

DE-RRD: A Knowledge Distillation Framework for Recommender System

SeongKu Kang, Junyoung Hwang, Wonbin Kweon, Hwanjo Yu*

Dept. of Computer Science and Engineering, POSTECH, South Korea
{seongku, jyhwang, kwk4453, hwanjoyu}@postech.ac.kr

ABSTRACT

Recent recommender systems have started to employ knowledge distillation, which is a model compression technique distilling knowledge from a cumbersome model (teacher) to a compact model (student), to reduce inference latency while maintaining performance. The state-of-the-art methods have only focused on making the student model accurately imitate the predictions of the teacher model. They have a limitation in that the prediction results incompletely reveal the teacher’s knowledge. In this paper, we propose a novel knowledge distillation framework for recommender system, called *DE-RRD*, which enables the student model to learn from the latent knowledge encoded in the teacher model as well as from the teacher’s predictions. Concretely, *DE-RRD* consists of two methods: 1) *Distillation Experts (DE)* that directly transfers the latent knowledge from the teacher model. DE exploits “experts” and a novel expert selection strategy for effectively distilling the vast teacher’s knowledge to the student with limited capacity. 2) *Relaxed Ranking Distillation (RRD)* that transfers the knowledge revealed from the teacher’s prediction with consideration of the relaxed ranking orders among items. Our extensive experiments show that *DE-RRD* outperforms the state-of-the-art competitors and achieves comparable or even better performance to that of the teacher model with faster inference time.

CCS CONCEPTS

• **Information systems** → **Learning to rank**; *Collaborative filtering*; Retrieval efficiency.

KEYWORDS

Recommender System; Knowledge Distillation; Learning to Rank; Model Compression; Retrieval efficiency

ACM Reference Format:

SeongKu Kang, Junyoung Hwang, Wonbin Kweon, Hwanjo Yu. 2020. DE-RRD: A Knowledge Distillation Framework for Recommender System. In *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM '20)*, October 19–23, 2020, Virtual Event, Ireland. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3340531.3412005>

*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 ACM 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](https://permissions.acm.org).

CIKM '20, October 19–23, 2020, Virtual Event, Ireland

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-6859-9/20/10...\$15.00

<https://doi.org/10.1145/3340531.3412005>

1 INTRODUCTION

In recent years, recommender system (RS) has been broadly adopted in various industries, helping users’ decisions in the era of information explosion, and playing a key role in promoting corporate profits. However, a growing scale of users (and items) and sophisticated model architecture to capture complex patterns make the size of the model continuously increasing [13, 25, 28, 30]. A large model with numerous parameters has a high capacity, and thus usually has better recommendation performance. On the other hand, it requires a large computational time and memory costs, and thus incurs a high latency during the inference phase, which makes it difficult to apply such large model to real-time platform.

Motivated by the significant success of *knowledge distillation* (KD) in the computer vision field, a few work [13, 25] have employed KD for RS to reduce the size of models while maintaining the performance. KD is a model-agnostic strategy to accelerate the learning of a new compact model (student) by transferring knowledge from a previously trained large model (teacher) [7]. The knowledge transfer is conducted as follows: First, the teacher model is trained with the user-item interactions in the training set which has binary labels – ‘1’ for observed interactions, and ‘0’ for unobserved interactions. Then, the student model is trained with the “soft” labels generated by the teacher model (i.e., teacher’s predictions) along with the available binary labels. The student model trained with KD has comparable performance to that of the teacher, and also has a lower latency due to its small size [13, 25].

The core idea behind this process is that the soft labels predicted by the teacher model reveal hidden relations among entities (i.e., users and items) not explicitly included in the training set, so that they accelerate and improve the learning of the student model. Specifically, the items ranked near the top of a user’s recommendation list would have strong correlations to the items that the user has interacted before [25]. Also, the soft labels provide guidance for distinguishing the items that each user would like and the items that each user would not be interested in among numerous unobserved items only labeled as ‘0’ [13]. By using the additional supervisions from the teacher model, the state-of-the-art methods [13, 25] have achieved comparable or even better performance to the teacher models with faster inference time.

However, there are still limitations in existing methods [13, 25]. First, the learning of the student is only guided by the teacher’s prediction results, which is not sufficient to fully take advantage of the knowledge stored in the teacher. This is because the prediction results incompletely reveal the teacher’s knowledge. As illustrated in Figure 1, the recommendation list from the teacher only shows that a user has a similar degree of preference on the two items (0.99 and 0.98). However, *latent knowledge* in the teacher, which is used to make such predictions, contains more detailed information that

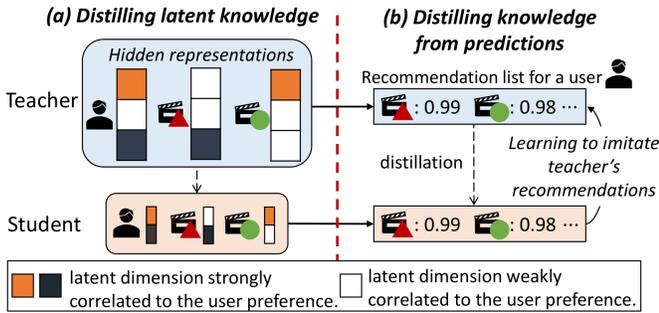


Figure 1: The existing methods [13, 25] distill the knowledge only based on the teacher’s predictions (b). The proposed framework directly distills the latent knowledge stored in the teacher (a) along with the knowledge revealed from the predictions (b).

the user likes different aspects of the two items (marked as navy blue and orange, respectively). In this regard, we argue that the training process and the performance of the student can be further improved by directly distilling such latent knowledge stored in the teacher model. Second, they distill the knowledge from the teacher’s predictions in a point-wise manner that considers a single item at a time. Because the point-wise approach does not consider multiple items simultaneously, it has a limitation in accurately maintaining the ranking orders predicted by the teacher model [24], which leads to degraded recommendation performance.

In this paper, we propose a novel knowledge distillation framework for RS, named DE-RRD, which distills both the latent knowledge stored in the teacher model (Fig. 1a) and the knowledge revealed from teacher’s predictions (Fig. 1b). By learning both the teacher’s final predictions and the detailed knowledge that provides the bases for such predictions, the student model can be further improved. The proposed framework consists of two methods: 1) Distillation Experts (DE) and 2) Relaxed Ranking Distillation (RRD). The main contributions of this paper lie in the following aspects:

Distilling latent knowledge in the teacher model. We propose a novel method—DE—for directly distilling latent knowledge stored in the teacher model. Specifically, DE transfers the knowledge from hidden representation space (i.e., the output of the intermediate layer) of the teacher to the representation space of the student. Due to the limited capacity, the student model cannot learn all the knowledge in the teacher representation space. DE first introduces an “expert”, which is a small feed-forward network, to distill the summarized knowledge that can restore the detailed knowledge of each entity in the teacher. However, distilling the knowledge of all entities with a single expert intermingles the information of weakly correlated entities and further hinders the entire distillation process. To tackle this problem, DE adopts the multiple experts and a novel expert selection strategy that clearly distinguishes the knowledge that each expert distills based on the correlations among the entities in the teacher representation space. To the best of our knowledge, our approach is the first attempt to directly distill the latent knowledge in the teacher model for RS. We demonstrate its rationality and superiority through extensive experiments and comprehensive analyses.

Relaxed Ranking Distillation from the teacher’s predictions.

We propose a new method—RRD—that transfers the knowledge from the teacher’s predictions with direct consideration of ranking orders among items. Unlike the existing methods [13, 25] that distill the knowledge of an item at a time, RRD formulates this as a ranking matching problem between the recommendation list of the teacher and that of the student. To this end, RRD adopts the list-wise learning-to-rank approach [29] and learns to ensure the student to preserve the ranking orders predicted by the teacher. However, directly applying the list-wise approach can have adverse effects on the recommendation performance. Since a user is interested in only a few items among the numerous total items [10], learning the detailed ranking orders of all items is not only daunting but also ineffective. To tackle this challenge, RRD reformulates the daunting task to a *relaxed ranking matching* problem. Concretely, RRD matches the recommendation list from the teacher and that from the student, *ignoring* the detailed ranking orders among the uninteresting items that the user would not be interested in. RRD achieves superior recommendation performance compared to the state-of-the-art methods [13, 25].

An unified framework. We propose a novel framework—DE-RRD—which enables the student model to learn both from the teacher’s predictions and from the latent knowledge stored in the teacher model. Our extensive experiments on real-world datasets show that DE-RRD considerably outperforms the state-of-the-art competitors. DE-RRD achieves comparable performance to that of the teacher with a smaller number of learning parameters than all the competitors. Also, DE-RRD shows the largest performance gain when the student has the identical structure to the teacher model (i.e., self-distillation [5]). Furthermore, we provide both qualitative and quantitative analyses to further investigate the superiority of each proposed component. The source code of DE-RRD is publicly available¹.

2 RELATED WORK

Balancing *effectiveness* and *efficiency* is a key requirement for real-time recommender system (RS); the system should provide *accurate recommendations* with *fast inference time*. Recently, the size of the recommender model is continuously increasing, and the computational time and memory cost required for the inference are also increasing accordingly [13, 25, 28, 30]. Due to the high latency, it becomes difficult to apply such large recommender to the real-time large-scale platform. In this section, we review several approaches to alleviate this problem.

Balancing Effectiveness and Efficiency. Several methods have adopted hash techniques to reduce the inference cost [10, 15, 16, 30]. They first learn binary representations of users and items, then construct the hash table. Although exploiting the binary representation can significantly reduce the inference costs, due to the constrained capability, their recommendation performance is limited compared to models that use real-values representations. In addition, several work has focused on accelerating the inference of the existing recommenders [1, 14, 26]. Specifically, tree-based data structures [2], data compression techniques [26], and approximated nearest neighbor search techniques [4, 23] have been successfully adopted to

¹https://github.com/SeongKu-Kang/DE-RRD_CIKM20

reduce the search costs. However, they still have problems such as applicable only to specific models (e.g., k-d tree for metric learning-based models [12]), or easily falling into a local optimum due to the local search.

Knowledge Distillation. Knowledge distillation (KD) is a model-agnostic strategy to improve the learning and the performance of a new “compact” model (student) by transferring knowledge from a previously trained “large” model (teacher) [3, 5, 7, 22]. The student model trained with KD has comparable performance to that of the teacher model, and also has lower inference latency due to its small size. Most KD methods have focused on the image classification problem. An early work [7] matches the softmax distribution of the teacher and the student. The predicted label distribution contains more rich information (e.g., inter-class correlation) than the one-hot class label, which leads to improved learning of the student model. Subsequent methods [3, 22] have focused on distilling knowledge from intermediate layers. Because teacher’s intermediate layers are generally bigger than that of the student, they [3, 22] utilize additional layers to bridge the different dimensions. Interestingly, KD has turned out to be effective in improving the teacher model itself by self-distillation [5].

Knowledge Distillation in Recommender System. Recently, inspired by the huge success of KD in the computer vision field, a few work [13, 25] have adopted KD to RS. A pioneer work is Ranking Distillation (RD) [25] which applies KD for the ranking problem; Providing recommendations of top- N unobserved items that have not interacted with a user. RD jointly optimizes a base recommender’s loss function with a *distillation loss*.

$$\min_{\theta_s} \mathcal{L}_{Base} + \lambda \mathcal{L}_{RD} \quad (1)$$

where θ_s is the learning parameters of the student model, λ is a hyperparameter that controls the effects of RD. The base recommender can be any existing RS model such as BPR [21], NeuMF [6], and \mathcal{L}_{Base} is its loss function (e.g., binary cross-entropy). The distillation loss of RD for user u is defined as follows:

$$\mathcal{L}_{RD} = - \sum_{\pi_k \in \pi} w_{\pi_k} \log(P(\text{rel} = 1|u, \pi_k)) \quad (2)$$

where π is a ranked list of top- K unobserved items for user u predicted by the teacher, π_k is the k -th item in this ranking, and $P(\text{rel} = 1|u, \pi_k)$ is the relevance probability of user u to π_k predicted by the student model. w_{π_k} is the weight, which is computed based on each item’s ranking from the student and the teacher, for reflecting relative importance among top- K items.

A subsequent work Collaborative Distillation (CD) [13] first samples unobserved items from the teacher’s recommendation list according to their ranking; high-ranked items are more frequently sampled, then trains the student to mimic the teacher’s prediction score (e.g., relevance probability) on the sampled items. The distillation loss of CD for user u is defined as follows:

$$\mathcal{L}_{CD} = - \left(\sum_{\pi_k \in \pi} q_{\pi_k} \log(P(\text{rel} = 1|u, \pi_k)) + (1 - q_{\pi_k}) \log(1 - P(\text{rel} = 1|u, \pi_k)) \right) \quad (3)$$

where π is a ranked list of K unobserved items sampled from teacher’s recommendations for user u , q_{π_k} is the weight, which

is computed based on teacher’s prediction score on each item, for reflecting relative importance among the sampled items.

In summary, the distillation loss of the existing methods makes the student model follow the teacher’s predictions on unobserved items with particular emphasis on the high-ranked items. In RS, only high-ranked items in the recommendation list are matter. Also, such high-ranked items reveal hidden patterns among entities (i.e., users and items); the high-ranked items in the recommendation list would have strong correlations to the user [25]. By using such additional supervisions from the teacher, they have achieved the comparable performance to the teacher with faster inference time.

However, the existing methods still have room for improvement by the following reasons: First, the student can be further improved by directly distilling the *latent knowledge* stored in the teacher model. Latent knowledge refers to all information of users, items, and relationships among them that is discovered and stored in the teacher model. Such knowledge is valuable for the student because it provides detailed explanations on the final prediction of the teacher. Second, they transfer the knowledge from the teacher’s predictions with a point-wise approach that considers a single item at a time. Since the point-wise approach does not take into account multiple items simultaneously, it has a limitation in accurately maintaining the ranking orders in the teacher’s ranking list [24]. This can lead to limited recommendation performance.

3 PROBLEM FORMULATION

In this work, we focus on top- N recommendations for implicit feedback. Let \mathcal{U} and \mathcal{I} denote the set of users and items, respectively. Given collaborative filtering (CF) information (i.e., implicit interactions between users and items), we build a binary matrix $R \in \{0, 1\}^{|\mathcal{U}| \times |\mathcal{I}|}$. Each element of R has a binary value indicating whether a user has interacted with an item (1) or not (0). Note that an unobserved interaction does not necessarily mean a user’s negative preference on an item, it can be that the user is not aware of the item. For each user, a recommender model ranks all items that have not interacted with the user (i.e., unobserved items) and provides a ranked list of top- N unobserved items.

The knowledge distillation is conducted as follows: First, a teacher model with a large number of learning parameters is trained with the training set which has binary labels. Then, a student model with a smaller number of learning parameters is trained with the help from the teacher model in addition to the binary labels. The goal of KD is to fully improve the inference efficiency without compromising the effectiveness; We aim to design a KD framework that enables the student model to maintain the recommendation performance of the teacher with a small number of learning parameters.

4 DE-RRD: THE PROPOSED FRAMEWORK

We propose DE-RRD framework which enables the student model to learn both from the teacher’s predictions and from the latent knowledge encoded in the teacher model. DE-RRD consists of two methods: 1) *Distillation Experts* (DE) that directly transfers the latent knowledge from the teacher, 2) *Relaxed Ranking Distillation* (RRD) that transfers the knowledge revealed from the teacher’s predictions with direct consideration of ranking orders among items. This section is organized as follows. We first describe each

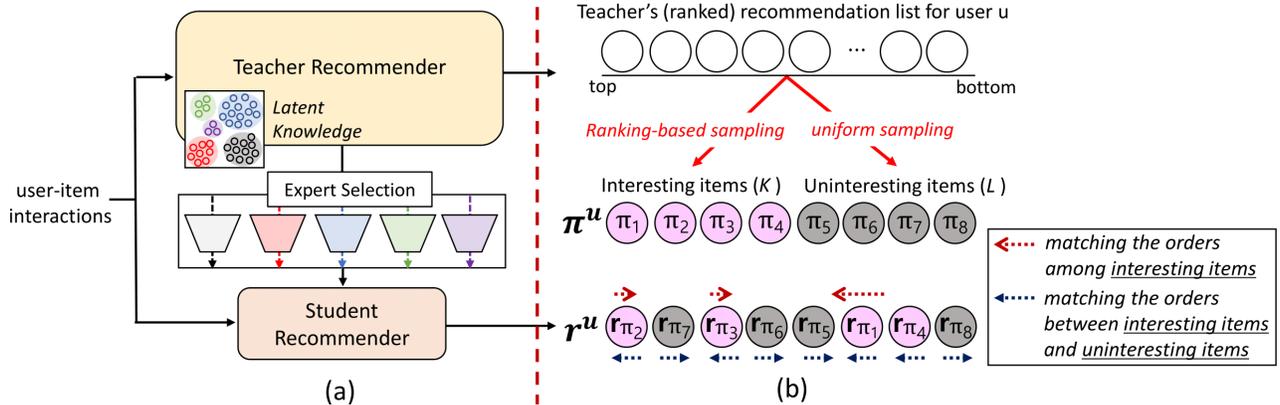


Figure 2: Illustration of DE-RRD framework. (a) Distillation Experts (DE) directly distills the teacher’s latent knowledge with the experts and the selection strategy. (b) Relaxed Ranking Distillation (RRD) distills the knowledge from the teacher’s prediction based on the relaxed ranking approach that ignores orders among the uninteresting items. Best viewed in color.

component of the proposed framework: DE in Section 4.1, RRD in Section 4.2. Then, we explain the end-to-end optimization process in Section 4.3. The overview of DE-RRD is provided in Figure 2.

4.1 Distillation Experts (DE)

In this section, we provide the details of DE which distills the latent knowledge from the hidden representation space (i.e., the output of the intermediate layer) of the teacher to the corresponding representation space of the student. We first introduce “expert” to distill the summarized knowledge that can restore the detailed teacher’s knowledge of each entity. Then, we introduce a novel expert selection strategy for effectively distilling CF knowledge that contains information of all the entities having diverse preferences and characteristics.

4.1.1 Expert for distillation. DE exploits “expert” to distill knowledge from the teacher’s hidden representation space. An expert, which is a small feed-forward network, is trained to *reconstruct* the representation on a selected intermediate layer of the teacher from the representation on the corresponding intermediate layer of the student. Let $h_t(\cdot)$ denote a mapping function to the representation space ($\in \mathbb{R}^{d_t}$) of the teacher model (i.e., a nested function up to the intermediate layer of the teacher). Similarly, let $h_s(\cdot)$ denote a mapping function to the student’s representation space ($\in \mathbb{R}^{d_s}$). The output of the mapping function can be a separate representation of a user, an item (e.g., BPR [21]) or their combined representation (e.g., NeuMF [6]) based on the base model’s structure and the type of selected layer. Here, we use user u as an example for convenience. An expert E is trained to reconstruct $h_t(u)$ from $h_s(u)$ as follows:

$$\mathcal{L}(u) = \|h_t(u) - E(h_s(u))\|_2 \quad (4)$$

Note that in the KD process, the teacher model is already trained and frozen. By minimizing the above equation, parameters in the student model (i.e., $h_s(\cdot)$) and the expert are updated.

The student model has smaller capacity compared to the teacher ($d_s \ll d_t$). By minimizing the equation 4, the student learns compressed information on the user’s preference that can restore more detailed knowledge in the teacher as accurate as possible. This approach provides a kind of filtering effect and improves the learning of the student model.

4.1.2 Expert selection strategy. Training a single expert to distill all the CF knowledge in the teacher is not sufficient to achieve satisfactory performance. The CF knowledge contains vast information of user groups with various preferences and item groups with diverse characteristics. When a single expert is trained to distill the knowledge of all the diverse entities, the information of the weakly correlated entities (e.g., users that have dissimilar preferences) is mixed and reflected in the expert’s weights. This leads to the adulterated distillation that hinders the student model from discovering some users’ preferences.

To alleviate the problem, DE puts multiple experts in parallel and clearly distinguishes the knowledge that each expert distills. The key idea is to divide the representation space into exclusive divisions based on the teacher’s knowledge and make each expert to be specialized in distilling the knowledge in a division (Fig. 2a). The representations belonging to the same division has strong correlations with each other, and they are distilled by the same expert without being mixed with weakly correlated representations belonging to the different divisions. The knowledge transfer of DE is conducted in the two steps: 1) a selection network first computes each expert’s degree of specialization for the knowledge to be distilled. 2) DE selects an expert based on the computed distribution, then distills the knowledge through the selected expert.

Concretely, DE has M experts (E_1, E_2, \dots, E_M) and a selection network S whose output is M -dimensional vector. To distill user u ’s knowledge from the teacher, the selection network S first computes the normalized specialization score vector $\alpha^u \in \mathbb{R}^M$ as follows:

$$\begin{aligned} e^u &= S(h_t(u)), \\ \alpha_m^u &= \frac{\exp(e_m^u)}{\sum_{i=1}^M \exp(e_i^u)} \quad \text{for } m = 1, \dots, M \end{aligned} \quad (5)$$

Then, DE selects an expert based on the computed distribution. We represent the selection variable s^u that determines which expert to be selected for distilling $h_t(u)$. s^u is a M -dimensional one-hot vector where an element is set to 1 if the corresponding expert is selected for distillation. DE samples this selection variable s^u from a multinoulli distribution parameterized by $\{\alpha_m^u\}$ i.e.,

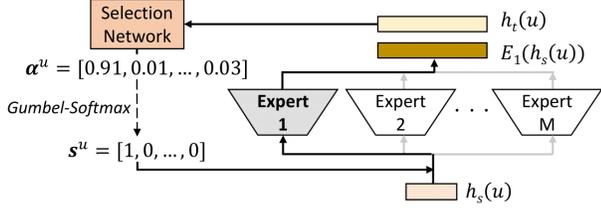


Figure 3: Illustration of the expert selection process of DE. During the training, s^u becomes a one-hot vector and selects the most specialized expert in the knowledge to be distilled.

$p(s_m^u = 1 | S, h_t(u)) = \alpha_m^u$, then reconstructs teacher’s representation as follows:

$$s^u \sim \text{Multinoulli}_M(\{\alpha_m^u\})$$

$$\mathcal{L}(u) = \|h_t(u) - \sum_{m=1}^M s_m^u \cdot E_m(h_s(u))\|_2 \quad (6)$$

However, the sampling process is non-differentiable, which would block the gradient flows and disable the end-to-end training. As a workaround, we adopt a continuous relaxation of the discrete distribution by using Gumbel-Softmax [8]. The Gumbel-Softmax is a continuous distribution on the simplex that can approximate samples from a categorical distribution; it uses the Gumbel-Max trick [18] to draw samples from the categorical distribution, then uses the softmax function as a continuous approximation of argmax operation to get the approximated one-hot representation. With the relaxation, the selection network can be trained by the back-propagation.

DE gets the approximated one-hot selection variable s^u by using the Gumbel-Softmax and reconstructs the teacher’s representation as follows:

$$s_m^u = \frac{\exp((\log \alpha_m^u + g_m) / \tau)}{\sum_{i=1}^M \exp((\log \alpha_i^u + g_i) / \tau)} \quad \text{for } m = 1, \dots, M \quad (7)$$

$$\mathcal{L}(u) = \|h_t(u) - \sum_{m=1}^M s_m^u \cdot E_m(h_s(u))\|_2$$

where g_i is i.i.d drawn from Gumbel(0, 1) distribution². The extent of relaxation is controlled by a temperature parameter τ . In the beginning of the training, we set a large value on τ , and gradually decreases its value during the training. As τ is decreased to 0, s^u smoothly becomes one-hot vector where $s_m^u = 1$ with probability α_m^u . In other words, during the training, each expert gradually gets specialized on certain information that has strong correlations. This process is illustrated in Figure 3.

Discussion: Effects of expert selection. As the expert selection is based on the teacher’s knowledge, correlations among the entities in the teacher representation space are naturally reflected in the expert selection; the user representations with very similar preferences (i.e., located closely in the space) would be distilled by the same expert with a high probability. This allows each expert to be trained to distill only the knowledge of strongly correlated entities, and thus each expert can provide better guidance that does not include the information of weakly correlated entities.

² $g_i = -\log(-\log(r))$, where r is sampled from $Uniform(0, 1)$.

Discussion: selection vs. attention. Instead of selecting one expert, the attention mechanism (i.e., the softmax function) can be adopted. However, we think the selection is a more appropriate choice to distill the CF knowledge containing all the entities having diverse preferences and characteristics. This is because the attention makes every expert involved in distilling the knowledge of each entity. In other words, like in the case of a single expert, all the experts and attention network are trained to minimize the overall reconstruction errors of all the diverse entities. By doing so, information of weakly relevant entities gets mixed together, and this leads to performance degrade in some user groups. We provide experiment results to support our claims. Please refer to Section 5.3.

4.1.3 Optimization of DE. DE is jointly optimized with the base model’s loss function in the end-to-end manner as follows:

$$\min_{\theta_s, \theta_{DE}} \mathcal{L}_{Base} + \lambda_{DE} \cdot \mathcal{L}_{DE} \quad (8)$$

where θ_s is the learning parameters of the student model, θ_{DE} is the learning parameters of DE (i.e., the selection network and the experts), and λ_{DE} is a hyperparameter that controls the effects of DE. The base model can be any existing recommender (e.g., BPR, NeuMF), and \mathcal{L}_{Base} corresponds to its loss function. Note that the experts are not used in the inference phase.

The loss function of DE can be flexibly defined based on the base model’s structure and the types of hidden layer chosen for the distillation. Concretely, for NeuMF [6], which is a state-of-the-art deep recommender, the loss function can be defined to 1) separately distill knowledge of users and items in a mini-batch (i.e., $\sum_{u \in B} \mathcal{L}(u) + \sum_{i \in B} \mathcal{L}(i)$) or 2) distill the combined knowledge (i.e., $\sum_{(u,i) \in B} \mathcal{L}(u, i)$). Also, we adopt a simple temperature annealing schedule, which gradually decays the temperature from τ_0 to τ_P as done in [11]: $\tau(p) = \tau_0 (\tau_P / \tau_0)^{p/P}$ where $\tau(p)$ is the temperature at epoch p , and P is the total training epochs.

4.2 Relaxed Ranking Distillation (RRD)

We propose RRD, a new method to distill the knowledge revealed from the teacher’s predictions with direct consideration of ranking orders among items. RRD formulates this as a ranking matching problem between the recommendation list of the teacher model and that of the student model. To this end, RRD adopts the classical list-wise learning-to-rank approach [29]. Its core idea is to define a probability of a permutation (i.e., a ranking order) based on the ranking score predicted by a model, and train the model to maximize the likelihood of the ground-truth ranking order. For more details about the list-wise approach, please refer to [29].

However, merely adopting the list-wise loss can have adverse effects on the ranking performance. Because a user is interested in only a few items among the numerous total items [10], learning the detailed ranking orders of all the unobserved items is not only daunting but also ineffective. The recommendation list from the teacher model contains information about a user’s potential preference on each unobserved item; A few items that the user would be interested in (i.e., interesting items) are located near the top of the list, whereas the majority of items that the user would not be interested in (i.e., uninteresting items) are located far from the top.

Based on this information, RRD reformulates the daunting task of learning all the precise ranking orders to a *relaxed ranking matching*

problem. In other words, RRD aims to match the recommendation list from the teacher and that from the student, *ignoring* the detailed ranking orders among the uninteresting items. Concretely, RRD distills the information of 1) the detailed ranking orders among the interesting items, 2) the relative ranking orders between the interesting items and the uninteresting items. The overview of RRD is provided in Figure 2b.

4.2.1 Sampling interesting/uninteresting items. The first step of RRD is to sample items from the teacher’s recommendation list. In specific, RRD samples K interesting items and L uninteresting items for each user. As a user would not be interested in the vast majority of items, the interesting items should be sampled from a very narrow range near the top of the list, whereas the uninteresting items should be sampled from the wide range of the rest. To sample the interesting items, we adopt a ranking position importance scheme [20, 25] that places more emphasis on the higher positions in the ranking list. In the scheme, the probability of the k -th ranked item to be sampled is defined as: $p_k \propto e^{-k/T}$ where T is the hyperparameter that controls emphasis on top positions. With the scheme, RRD samples K interesting items according to the user’s potential preference on each item (i.e., item’s ranking) predicted by the teacher. To sample the uninteresting items that corresponds the majority of items, we use a simple uniform sampling. Concretely, RRD uniformly samples L uninteresting items from a set of items that have lower rankings than the previously sampled interesting items.

4.2.2 Relaxed permutation probability. Then, RRD defines a relaxed permutation probability motivated by [29]. For user u , π^u denotes a ranked list of all the sampled items ($K + L$) sorted by the original order in the teacher’s recommendation list. r^u denotes ranking scores on the sampled items predicted by the student model. The relaxed permutation probability is formulated as follows:

$$p(\pi_{1:K}^u | r^u) = \prod_{k=1}^K \frac{\exp(r_{\pi_k}^u)}{\sum_{i=k}^K \exp(r_{\pi_i}^u) + \sum_{j=K+L}^{K+L} \exp(r_{\pi_j}^u)} \quad (9)$$

where $r_{\pi_k}^u$ denotes a ranking score predicted by the student for the k -th item in π^u , $\pi_{1:K}^u$ denotes the partial list that contains the interesting items. RRD learns to maximize the log-likelihood $\log p(\pi_{1:K} | r)$ for all users. The proposed permutation probability is not affected by the detailed ranking orders among the uninteresting items (L). By maximizing the log-likelihood, the student model is trained to locate all the interesting items (K) higher than all the uninteresting items (L) in the recommendation list, while maintaining the detailed ranking orders (from the teacher’s recommendation list) among the interesting items.

4.2.3 Optimization of RRD. RRD is jointly optimized with the base model’s loss function in the end-to-end manner as follows:

$$\min_{\theta_s} \mathcal{L}_{Base} + \lambda_{RRD} \cdot \mathcal{L}_{RRD} \quad (10)$$

where θ_s is the learning parameters of the student model and λ_{RRD} is a hyperparameter that controls the effects of RRD. The base model can be any existing recommender, and \mathcal{L}_{Base} corresponds to its loss function. The sampling process is conducted at every epoch. The loss function of RRD is defined to distill the knowledge of users in the mini-batch: $-\frac{1}{|B|} \sum_{u \in B} \log p(\pi_{1:K}^u | r^u)$.

4.3 Optimization of DE-RRD

The proposed DE-RRD framework is optimized in the end-to-end manner as follows:

$$\min_{\theta_s, \theta_{DE}} \mathcal{L}_{Base} + \lambda_{DE} \cdot \mathcal{L}_{DE} + \lambda_{RRD} \cdot \mathcal{L}_{RRD} \quad (11)$$

where θ_s is the learning parameters of the student model, θ_{DE} is the learning parameters of DE (i.e., the selection network and the experts). The base model can be any existing recommender, and \mathcal{L}_{Base} corresponds to its loss function.

5 EXPERIMENTS

We validate the superiority of DE-RRD on 12 experiment settings (2 real-world datasets \times 2 base models \times 3 different student model sizes). We first provide extensive experiment results supporting that DE-RRD outperforms the state-of-the-art competitors (Section 5.2). We also provide both quantitative and qualitative analyses to verify the rationality and superiority of each proposed component (Section 5.3). Lastly, we provide hyperparameter study (Section 5.4).

5.1 Experimental Setup

Datasets. We use two public real-world datasets: CiteULike [27], Foursquare [17]. We remove users and items having fewer than five ratings for CiteULike, twenty ratings for Foursquare as done in [6, 9, 21]. Data statistics are summarized in Table 1.

Table 1: Data Statistics (after preprocessing)

Dataset	#Users	#Items	#Interactions	Sparsity
CiteULike	5,220	25,182	115,142	99.91%
Foursquare	19,466	28,594	609,655	99.89%

Base Models. We validate the proposed framework on base models that have different architectures and optimization strategies. We choose a latent factor model and a deep learning model that are broadly used for top- N recommendation with implicit feedback.

- **BPR [21]:** A learning-to-rank model for implicit feedback. It assumes that observed items are more preferred than unobserved items and optimizes Matrix Factorization (MF) with the pair-wise ranking loss function.
- **NeuMF [6]:** The state-of-the-art deep model for implicit feedback. NeuMF combines MF and Multi-Layer Perceptron (MLP) to learn the user-item interaction, and optimizes it with the point-wise objective function (i.e., binary cross-entropy).

Teacher/Student. For each base model and dataset, we increase the number of learning parameters until the recommendation performance is no longer increased, and use the model with the best performance as Teacher model. For each base model, we build three student models by limiting the number of learning parameters. We adjust the number of parameters based on the size of the last hidden layer. The limiting ratios (ϕ) are {0.1, 0.5, 1.0}. Following the notation of the previous work [13, 25], we call the student model trained without the help of the teacher model (i.e., no distillation) as “Student” in this experiment sections.

Comparison Methods. The proposed framework is compared with the following methods:

- **Ranking Distillation (RD)** [25]: A KD method for recommender system that uses items with the highest ranking from the teacher’s predictions for distilling the knowledge.
- **Collaborative Distillation (CD)** [13]: The state-of-the-art KD method for recommender system. CD samples items from teacher’s predictions based on their ranking, then uses them for distillation. As suggested in the paper, we use unobserved items only for distilling the knowledge.

Finally, **DE-RRD** framework consists of the following two methods:

- **Distillation Experts (DE)**: A KD method that directly distills the latent knowledge stored in the teacher model. It can be combined with any *prediction-based* KD methods (e.g., RD, CD, RRD).
- **Relaxed Ranking Distillation (RRD)**: A KD method that distills the knowledge revealed from the teacher’s predictions with consideration of relaxed ranking orders among items.

Evaluation Protocol. We follow the widely used *leave-one-out* evaluation protocol [6, 9, 19]. For each user, we leave out a single interacted item for testing, and use the rest for training. In our experiments, we leave out an additional interacted item for the validation. To address the time-consuming issue of ranking all the items, we randomly sample 499 items from a set of unobserved items of the user, then evaluate how well each method can rank the test item higher than these sampled unobserved items. We repeat this process of sampling a test/validation item and unobserved items five times and report the average results.

As we focus on the top- N recommendation task based on implicit feedback, we evaluate the performance of each method with widely used three ranking metrics [6, 9, 10]: hit ratio ($H@N$), normalized discounted cumulative gain ($N@N$), and mean reciprocal rank ($M@N$). $H@N$ measures whether the test item is present in the top- N list, while $N@N$ and $M@N$ are position-aware ranking metrics that assign higher scores to the hits at upper ranks.

Implementation Details for Reproducibility. We use PyTorch to implement the proposed framework and all the baselines, and use Adam optimizer to train all the methods. For RD, we use the public implementation provided by the authors. For each dataset, hyperparameters are tuned by using grid searches on the validation set. The learning rate for the Adam optimizer is chosen from $\{0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001\}$, the model regularizer is chosen from $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$. We set the total number of epochs as 1000, and adopt early stopping strategy; stopping if $H@5$ on the validation set does not increase for 30 successive epochs. For all base models (i.e., BPR, NeuMF), the number of negative sample is set to 1, and no pre-trained technique is used. For NeuMF, the number of the hidden layers is chosen from $\{1, 2, 3, 4\}$.

For all the distillation methods (i.e., RD, CD, DE, RRD), weight for KD loss (λ) is chosen from $\{1, 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$. For DE, the number of experts (M) is chosen from $\{5, 10, 20, 30\}$, MLP is employed for the experts and the selection network. The shape of the layers of an expert is $[d_s \rightarrow (d_s + d_t)/2 \rightarrow d_t]$ with *relu* activation, and that of the selection network is $[d_t \rightarrow M]$. We select the last hidden layer of all the base models to distill latent knowledge. We put the experts according to the structure of the selected layer; For the layer where user and item are separately encoded (i.e., BPR), we put M user-side experts and M item-side experts, and for the layer where user and items are jointly encoded

(i.e., NeuMF), we put M experts to distill the combined information. τ_0 and τ_p are set to $1, 10^{-10}$, respectively. For prediction-based KD methods (i.e., RD, CD, RRD), the number of high-ranked (or interesting) items (K) for distillation is chosen from $\{10, 20, 30, 40, 50\}$, weight for controlling the importance of top position (T) is chosen from $\{1, 5, 10, 20\}$. For RRD, the number of uninteresting items (L) is set to the same with K , but it can be further tuned. For RD, the number of the warm-up epoch is chosen from $\{30, 50, 100\}$, the number of negative items in the dynamic weight is chosen from $\{50, 100\}$. Also, RD and CD have additional hyperparameters for reflecting the relative importance of the items used for distillation. We follow the recommended values from the public implementation and from the original papers.

5.2 Performance Comparison

Table 2 shows top- N recommendation accuracy of different methods in terms of various ranking metrics. In summary, DE-RRD shows the significant improvement compared to the state-of-the-art KD methods on two base models that have different architectures and optimization strategies. Also, DE-RRD consistently outperforms the existing methods on three different sizes of the student model in Figure 4. We analyze the results from various perspectives.

We first observe that the two methods of the proposed framework (i.e., DE, RRD) improve the performance of the student model. DE directly distills the teacher’s latent knowledge that includes detailed information on users, items, and the relationships among them. This enables the student to be more effectively trained than finding such information from scratch with a limited capacity. RRD distills the knowledge from the teacher’s predictions based on the relaxed ranking approach which makes the student to effectively maintain the ranking orders of interesting items predicted by the teacher. Unlike the existing methods (i.e., RD, CD), it directly handles the ranking violations among the sampled items, which can lead to better ranking performance.

Also, we observe that RRD achieves large performance gain particularly in NeuMF ($\phi = 0.1$). One possible reason is that NeuMF is trained with the point-wise loss function (i.e., binary cross-entropy) which considers only one item at a time. In general, it is known that the approaches considering the preference orders between items (e.g., pair-wise, list-wise) can achieve better ranking performance than the point-wise approach [24]. RRD enables the model to capture the ranking orders among the unobserved items, so that it can lead to the large performance gain. Interestingly, we observe that the prediction-based KD methods (i.e., RD, CD, RRD) can have an adverse effect when the model size is large (NeuMF with $\phi = 0.5, 1.0$ in Figure 4). We conjecture that this is because when a model has sufficient capacity to achieve comparable performance to the teacher, enforcing it to exactly mimic the teacher’s prediction results can act as a strong constraint that rather hinders its learning.

In addition, we observe that DE-RRD achieves the best performance among all the methods in general. DE-RRD enables the student to learn both from the teacher’s prediction and from the latent knowledge that provides the bases for such predictions. Interestingly, DE-RRD also shows a large performance gain when the student model has the identical structure to the teacher model (i.e., self-distillation with $\phi = 1.0$ in Figure 4). This result shows that

Table 2: Recommendation performances ($\phi = 0.1$). *Improv.b* and *Improv.s* denote the improvement of DE-RRD over the best baseline and student respectively. *, **, *, and **** indicate $p \leq 0.05$, $p \leq 0.005$, $p \leq 0.0005$, and $p \leq 0.00005$ for the paired t-test of vs. the best baseline (for RRD, DE-RRD), vs. Student (for DE) on H@5.**

Dataset	Base Model	KD Method	H@5	M@5	N@5	H@10	M@10	N@10	H@20	M@20	N@20
CiteULike	BPR	Teacher	0.5135	0.3583	0.3970	0.6185	0.3724	0.4310	0.7099	0.3788	0.4541
		Student	0.4441	0.2949	0.3319	0.5541	0.3102	0.3691	0.6557	0.3133	0.3906
		RD	0.4533	0.3019	0.3395	0.5601	0.3161	0.3740	0.6633	0.3232	0.3993
		CD	0.4550	0.3025	0.3404	0.5607	0.3167	0.3746	0.6650	0.3240	0.4011
		DE **	0.4817	0.3230	0.3625	0.5916	0.3372	0.3977	0.6917	0.3441	0.4229
	RRD **	0.4622	0.3076	0.3461	0.5703	0.3220	0.3809	0.6746	0.3293	0.4074	
	DE-RRD ***	0.4843	0.3231	0.3632	0.5966	0.3373	0.3989	0.6991	0.3447	0.4251	
	<i>Improv.b</i>	6.44%	6.81%	6.7%	6.4%	6.47%	6.47%	5.12%	6.4%	5.98%	
	<i>Improv.s</i>	9.06%	9.57%	9.44%	7.66%	8.7%	8.06%	6.62%	10.02%	8.83%	
	Foursquare	BPR	Teacher	0.4790	0.3318	0.3684	0.5827	0.3457	0.4020	0.6748	0.3521
Student			0.3867	0.2531	0.2865	0.4909	0.2670	0.3202	0.5833	0.2738	0.3436
RD			0.4179	0.2760	0.3113	0.5211	0.2896	0.3444	0.6227	0.2958	0.3696
CD			0.4025	0.2633	0.2979	0.5030	0.2769	0.3306	0.6053	0.2822	0.3550
DE **			0.4079	0.2625	0.2986	0.5139	0.2766	0.3328	0.6238	0.2843	0.3607
RRD ***		0.4737	0.3086	0.3497	0.5800	0.3236	0.3847	0.6765	0.3305	0.4094	
DE-RRD ****		0.4758	0.3108	0.3518	0.5805	0.3246	0.3856	0.6770	0.3312	0.4099	
<i>Improv.b</i>		13.83%	12.6%	13.03%	11.42%	12.09%	11.95%	8.72%	11.95%	10.9%	
<i>Improv.s</i>		23.03%	22.79%	22.8%	18.26%	21.58%	20.42%	16.07%	20.95%	19.28%	

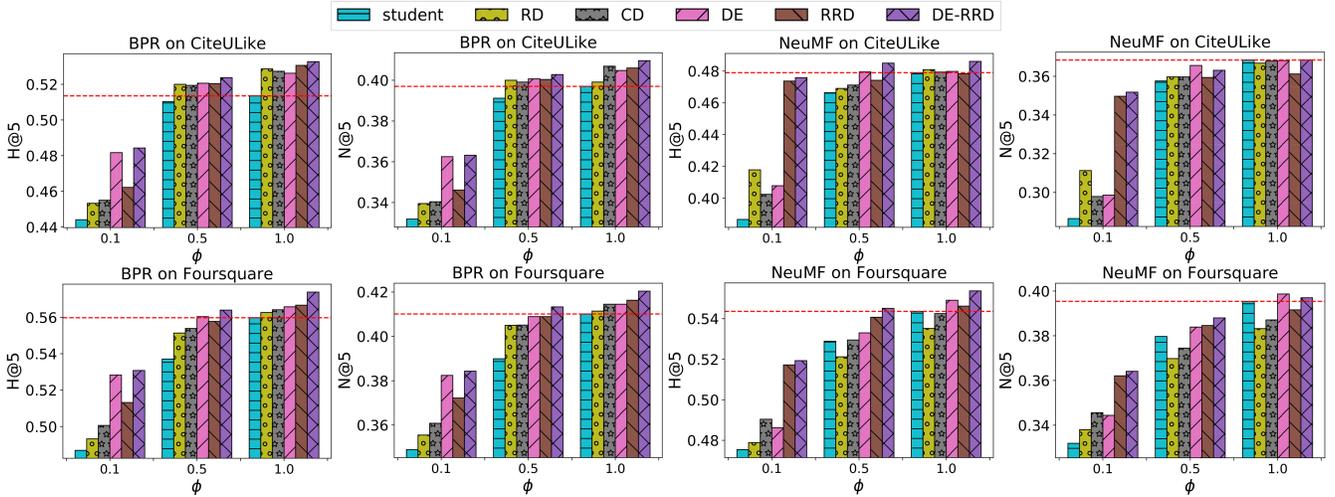


Figure 4: Recommendation Performance on across three different student model sizes. (Red dotted line: Teacher)

it can be also used to maximize the performance of the existing recommender.

Lastly, we provide the result of the online inference efficiency test in Table 3. All inferences are made using PyTorch with CUDA from Tesla P40 GPU and Xeon on Gold 6148 CPU. The student model trained with DE-RRD achieves comparable performance with only 10-50% of learning parameters compared to the teacher.

The smaller model requires less computations and memory costs, so it can achieve lower latency. In particular, deep recommender (i.e., NeuMF) which has a large number of learning parameters and complex structures takes more benefits from the smaller model size. On real-time RS application that has larger numbers of users (and items) and has a more complex model structure, DE-RRD can lead to a larger improvement in online inference efficiency.

Table 3: Model compactness and online inference efficiency. Time (seconds) indicates the wall time used for generating recommendation list for every user. H@5 Ratio denotes the ratio of H@5 from DE-RRD over that from Teacher.

Dataset	Base Model	ϕ	Time (s)	#Params.	H@5 Ratio
CiteULike	BPR	1.0	59.27s	6.08M	1.03
		0.5	57.53s	3.04M	1.01
		0.1	55.39s	0.61M	0.94
	NeuMF	1.0	79.27s	15.33M	1.01
		0.5	68.37s	7.63M	1.01
		0.1	58.27s	1.52M	0.99
Foursquare	BPR	1.0	257.28s	9.61M	1.03
		0.5	249.19s	4.81M	1.01
		0.1	244.23s	0.96M	0.95
	NeuMF	1.0	342.84s	24.16M	1.02
		0.5	297.34s	12.05M	1.01
		0.1	255.24s	2.40M	0.95

5.3 Design Choice Analysis

We provide both quantitative and qualitative analyses on the proposed methods and alternative design choices (i.e., ablations) to verify the superiority of our design choice. The performance comparisons with the ablations are summarized in Table 4.

For DE, we consider three ablations: (a) Attention (b) One expert (large) (c) One expert (small). As discussed in Section 4.1.2, instead of the selection strategy, attention mechanism can be adopted. We also compare the performance of one large expert³ and one small expert. Note that DE, attention, and one expert (large) has the exact same number of learning parameters for experts. We observe that the increased numbers of learning parameters do not necessarily contribute to performance improvement ((a) vs. (c) in BPR).

We also observe that the selection shows the best performance among all the ablations. To further investigate this result, we conduct qualitative analysis on user representation spaces induced by each design choice. Specifically, we first perform clustering⁴ on user representation space from the teacher model to find user groups that have strong correlations (or similar preferences). Then, we visualize the average performance gain (per group) map in Figure 5. We observe that distilling the knowledge by the attention, one large expert can cause performance decreases in many user groups (blue clusters), whereas the selection improves the performance in more numbers of user groups (red clusters). In the ablations (a)-(c), the experts are trained to minimize the overall reconstruction errors on all the diverse entities. This makes the information of weakly correlated entities to be mixed together and further hinders discovering the preference of a particular user group. Unlike the ablations, DE clearly distinguishes the knowledge that each expert distills, and makes each expert to be trained to distill only the knowledge of strongly correlated entities. So, it can alleviate such problem. The expert selection map of DE is visualized in Figure 6. We can observe that each expert gets gradually specialized in certain user groups that share similar preferences during the training.

For RRD, we consider two ablations: (d) and (e). The ablations are intended to show the effects of the proposed relaxed ranking. Concretely, we apply the list-wise loss (i.e., no relaxation) on all the sampled items (interesting and uninteresting items) for (d), on the top-ranked items (interesting items) for (e). Note that all the methods use the same number of items for distillation. We observe

³We make one large expert by adopting the average pooling.

⁴We use k -Means clustering in Scikit-learn. k is set to 20.

Table 4: Performance comparison for alternative design choices on Foursquare ($\phi = 0.1$).

Base Model	Design choices	H@5	N@5	H@10	N@10
BPR	DE	0.5283	0.3824	0.6810	0.4316
	(a) Attention	0.5019	0.3625	0.6575	0.4131
	(b) One expert (large)	0.5151	0.3716	0.6733	0.4230
	(c) One expert (small)	0.5136	0.3717	0.6683	0.4213
	RRD	0.5132	0.3722	0.6616	0.4202
NeuMF	(d) Full ranking	0.4983	0.3595	0.6474	0.4080
	(e) <i>Interesting</i> ranking	0.4814	0.3479	0.6416	0.3999
	DE	0.4862	0.3444	0.6413	0.3938
	(a) Attention	0.4770	0.3364	0.6364	0.3903
	(b) One expert (large)	0.4741	0.3367	0.6341	0.3885
(c) One expert (small)	0.4740	0.3339	0.6316	0.3860	
RRD	(d) Full ranking	0.5172	0.3621	0.6739	0.4132
	(e) <i>Interesting</i> ranking	0.4799	0.3457	0.6324	0.3949
	(e) <i>Interesting</i> ranking	0.4641	0.3294	0.6228	0.3809

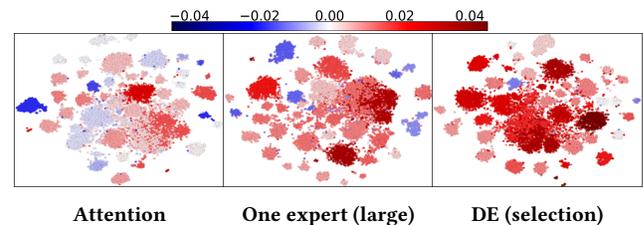


Figure 5: Performance (N@20) gain map (BPR with $\phi = 0.1$ on Foursquare).

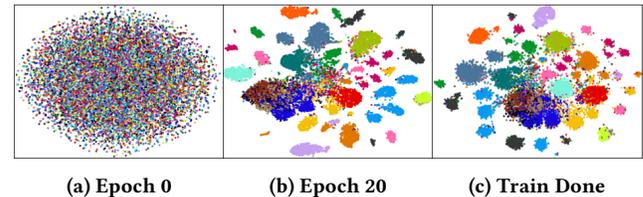


Figure 6: Expert selection map of DE. Each color corresponds to an expert (BPR with $\phi = 0.1$ on Foursquare).

that merely adopting the list-wise loss has adverse effects on the ranking performance. First, (d) learns to match the full ranking order among all the sampled items. Learning the detailed order among the uninteresting items is not necessarily helpful to improve the ranking performance, and may further interfere with focusing on the interesting items. Also, (e), which only considers the interesting items, shows even worse performance than Student. The list-wise loss does not take into account the absolute ranking positions of the items; a ranking order can be satisfied regardless of the items' absolute positions. Since (e) does not consider the relative orders between the interesting items and the uninteresting items, it may push such interesting items far from the top of the ranking list. Unlike the ablations, RRD adopts the relaxed ranking approach, which enables the student to better focus on the interesting items while considering the relative orders with the uninteresting items.

5.4 Hyperparameter Analysis

We provide analyses to offer guidance of hyperparameter selection of DE-RRD. For the sake of space, we report the results on Foursquare dataset with $\phi = 0.1$. We observe similar tendencies on CiteULike dataset. For DE, we show the effects of two hyperparameters: λ_{DE} that controls the importance of DE and the number

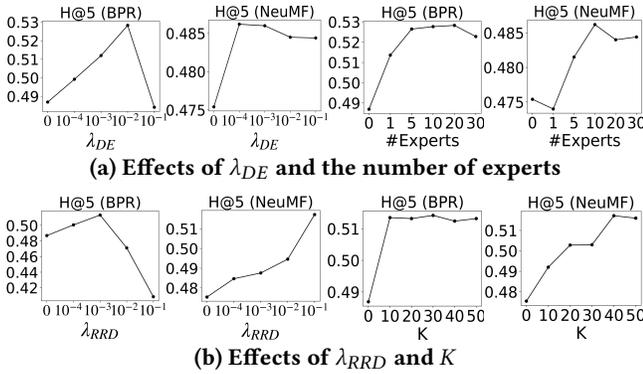


Figure 7: Effects of the hyperparameters. (a) DE (b) RRD.

Table 5: Effects of λ_{DE} and λ_{RRD} in DE-RRD framework.

Foursquare (H@5)		BPR				NeuMF			
		λ_{RRD}				λ_{RRD}			
λ_{DE}		10^{-4}	10^{-3}	10^{-2}	10^{-1}	10^{-4}	10^{-3}	10^{-2}	10^{-1}
10^{-4}		0.5081	0.5201	0.4590	0.3901	0.4774	0.4896	0.5014	0.5193
10^{-3}		0.5186	0.5276	0.4688	0.3906	0.4774	0.4858	0.4942	0.5112
10^{-2}		0.5261	0.5308	0.4791	0.3977	0.4846	0.4868	0.4892	0.5110
10^{-1}		0.5269	0.5308	0.4928	0.4154	0.4848	0.4881	0.4908	0.5055

of experts in Figure 7a. For RRD, we show the effects of two hyperparameters: λ_{RRD} that controls the importance of RRD and the number of interesting items (K) in Figure 7b. In our experiment, the number of uninteresting items is set to the same with K . Note that for all graphs value ‘0’ corresponds to Student (i.e., no distillation).

Because the types of loss function of the proposed methods are different from that of the base models, it is important to properly balance the losses by using λ . For DE, the best performance is achieved when the magnitude of DE loss is approximately 20% (BPR), 2-5% (NeuMF) compared to that of the base model’s loss. For RRD, the best performance is achieved when the magnitude of RRD loss is approximately 7-10% (BPR), 1000% (NeuMF) compared to that of the base model’s loss. For the number of experts and K , the best performance is achieved near 10-20 and 30-40, respectively. Lastly, we show the effects of combinations of λ_{DE} and λ_{RRD} in DE-RRD framework in Table 5. Generally, the best performance of DE-RRD is observed in the ranges where each method (i.e., DE, RRD) achieves the best performance.

6 CONCLUSION

This paper proposes a novel knowledge distillation framework for recommender system, DE-RRD, that enables the student model to learn both from the teacher’s predictions and from the latent knowledge stored in a teacher model. To this end, we propose two novel methods: 1) DE that directly distills latent knowledge from the representation space of the teacher. DE adopts the experts and the expert selection strategy to effectively distill the vast CF knowledge to the student. 2) RRD that distills knowledge revealed from teacher’s predictions with direct considerations of ranking orders among items. RRD adopts the relaxed ranking approach to better focus on the interesting items. Extensive experiment results demonstrate that DE-RRD significantly outperforms the state-of-the-art competitors.

Acknowledgments: This research was supported by the NRF grant funded by the MSIT: (No. 2017M3C4A7063570), the IITP grant funded by the MSIT: (No. 2018-0-00584), the IITP grant funded

by the MSIP (No. 2019-0-01906, Artificial Intelligence Graduate School Program (POSTECH)) and the MSIT under the ICT Creative Consilience program (IITP-2020-2011-1-00783) supervised by the IITP.

REFERENCES

- [1] Yoram Bachrach, Yehuda Finkelstein, Ran Gilad-Bachrach, Liran Katzir, Noam Koenigstein, Nir Nice, and Ulrich Paquet. 2014. Speeding up the xbox recommender system using a euclidean transformation for inner-product spaces. In *RecSys*.
- [2] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* (1975).
- [3] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. 2017. Learning efficient object detection models with knowledge distillation. In *NeurIPS*. 742–751.
- [4] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S Mirrokni. 2004. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*. 253–262.
- [5] Tommaso Furlanello, Zachary C Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. 2018. Born again neural networks. *arXiv preprint arXiv:1805.04770* (2018).
- [6] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *WWW*.
- [7] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).
- [8] Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144* (2016).
- [9] SeongKu Kang, Junyoung Hwang, Dongha Lee, and Hwanjo Yu. 2019. Semi-supervised learning for cross-domain recommendation to cold-start users. In *CIKM*.
- [10] Wang-Cheng Kang and Julian McAuley. 2019. Candidate Generation with Binary Codes for Large-Scale Top-N Recommendation. In *CIKM*.
- [11] Wonbin Kweon, Seongku Kang, Junyoung Hwang, and Hwanjo Yu. 2020. Deep Rating Elicitation for New Users in Collaborative Filtering. In *WWW*.
- [12] Dongha Lee, Chanyoung Park, Hyunjun Ju, Junyoung Hwang, and Hwanjo Yu. 2019. Action Space Learning for Heterogeneous User Behavior Prediction. In *IJCAI*.
- [13] Jaewoong Lee, Minjin Choi, Jongwuk Lee, and Hyunjung Shim. 2019. Collaborative Distillation for Top-N Recommendation. *ICDM* (2019).
- [14] Hui Li, Tsz Nam Chan, Man Lung Yiu, and Nikos Mamoulis. 2017. FEXIPRO: fast and exact inner product retrieval in recommender systems. In *Proceedings of the 2017 ACM International Conference on Management of Data*.
- [15] Defu Lian, Rui Liu, Yong Ge, Kai Zheng, Xing Xie, and Longbing Cao. 2017. Discrete Content-Aware Matrix Factorization. In *KDD*.
- [16] Han Liu, Xiangnan He, Fuli Feng, Liqiang Nie, Rui Liu, and Hanwang Zhang. 2018. Discrete factorization machines for fast feature-based recommendation. *arXiv preprint arXiv:1805.02232* (2018).
- [17] Yiding Liu, Tuan-Anh Nguyen Pham, Gao Cong, and Quan Yuan. 2017. An experimental evaluation of point-of-interest recommendation in location-based social networks. (2017).
- [18] Chris J Maddison, Daniel Tarlow, and Tom Minka. 2014. A* sampling. In *NeurIPS*.
- [19] C. Park, D. Kim, X. Xie, and H. Yu. 2018. Collaborative Translational Metric Learning. In *ICDM*.
- [20] Steffen Rendle and Christoph Freudenthaler. 2014. Improving pairwise learning for item recommendation from implicit feedback. In *WSDM*.
- [21] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *UAI*.
- [22] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2014. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550* (2014).
- [23] Anshumali Shrivastava and Ping Li. 2014. Asymmetric LSH (ALSH) for sublinear time maximum inner product search (MIPS). In *NeurIPS*.
- [24] Bo Song, Xin Yang, Yi Cao, and Congfu Xu. 2018. Neural collaborative ranking. In *CIKM*.
- [25] Jiayi Tang and Ke Wang. 2018. Ranking distillation: Learning compact ranking models with high performance for recommender system. In *KDD*.
- [26] Saúl Vargas, Craig Macdonald, and Iadh Ounis. 2015. Analysing compression techniques for in-memory collaborative filtering. (2015).
- [27] Hao Wang, Binyi Chen, and Wu-Jun Li. 2013. Collaborative topic regression with social regularization for tag recommendation. In *IJCAI*.
- [28] Haoyu Wang, Defu Lian, and Yong Ge. 2019. Binarized collaborative filtering with distilling graph convolutional networks. *IJCAI* (2019).
- [29] Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. Listwise approach to learning to rank: theory and algorithm. In *ICML*.
- [30] Hanwang Zhang, Fumin Shen, Wei Liu, Xiangnan He, Huanbo Luan, and Tat-Seng Chua. 2016. Discrete Collaborative Filtering. In *SIGIR*.