.3903	0.3508	0.4245	0.3629	0.4554	0.3742

Table 4: Performance under different set to set comparisons, with “easy” means using the “easy” observed sample (Eq.(11)).

Dataset	Models	N=10		N=20		N=30		N=40		N=50	
		HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG
Amazon Books	Set2setRank_BPR(L=3,K=5)easy	0.02119	0.01919	0.03334	0.02476	0.04458	0.02934	0.05476	0.03319	0.06385	0.03645
	Set2setRank_BPR(L=3,K=5)	0.02144	0.0197	0.03341	0.02517	0.04448	0.02968	0.05459	0.03350	0.06378	0.03679
	Set2setRank_BPR(L=4,K=5)easy	0.01946	0.01780	0.03072	0.02291	0.04062	0.02693	0.04972	0.03037	0.05811	0.03336
	Set2setRank_BPR(L=4,K=5)	0.02095	0.01904	0.03256	0.02432	0.04366	0.02883	0.05378	0.03265	0.06277	0.03588
MovieLens-20M	Set2setRank_BPR(L=3,K=5)easy	0.3075	0.3196	0.3265	0.3177	0.3591	0.3263	0.3903	0.3362	0.4183	0.3456
	Set2setRank_BPR(L=3,K=5)	0.3071	0.3204	0.3235	0.3169	0.3550	0.3248	0.3857	0.3344	0.4141	0.3439
	Set2setRank_BPR(L=4,K=5)easy	0.2994	0.2976	0.3267	0.3028	0.3615	0.3134	0.3933	0.3241	0.4220	0.3342
	Set2setRank_BPR(L=4,K=5)	0.3052	0.3067	0.3261	0.3077	0.3584	0.3168	0.3899	0.3270	0.4188	0.3369

Table 5: Performance of the two modules of Set2setRank on AmazonBooks: item to set comparison and set to set comparison.

Models	Item to set comparison	Set to set comparison	N=10		N=20		N=30		N=40		N=50	
			HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG
BPR	/	/	0.01851	0.01710	0.02853	0.02169	0.03821	0.02564	0.04737	0.02911	0.05556	0.03205
Set2setRank_BPR	✓	✗	0.01999	0.01835	0.03128	0.02346	0.04150	0.02762	0.05110	0.03124	0.05973	0.03435
Set2setRank_BPR	✗	✓	0.02019	0.01872	0.02978	0.02309	0.03946	0.02704	0.04835	0.03042	0.05651	0.03335
Set2setRank_BPR	✓	✓	0.02084	0.01900	0.03221	0.02417	0.04276	0.02846	0.05288	0.03228	0.06189	0.03550
LR-GCCF	/	/	0.02209	0.02040	0.03407	0.02583	0.04532	0.03039	0.05532	0.03416	0.06498	0.03761
Set2setRank_GCN	✓	✗	0.02237	0.02052	0.03459	0.02605	0.04658	0.03092	0.05739	0.03503	0.06770	0.03871
Set2setRank_GCN	✗	✓	0.02250	0.02076	0.03465	0.02628	0.04603	0.03093	0.05628	0.03480	0.06592	0.03824
Set2setRank_GCN	✓	✓	0.02247	0.02079	0.03454	0.02629	0.04567	0.03089	0.05648	0.03490	0.06597	0.03830

- Our proposed Set2setRank achieves the best performance across all datasets with different evaluation metrics. Specifically, Set2setRank_BPR(A) outperforms the best baselines by an average relative boost of 9%, 11% and 4% on Amazon Books MovieLens-20M, and Yelp, respectively. Set2setRank_GCN(A) also outperforms all the baselines. The adaptive setting usually performs better than the counterpart without adaptive sampling. These observations demonstrate that our proposed two-level set comparison can not only fully exploit implicit feedback to model user preference, but also provide accurate recommendation results. Moreover, based on Set2setRank, the improvement of BPR is better than LR-GCCF. The reason is that LR-GCCF uses more characteristics of implicit feedback than BPR, which can further prove the

effectiveness of our proposed Set2setRank framework on implicit feedback utilization.

- Set2setRank_GCN (Set2setRank_BPR) also achieves promising performance on all datasets. However, the corresponding improvements are not so obvious than Set2setRank_GCN (Set2setRank_BPR), demonstrating that our newly designed adaptive sampling technique can further improve the model performance. Moreover, we obtain that the improvement is not much. We speculate one possible reason is that the length of unobserved samples set and observed samples set in this experiment is small, which limits the ability of our proposed adaptive sampling technique.
- Among all baselines, CLiMF does not achieve satisfactory performance on all datasets. This is because CLiMF focuses

Table 6: Performance of the two modules of Set2setRank on MovieLens-20M: item to set comparison and set to set comparison.

Models	Item to set comparison	Set to set comparison	N=10		N=20		N=30		N=40		N=50	
			HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG	HR	NDCG
BPR	/	/	0.2380	0.2504	0.2512	0.2464	0.2767	0.2522	0.3012	0.2593	0.3240	0.2668
Set2setRank_BPR	✓	✗	0.2899	0.3047	0.3046	0.3005	0.3333	0.3077	0.3623	0.3168	0.3886	0.3257
Set2setRank_BPR	✗	✓	0.2869	0.3057	0.2984	0.2979	0.3246	0.3033	0.3510	0.3109	0.3761	0.3190
Set2setRank_BPR	✓	✓	0.3072	0.3234	0.3245	0.3195	0.3539	0.3266	0.3838	0.3360	0.4112	0.3453
LR-GCCF	/	/	0.2672	0.2764	0.2832	0.2752	0.3125	0.2837	0.3426	0.2929	0.3715	0.3042
Set2setRank_GCN	✓	✗	0.2671	0.2777	0.2832	0.2760	0.3125	0.2843	0.3422	0.2942	0.3700	0.3040
Set2setRank_GCN	✗	✓	0.2650	0.2793	0.2773	0.2746	0.3042	0.2811	0.3318	0.2896	0.3583	0.2984
Set2setRank_GCN	✓	✓	0.2689	0.2801	0.2832	0.2771	0.3114	0.2847	0.3403	0.2942	0.3675	0.3036

on the first relevant item of recommendation lists. This situation also appears in [18, 50]. Meanwhile, Deep-SetRank outperforms the pairwise learning and listwise learning method, this is also consistent with the results in Deep-SetRank [38].

5.2.2 Influence of Set Size. Since we propose to compare the observed samples and unobserved samples at set level, the sizes of these two sets play an important role. In order to explore the influence of different negative set size K and positive set size L , we make additional experiments on Amazon Books and MovieLens-20M datasets to verify the model performance. Moreover, we select Set2setRank_BPR to avoid the influence of our proposed adapted sampling technique. The results are reported in Table 2 and 3. From these two tables, we can observe that with the increase of K and L , the results first increase and then decrease. The best setting is $L = 3, K = 50$ for Amazon Book dataset, and $L = 3, K = 50$ for MovieLens-20M dataset. On the one hand, with the increase of K and L , the model can compare more observed samples and unobserved samples, which is very helpful for hidden structure information exploration and implicit feedback utilization. On the other hand, since the size of unobserved items is much larger than observed items, continually increasing K and L will lead to the inconsistency of sampled observed set and sampled unobserved set, which will result model performance drop. This observation inspires us to carefully determine the value of K and L for the best performance of our proposed framework.

5.2.3 Set to Set Comparison Strategy. In Section 4.2.3, we have mentioned that the set to set comparison in our proposed two-level comparison has two strategies: 1) in Eq.(8), considering the relationships of all observed pairs as the summary of the observed set. 2) in Eq.(11), selecting the “easy” observed sample with the consideration of the observed set. In this part, we compare the performance of these two set to set comparison strategies. We employ Set2setRank_BPR as the base model and make additional experiments to verify the effectiveness of each strategy. Table 4 summaries the corresponding results on Amazon Books and MovieLens-20M datasets. From the results, we can conclude that both strategies achieve relatively high performance on two datasets with different evaluation metrics, in which summarizing all observed pairs achieves better performance. We speculate a possible reason is that selecting the “hard” unobserved sample already fully exploits the observed set and is also an effective way to improve performance, as shown by the negative sampling model [17, 30]. If we further select “easy” positive item, this positive item information has already been calculated from item to set level comparison, and set to set comparison degenerates to increase weights of “easy” positive item. In contrast, the summary information of the observed set is not well

exploited. Therefore, considering the relationships of all observed pairs in Eq.(8) shows the best performance. To this end, we select the strategy in Eq.(8) to implement our proposed framework.

5.2.4 Ablation Study. In this part, we investigate the effectiveness of each component: item to set comparison (Eq.(7)) and set to set comparison (Eq.(10)) of our proposed Set2set framework. The results are illustrated in Tables 5 and 6. From the two tables, we can obtain the following observations. First, each single component (item to set comparison or set to set comparison) can still help model achieve comparable performance, indicating the usefulness of our proposed comparison at set level. Moreover, with single component, the model achieves similar performance, indicating that both of them are very important for user preference modeling and model performance improvement. Besides, compared with the performance of models with single component, models with both of them (entire Set2setRank framework) have better performance on all datasets, demonstrating that necessity of both components in our proposed two-level comparison. As these two components consider different self-supervised information for implicit feedback, combining them together reaches the best performance.

6 CONCLUSION

In this paper, we argued that current ranking models for CF are still far from satisfactory due to the sparsity of the observed user-item interactions for each user. To this end, we proposed to compare the observed items and unobserved items at set level to fully explore the structure information hidden in the implicit feedback. Specifically, we first constructed the sampled observed set and sampled unobserved set. Then, we proposed a Set2setRank framework for the utilization of implicit feedback. Set2setRank mainly consisted of two parts: item to set comparison to encourage unobserved samples are ranked lower than the observed set, and set to set comparison to fully explore the limited observed set to be far away from the hard negative sample from the unobserved set. Moreover, we designed an adaptive technique to further improve the performance of Set2setRank framework. Extensive experiments on three real-world datasets clearly showed the advantages of our proposed framework.

ACKNOWLEDGMENTS

This work was supported in part by grants from the National Natural Science Foundation of China (Grant No. 61972125, U1936219, 61932009, 61732008, 62006066), the Young Elite Scientists Sponsorship Program by CAST and ISZS, CAAI-Huawei MindSpore Open Fund, and the Fundamental Research Funds for the Central Universities, HFUT.

REFERENCES

- [1] Suhrid Balakrishnan and Sumit Chopra. 2012. Collaborative ranking. In *WSDM*. 143–152.
- [2] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *ICML*. 129–136.
- [3] Dong-Kyu Chae, Jihoo Kim, Duen Horng Chau, and Sang-Wook Kim. 2020. AR-CF: Augmenting Virtual Users and Items in Collaborative Filtering for Addressing Cold-Start Problems. In *SIGIR*. 1251–1260.
- [4] Chong Chen, Min Zhang, Weizhi Ma, Yiqun Liu, and Shaoping Ma. 2020. Jointly non-sampling learning for knowledge graph enhanced recommendation. In *SIGIR*. 189–198.
- [5] Lei Chen, Le Wu, Richang Hong, Kun Zhang, and Meng Wang. 2020. Revisiting Graph Based Collaborative Filtering: A Linear Residual Graph Convolutional Network Approach. In *AAAI*, Vol. 34. 27–34.
- [6] Zhihong Chen, Rong Xiao, Chenliang Li, Gangfeng Ye, Haochuan Sun, and Hongbo Deng. 2020. Esam: Discriminative domain adaptation with non-displayed items to improve long-tail performance. In *SIGIR*. 579–588.
- [7] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Ziqing Yang, Shijin Wang, and Guoping Hu. 2020. Pre-training with whole word masking for chinese bert. In *EMNLP*. 657–668.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*. 4171–4186.
- [9] Jingtao Ding, Fuli Feng, Xiangnan He, Guanghui Yu, Yong Li, and Depeng Jin. 2018. An improved sampler for bayesian personalized ranking by leveraging view data. In *WWW*. 13–14.
- [10] Zeno Gantner, Lucas Drumond, Christoph Freudenthaler, and Lars Schmidt-Thieme. 2012. Personalized ranking for non-uniformly sampled items. In *SIGKDD*. 231–247.
- [11] Mihajlo Grbovic and Haibin Cheng. 2018. Real-time personalization using embeddings for search ranking at airbnb. In *SIGKDD*. 311–320.
- [12] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *SIGIR*. 639–648.
- [13] Cheng-Kang Hsieh, Longqi Yang, Yin Cui, Tsung-Yi Lin, Serge Belongie, and Deborah Estrin. 2017. Collaborative metric learning. In *WWW*. 193–201.
- [14] Shanshan Huang, Shuaiqiang Wang, Tie-Yan Liu, Jun Ma, Zhumin Chen, and Jari Veijalainen. 2015. Listwise collaborative filtering. In *SIGIR*. 343–352.
- [15] Mohsen Jamali and Martin Ester. 2010. A matrix factorization technique with trust propagation for recommendation in social networks. In *RecSys*. 135–142.
- [16] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* 42, 8 (2009), 30–37.
- [17] Defu Lian, Qi Liu, and Enhong Chen. 2020. Personalized Ranking with Importance Sampling. In *WWW*. 1093–1103.
- [18] Junjie Liang, Jinlong Hu, Shoubin Dong, and Vasant Honavar. 2018. Top-N-Rank: A Scalable List-wise Ranking Method for Recommender Systems. In *Big Data*. 1052–1058.
- [19] Nathan N Liu, Evan W Xiang, Min Zhao, and Qiang Yang. 2010. Unifying explicit and implicit feedback for collaborative filtering. In *ICKM*. 1445–1448.
- [20] Qi Liu, Zhenya Huang, Yu Yin, Enhong Chen, Hui Xiong, Yu Su, and Guoping Hu. 2019. Ekt: Exercise-aware knowledge tracing for student performance prediction. *TKDE* 33, 1 (2019), 100–115.
- [21] Yong Liu, Peilin Zhao, Aixin Sun, and Chunyan Miao. 2015. A boosting algorithm for item recommendation with implicit feedback. In *IJCAI*. 1792–1798.
- [22] Hongyu Lu, Min Zhang, Weizhi Ma, Ce Wang, Feng xia, Yiqun Liu, Leyu Lin, and Shaoping Ma. 2019. Effects of User Negative Experience in Mobile News Streaming. In *SIGIR*. 705–714.
- [23] Mindspore. 2020. . <https://www.mindspore.cn/>
- [24] Shan Ouyang, Lin Li, WeiKe Pan, and Zhong Ming. 2019. Asymmetric Bayesian personalized ranking for one-class collaborative filtering. In *RecSys*. 373–377.
- [25] Rong Pan, Yunhong Zhou, Bin Cao, Nathan N Liu, Rajan Lukose, Martin Scholz, and Qiang Yang. 2008. One-class collaborative filtering. In *ICDM*. 502–511.
- [26] WeiKe Pan and Li Chen. 2013. Cofiset: Collaborative filtering via learning pairwise preferences over item-sets. In *SIAM*. 180–188.
- [27] WeiKe Pan and Li Chen. 2013. Gbpr: Group preference based bayesian personalized ranking for one-class collaborative filtering. In *IJCAI*. 2691–2697.
- [28] Liang Pang, Jun Xu, Qingyao Ai, Yanyan Lan, Xueqi Cheng, and Jirong Wen. 2020. Setrank: Learning a permutation-invariant ranking model for information retrieval. In *SIGIR*. 499–508.
- [29] Lianyong Qi, Ruili Wang, Chunhua Hu, Shancang Li, Qiang He, and Xiaolong Xu. 2019. Time-aware distributed service recommendation with privacy-preservation. *Information Sciences* 480 (2019), 354–364.
- [30] Steffen Rendle and Christoph Freudenthaler. 2014. Improving pairwise learning for item recommendation from implicit feedback. In *WSDM*. 273–282.
- [31] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *UAI*. 452–461.
- [32] Shaoyun Shi, Weizhi Ma, Min Zhang, Yongfeng Zhang, Xinxing Yu, Houzhi Shan, Yiqun Liu, and Shaoping Ma. 2020. Beyond User Embedding Matrix: Learning to Hash for Modeling Large-Scale Users in Recommendation. In *SIGIR*. 319–328.
- [33] Yue Shi, Alexandros Karatzoglou, Linas Baltrunas, Martha Larson, Alan Hanjalic, and Nuria Oliver. 2012. Tfmap: optimizing map for top-n context-aware recommendation. In *SIGIR*. 155–164.
- [34] Yue Shi, Alexandros Karatzoglou, Linas Baltrunas, Martha Larson, Nuria Oliver, and Alan Hanjalic. 2012. CLIMF: learning to maximize reciprocal rank with collaborative less-is-more filtering. In *RecSys*. 139–146.
- [35] Yue Shi, Martha Larson, and Alan Hanjalic. 2010. List-wise learning to rank with matrix factorization for collaborative filtering. In *RecSys*. 269–272.
- [36] Harald Steck. 2011. Item popularity and recommendation accuracy. In *RecSys*. 125–132.
- [37] Peijie Sun, Le Wu, and Meng Wang. 2018. Attentive recurrent social recommendation. In *SIGIR*. 185–194.
- [38] Chao Wang, Hengshu Zhu, Chen Zhu, Chuan Qin, and Hui Xiong. 2020. SetRank: A Setwise Bayesian Approach for Collaborative Ranking from Implicit Feedback. In *AAAI*. 6127–6136.
- [39] Huazheng Wang, Sonwoo Kim, Eric McCord-Snook, Qingyun Wu, and Hongning Wang. 2019. Variance reduction in gradient exploration for online learning to rank. In *SIGIR*. 835–844.
- [40] Pengfei Wang, Hanxiong Chen, Yadong Zhu, Huawei Shen, and Yongfeng Zhang. 2019. Unified collaborative filtering over graph embeddings. In *SIGIR*. 155–164.
- [41] Shuaiqiang Wang, Shanshan Huang, Tie-Yan Liu, Jun Ma, Zhumin Chen, and Jari Veijalainen. 2016. Ranking-oriented collaborative filtering: A listwise approach. *TOIS* 35, 2 (2016), 1–28.
- [42] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *SIGIR*. 165–174.
- [43] Ga Wu, Maksims Volkovs, Chee Loong Soon, Scott Sanner, and Himanshu Rai. 2019. Noise Contrastive Estimation for One-Class Collaborative Filtering. In *SIGIR*. 135–144.
- [44] Le Wu, Xiangnan He, Xiang Wang, Kun Zhang, and Meng Wang. 2021. A Survey on Neural Recommendation: From Collaborative Filtering to Content and Context Enriched Recommendation. *arXiv preprint arXiv:2104.13030* (2021).
- [45] Liwei Wu, Cho-Jui Hsieh, and James Sharpnack. 2018. SQL-Rank: A Listwise Approach to Collaborative Ranking. In *ICML*. 5315–5324.
- [46] Le Wu, Yonghui Yang, Kun Zhang, Richang Hong, Yanjie Fu, and Meng Wang. 2020. Joint item recommendation and attribute inference: An adaptive graph convolutional network approach. In *SIGIR*. 679–688.
- [47] Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. Listwise approach to learning to rank: theory and algorithm. In *ICML*. 1192–1199.
- [48] Hong-Jian Xue, Xinyu Dai, Jianbing Zhang, Shujian Huang, and Jiajun Chen. 2017. Deep Matrix Factorization Models for Recommender Systems. In *IJCAI*. 3203–3209.
- [49] Hsiang-Fu Yu, Mikhail Bilenko, and Chih-Jen Lin. 2017. Selection of negative samples for one-class matrix factorization. In *SIAM*. 363–371.
- [50] Runlong Yu, Qi Liu, Yuyang Ye, Mingyue Cheng, Enhong Chen, and Jianhui Ma. 2020. Collaborative List-and-Pairwise Filtering from Implicit Feedback. *TKDE* (2020).
- [51] Runlong Yu, Yunzhou Zhang, Yuyang Ye, Le Wu, Chao Wang, Qi Liu, and Enhong Chen. 2018. Multiple pairwise ranking with implicit feedback. In *CIKM*. 1727–1730.
- [52] Wenhui Yu and Zheng Qin. 2020. Sampler Design for Implicit Feedback Data by Noisy-label Robust Learning. In *SIGIR*. 861–870.
- [53] Hafed Zarzour, Ziad Al-Sharif, Mahmoud Al-Ayyoub, and Yaser Jararweh. 2018. A new collaborative filtering recommendation algorithm based on dimensionality reduction and clustering techniques. In *ICICS*. 102–106.
- [54] Weinan Zhang, Tianqi Chen, Jun Wang, and Yong Yu. 2013. Optimizing top-n collaborative filtering via dynamic negative item sampling. In *SIGIR*. 785–788.
- [55] Hao Zhong, WeiKe Pan, Congfu Xu, Zhi Yin, and Zhong Ming. 2014. Adaptive pairwise preference learning for collaborative recommendation with implicit feedbacks. In *CIKM*. 1999–2002.