

# Double Correction Framework for Denoising Recommendation

Zhuangzhuang He<sup>\*†</sup>  
School of Computer Science and  
Information Engineering, Hefei  
University of Technology  
hyicheng223@gmail.com

Yifan Wang<sup>†</sup>  
Department of Computer Science and  
Technology, Tsinghua University  
yf-wang21@mails.tsinghua.edu.cn

Yonghui Yang  
School of Computer Science and  
Information Engineering, Hefei  
University of Technology  
yyh.hfut@gmail.com

Peijie Sun  
Department of Computer Science and  
Technology, Tsinghua University  
sun.hfut@gmail.com

Le Wu<sup>‡</sup>  
School of Computer Science and  
Information Engineering, Hefei  
University of Technology  
Institute of Dataspace,  
Hefei Comprehensive National  
Science Center  
lewu.ustc@gmail.com

Haoyue Bai  
School of Computer Science and  
Information Engineering, Hefei  
University of Technology  
baihaoyue621@gmail.com

Jinqi Gong  
Department of Mathematics,  
University of Macau  
eggmangong@gmail.com

Richang Hong  
School of Computer Science and  
Information Engineering, Hefei  
University of Technology  
hongrc.hfut@gmail.com

Min Zhang<sup>‡</sup>  
Department of Computer Science and  
Technology, Tsinghua University  
z-m@tsinghua.edu.cn

## ABSTRACT

As its availability and generality in online services, implicit feedback is more commonly used in recommender systems. However, implicit feedback usually presents noisy samples in real-world recommendation scenarios (such as misclicks or non-preferential behaviors), which will affect precise user preference learning. To overcome the noisy sample problem, a popular solution is based on dropping noisy samples in the model training phase, which follows the observation that noisy samples have higher training losses than clean samples. Despite the effectiveness, we argue that this solution still has limits. (1) High training losses can result from model optimization instability or hard samples, not just noisy samples. (2) Completely dropping of noisy samples will aggravate the data sparsity, which lacks full data exploitation.

To tackle the above limitations, we propose a *Double Correction Framework for Denoising Recommendation (DCF)*, which contains two correction components from views of more precise sample dropping and avoiding more sparse data. In the sample dropping

correction component, we use the value of the sample loss over time to determine whether it is noise or not, increasing stability. Instead of averaging directly, we use the damping function to reduce the bias effect of outliers. Furthermore, due to the higher variance exhibited by hard samples, we derive a lower bound for the loss through concentration inequalities to identify and reuse hard samples. In progressive label correction, we iteratively re-label highly deterministic noisy samples and retrain them to further improve performance. Finally, extensive experimental results on three datasets and four backbones demonstrate the effectiveness and generalization of our proposed framework. Our code is available at <https://github.com/bruno686/DCF>.

## CCS CONCEPTS

• **Information systems** → **Recommender systems; Collaborative filtering.**

## 1 INTRODUCTION

Recommender systems often use implicit feedback (e.g., click, watch, and purchase) to learn user interests and recommend items, where users' implicit feedback is thought to reflect users' true preferences [1–3, 26, 27, 34, 35, 45, 48]. However, here are some exceptions. For example, a user may purchase and return an item or click on a video but the dwell time is very short, thus the purchasing and clicking behavior may not represent users' interests, which is referred to *noisy interaction* by recent research [32, 36]. Recent research [32, 36] refer to this type of interaction as *noisy interaction*. Treating these noisy interactions as clean interactions may get sub-optimal performance [32, 36, 39]. Therefore, how to eliminate the

<sup>\*</sup>This work is produced during a research intern at Tsinghua University.

<sup>†</sup>Equal Contribution.

<sup>‡</sup>Corresponding authors.

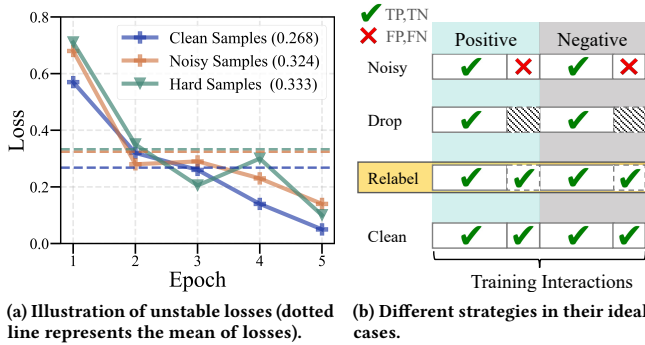
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components 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](mailto:permissions@acm.org).

KDD'24, August 25–29 2024, Barcelona, Spain

© 2024 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

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



**Figure 1: (a) Illustration of unstable losses.** We observe that clean samples do not exhibit low loss in every epoch. Similarly, noisy samples do not always exhibit high loss. Also, hard samples exhibit high losses. However, the noisy samples can be identified from the perspective of mean loss. **(b) Different strategies in their ideal cases.** Explain the difference between our relabeling strategy and other strategies under ideal conditions. Here, the TP, TN of T stands for true. Similarly, F stands for false.

adverse effects of noisy interactions has attracted a broad interest from the research community.

Although some noisy interactions may be detected by multiple behaviors (e.g., explicit feedback), we cannot guarantee that other behaviors validate every interaction due to the sparsity of user behavior. Additionally, the request for explicit information is expensive as it may hurt the user experience. Thus, it is important to discover noisy interactions based on their unique data pattern, which has also received much attention in recent years. A commonly observed pattern, which is also the basis of most current work [11, 21, 32], is that the loss values for noisy interactions are generally higher than for clean ones during training. Hence, a straightforward solution is to use loss values as a distinguishing feature to process samples with high losses. For instance, [21, 32] directly and dynamically drops samples with high losses during training, achieving notable results.

Despite the notable performance, we argue that current methods [11, 21, 32] have two limitations that may hinder effectiveness. (1) They ignore that losses may not be highly correlated with noise, i.e., the instability of optimization and the hardness of clean interactions. Firstly, the instability in the optimization process may lead to sharp changes in the immediate loss values of the sample, which may show a contrary phenomenon to the above basis. As illustrated in Fig. 1 (a), noisy samples also show lower training loss sometimes, and vice versa, which leads to incorrect dropping and may hurt the performance. Secondly, hard samples also typically show high loss and could be dropped along with noisy samples. Abandoning these hard samples may lead to suboptimal performance as they are informative and beneficial for recommendation performance [24]. (2) Simply dropping samples may lead to more sparse data. As illustrated in Fig. 1 (b), although the noisy interactions are mislabeled, they are still part of the training space. Merely dropping these noisy

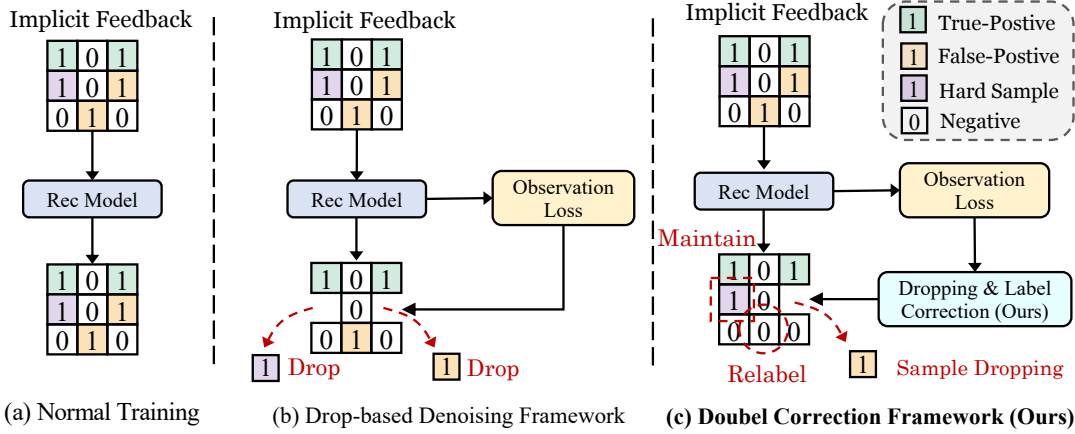
samples leads to sample wastage, as well as potentially causing the training space to be inconsistent with the ideal clean space.

To address these limitations, we explore potential solutions and the corresponding challenges. For the first point, we attribute the low correlation between loss and noise to limited observations of current model predictions and the neglect to consider hard samples. Hence, we expect to stabilize the model’s predictions by extending the observation interval and aggregating loss values from multiple training iterations. In addition, high-loss samples might be hard samples beneficial for training. Therefore, we aim to identify and retain these hard samples to enhance recommendation performance. However, how to efficiently and simply identify hard samples is still challenging. For the second point, we argue that even if noisy samples, have corresponding correct labels and cannot be merely dropped from the sample space, which would result in more sparse space. Hence, we would like to relabel and reintroduce part of noisy samples that are highly determined to be noise into the training process. And we illustrate in Fig. 1(b) how relabeling addresses the issue of sample waste caused by drops. However, there is a practical challenge in determining which samples need to be relabelled and the proportion of relabelling during model training.

In this paper, we propose a Double Correction Framework for Denoising Recommendation (**DCF**). The core of this framework is sample dropping correction and progressive label correction, which are designed for the above two limitations respectively. Sample dropping correction combines confirmed loss calculation and cautious hard sample search to accurately drop noisy samples and retain hard samples. Specifically, we calculate the mean loss of samples over a time interval and robustly compute each loss value using damping functions to mitigate the effects of occasional outliers. Inspired by [8], hard samples have a higher loss variance than clean and noisy samples. In other words, the loss value of hard samples has a lower bound in terms of the whole training process. Therefore, we use concentration inequalities to derive confidence intervals for each sample’s loss value. Then we use the lower bound of the calculated confidence interval as a hard sample search criterion. In this way, hard samples are retained for training instead of being dropped. In the progressive label correction component, we believe the stability of the model optimization process increases gradually [13], so the model’s relabeling strategy should adapt to this characteristic. Specifically, we initially relabel a small fraction of the samples and progressively increase the proportion of relabeling as training proceeds. Note that we still drop samples with high loss values and just relabel a fraction of the samples with highly determined to be noisy.

To summarize, our main contributions are as follows:

- We analyze two limitations of the loss-based dropping recommendations, (1) loss values are not highly correlated with noise, and (2) complete dropping of noisy samples, which leads to more sparse data space.
- We design two correction components to enhance the effectiveness of denoising recommendations from more precise sample dropping and reusing noisy sample perspectives, respectively.
- We experiment with our proposed DCF of four backbones on three popular datasets. And extensive experimental results demonstrate the effectiveness and generalization of our framework.



**Figure 2: Illustration comparison: (a) Normal training model without denoising, (b) Denoising model with drop strategy, (c) Double correction framework for denoising recommendation (Ours).**

## 2 TASK DESCRIPTION

We first give a formal description of the denoising recommendation task. We use  $u \in \mathcal{U}$  to denote users,  $p \in \mathcal{P}$  to denote the item, and observed interaction matrix  $\tilde{Y} \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{P}|}$ . And where  $\tilde{y}_{ui} = 1$  means that the user interacted with the item, and  $\tilde{y}_{ui} = 0$  means no interaction. In previous work, the default assumption is that whenever  $\tilde{y}_{ui} = 1$ , it means that the user likes the item. However, user interactions may be due to various noises (i.e., clicking on a video by mistake or clicking on a video for a very short time to exit), resulting in users not liking the interacted item. Therefore, the target of denoising recommendations is to develop a model  $f$  with parameters  $\theta_f$  to learn noise-free representation of users and items from noisy data  $\tilde{Y}$ . The training of denoising recommendation model is formulated as follows:

$$\theta_f^* = \arg \min_{\theta_f} \mathcal{L}_{\text{rec}}(\tilde{Y}),$$

where  $\mathcal{L}_{\text{rec}}$  is the recommendation loss, such as BPR loss [40] and BCE loss [32], and  $\theta_f^*$  is the optimal parameters of  $f$ . Here, following previous denoising recommendation works [21, 32, 36], we use the BCE loss as an instantiation of  $\mathcal{L}_{\text{rec}}$ :

$$\mathcal{L}_{\text{rec}} = - \mathbb{E}_{(u,p) \sim P_{\tilde{Y}}} [\tilde{y}_{u,p} \cdot \log(\hat{y}_{u,p}) + (1 - \tilde{y}_{u,p}) \cdot \log(1 - \hat{y}_{u,p})],$$

where  $P_{\tilde{Y}}(\cdot)$  denotes the distribution established over the interaction data. The  $\hat{y}_{u,p}$  represents the prediction made by the model for the user's preference. In addition, the loss of one interaction is  $\ell$ .

## 3 THE PROPOSED FRAMEWORK

### 3.1 Overview

Our framework DCF (as shown in Fig. 2 (c)) is developed on a common observation: noisy samples usually incur high losses [32]. However, on the one hand, the model optimization is usually unstable, and hard samples also show high losses; on the other hand, some methods [12, 32] directly drop samples with high loss values for simplicity. Such actions may lose feature information and degrade performance. To mitigate these limitations, we design two components (as shown in Fig. 3): sample dropping correction and

progressive label correction. The first component, sample dropping correction, aims to correct unstable loss values as model random initialization parameters and unstable optimization processes. In addition, we also calculate the confidence intervals for each sample loss and use the lower bound as our dropping strategy. This method aims to retain hard samples because hard samples enhance model performance by exposing edge cases and revealing more profound insights into data patterns [24]. Additionally, for those samples that are highly determined to be noise, the second component, progressive label correction, relabels these and adds them to the next training iteration to improve performance further. Based on these two components, our framework mitigates the adverse effects of noisy interactions more effectively. The complete algorithmic process is at Algorithm 1.

### 3.2 Sample Dropping Correction

Due to the instability of the optimization process and randomness of the initial weights. As shown in Fig. 1 (a), true positive interactions may exhibit high loss at some time points, while noisy interactions may exhibit low ones. Fortunately, we also observe that the sample loss can be more stable by replacing the loss value at a single time point with the mean loss over time. Thus, inspired by this, we estimate the sample loss using a confirmed approach. Besides, samples that perform with high loss may also be hard samples, which can enhance the model performance [24]. In a word, our goal in this subsection is to drop the noisy samples more accurately, find the hard samples, and retain them.

**3.2.1 Confirmed Loss Calculation.** Relying only on empirical losses calculated from a single training iteration may incorrectly distinguish between noise and clean, thus degrading performance. Therefore, we consider the loss values across previous training iterations to mitigate misjudge. Specific practice is to use the mean of losses at different iterations.

$$\mu_i = \frac{1}{v} \sum_{j=i-v+1}^i \ell(j), \quad (1)$$

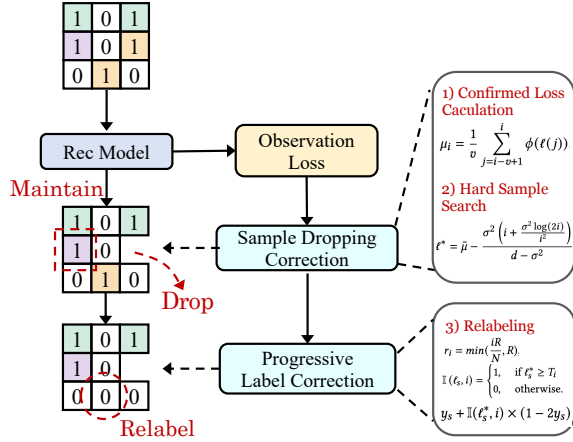


Figure 3: Details of Loss & Label Correction.

where we calculate the mean loss value  $\mu_i$  over  $v$  time intervals prior to the  $i$ -th training iteration (including the  $i$ -th). And if the current iteration is less than  $v - 1$  times,  $\mu_i$  is the mean of all losses up to the current iteration  $i$ . However, during the training process, the model may encounter extreme loss values. Although this possibility is unlikely, to avoid negative impacts when calculating means, we need to take defensive measures for robust computation. Thus, instead of calculating the mean directly, we robustly process each loss. After that, we obtain the mean as follows:

$$\mu_i = \frac{1}{v} \sum_{j=i-v+1}^i \phi(\ell(j)), \quad (2)$$

where, the non-decreasing damping function  $\phi(\cdot)$  is used to calculate mean values more robustly. The specific forms are as follows:  $\phi(\ell) = \log(1 + \ell + \ell^2/2)$ .

The advantages of the damping function include: (1) **Reduced impact of extremes**. The damping function grows slowly at very large loss values because it is logarithmic. These large values may be outliers due to unstable model predictions from noisy samples. Therefore, the effect of extreme values is reduced when the mean is calculated afterward. (2) **Approximate linearity at small values**. At small values,  $\phi(\cdot)$  can be approximated as linear function. This means that near small values,  $\phi(\cdot)$  does not have side effects on the loss values, preserving the original.

**3.2.2 Cautious Hard Samples Search.** We argue that dropping samples with higher mean loss still leaves room for potential improvements. This is because a higher mean sample loss does not necessarily indicate noisy samples; they could be hard samples. [32] point out that the hard samples contain more information than the simple ones. Thus we aim to search and retain hard samples. Our search strategy is based on the observation [8] that hard samples exhibit higher variance in loss values during training. Consequently, we compute the lower confidence interval bound for the loss by *Concentration Inequalities* [41]. Because, for noisy samples, the lower bound of the loss confidence interval usually closely matches the actual loss value. However, for hard samples with high variance, the confidence lower bound tends to be lower. Therefore, in our sample dropping and subsequent relabeling strategy, we use the

lower bound as our criterion rather than the loss value itself. We define the form of centralized inequality specified below.

**Theorem 1.** Let  $Z_n = \{z_1, \dots, z_n\}$  be an observation set with mean  $\mu_z$  and variance  $\sigma^2$ . By exploiting the non-decreasing damping function  $\phi(z) = \log(1 + z + z^2/2)$ . For any  $\epsilon > 0$ , we have

$$\left| \frac{1}{n} \sum_{i=1}^n \phi(z_i) - \mu_z \right| \leq \frac{\sigma^2 \left( n + \frac{\sigma^2 \log(\epsilon^{-1})}{n^2} \right)}{n - \sigma^2}, \quad (3)$$

with probability at least  $1 - 2\epsilon$ .

There are other methods to filter out the hard samples, such as gradient [25] and cluster [42]. Compared to these methods, our framework has the following advantages. (1) **Intuition**. Lower bound, as a basic statistical concept, can intuitively reflect the degree of data dispersion. For hard samples, their loss values fluctuate more during the training process, and thus the lower bound is lower. (2) **Computational efficiency**. The low complexity of computing lower bound is suitable for large-scale datasets. Compared to other selection methods, the lower bound of sample loss provides a simple and efficient strategy.

Let,  $\epsilon = \frac{1}{2i}$ , the following loss lower bound  $\ell^*$  form is obtained after derivation. The detailed derivation process is in Appendix A.1.

$$\ell^* = \tilde{\mu} - \frac{\sigma^2 \left( i + \frac{\sigma^2 \log(2i)}{i^2} \right)}{d - \sigma^2}, \quad (4)$$

where  $d$  represents the number of times a sample has not been dropped. Intuitively, the smaller  $d$  is, the lower the confidence interval bound will be. We prioritize reincorporating samples with lower confidence interval bounds (i.e., higher variance) into training. In addition, the search intensity is retained through an adjustment factor,  $\sigma^2$ , that reduces the noisy samples' adverse effects and increases the model performance. It is important to note that we are not selecting samples based on their loss values but rather based on the lower bound.

Through the sample dropping correction strategy, we achieve higher confirmed sample loss values and effectively search and retain hard samples, further enhancing the model performance.

### 3.3 Progressive Label Correction

Compared to directly dropping noisy samples. We argue that the small fraction of samples with the highest loss<sup>1</sup> can be leveraged in a relabeling way. Relabeling noisy labels preserves samples rather than dropping them, providing more training samples. This is valuable for sparse recommendation datasets, which makes the entire training sample space and test set space more consistent. Moreover, relabeling is a simple operation that does not require complex reweighting strategies or other processing techniques. Thereby further training the model adequately to improve its performance. In addition, since the predictive stability of the model improves incrementally during training, we argue it is reasonable to relabel more in the later stages and less in the earlier stages. From the beginning to the end of the training, we progressively increased the percentage of relabeling.

<sup>1</sup>For simplicity in description, we still use the loss to represent the lower bound of the loss value. But we have the notational distinction between  $\ell$  and  $\ell^*$

**Algorithm 1:** Our proposed DCF.

---

**Input** : Training set  $D$ , Maximum Epochs  $N$ , Time interval  $v$ , Relabel Ratio  $R$ ;  
**Output**: Trained Model  $MM$ .

- 1 Initialize model  $M$  with random weights;
- 2 Initialize empty loss history  $L$  for each sample in  $D$ ;
- 3 **for**  $i = 1$  **to**  $N$  **do**
- 4     **for** each mini-batch  $B$  from  $D$  **do**
- 5         // Sample Dropping Correction
- 6         **for** each sample  $s$  in  $B$  **do**
- 7             Compute loss  $\ell$  using model  $M$ ;
- 8             Append  $\ell$  to loss history  $L[s]$ ;
- 9             **if**  $\text{length}(L[s]) > v$  **then**
- 10                 Remove samples in  $L[s]$  with time interval greater than  $v$ ;
- 11             **else**
- 12                 Compute mean loss  $\mu$  using Eq. (2);
- 13                 Compute lower bound  $\ell^*$  using Eq. (4);
- 14             **end**
- 15         **end**
- 16         Update model  $M$  using lower bound  $\ell^*$ ;
- 17     **end**
- 18     // Progressive Label Correction
- 19     Compute relabel ratio  $r$  using Eq. (5);
- 20     **for** each sample  $s$  in  $D$  **do**
- 21         **if**  $\ell^* \geq T_i$  ( $T_i$  is computed based on  $r$  and  $\ell^*$ ) **then**
- 22             Flip label  $y_s$  using Eq. (7);
- 23         **else**
- 24             Keep label  $y_s$ ;
- 25         **end**
- 26     **end**
- 27     Training the model with the corrected samples at the next epoch;
- 28 **end**
- 29 **return** Trained Model  $M$ ;

---

We first define the percentage of progressive relabeling. Our relabel ratio  $r_i$  grows as epoch  $i$  increases. It is defined as follows:

$$r_i = \min\left(\frac{iR}{O}, R\right), \quad (5)$$

where  $O$  denotes that the flipping rate remains constant after the  $O$ -th epoch, and  $R$  is the percentage we end up relabeling. This way, we get each epoch's relabeled ratio  $r_i$ . And the corresponding loss threshold of the samples to be labeled is  $T_i = \lfloor \lfloor B(1 - r_i) \rfloor \rfloor$ .  $\lfloor \cdot \rfloor$  is a downward rounding operation. If the  $\ell_i$  is greater than  $T_i$ , the indicator  $\mathbb{I}(\ell_s, i)$  is 1.

$$\mathbb{I}(\ell_s^*, i) = \begin{cases} 1, & \text{if } \ell_s^* \geq T_i \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

Up to this point, we know which samples we should flip the labels on. We flip labels in the dataset so that the flipped labels are also used at subsequent training epochs.

$$y'_s = y_s + \mathbb{I}(\ell_s^*, i) \times (1 - 2y_s), \quad (7)$$

where we relabel the original label from  $y_s$  to  $y'_s$ . After careful processing, we relabel these noisy labels to ensure they are closer to the true ones. This step not only enhances the data quality but also provides more accurate labels for subsequent model optimization, which is expected to improve the effectiveness.

Through the progressive label correction strategy, we have avoided sample waste, thereby providing more data for the precise modeling of users and items.

### 3.4 Model Discussion

This section compare different recommendation models based on their space and time complexities. The comparison is summarized in Table 1. In addition, we discuss the improved features of DCF to the previous dropping approach T-CE [32].

**Space Complexity.** The space complexity of each model is determined by the number of parameters, denoted as  $M$ . The base model and T-CE model have the same space complexity of  $M$  since T-CE does not add any new parameters or structures. However, BOD includes a weight generator with an encoder layer  $EN \in \mathbb{R}^{2d \times d_g}$  and a decoder layer  $DE \in \mathbb{R}^{d_g \times 1}$ , where  $d$  is the embedding size of user and item in the recommendation model, and  $d_g$  is the hidden layer size of the generator. On the other hand, DCF introduces additional computations like confidence bounds of the loss but does not add new model parameters, meaning its space complexity is also  $M$ .

**Time Complexity.** Regarding time complexity, the base model computes the loss directly in  $O(n)$  time, where  $n$  is the number of samples. T-CE sorts the computed loss, resulting in a time complexity of  $O(n \log n)$ . BOD involves bi-level optimization, leading to a time complexity of  $O(n + d(d_g) + (d_g))$ . Our proposed method includes multiple steps like computing confidence bounds, soft processing, and sorting. Sorting is the dominant factor here, leading to a time complexity of  $O(n \log n)$ .

**Table 1: Model complexity comparison.**

Model	Space Complexity	Time Complexity
Base [19]	$M$	$O(n)$
T-CE [32]	$M$	$O(n \log n)$
BOD [37]	$M + EN + DE$	$O(n + d(d_g) + (d_g))$
<b>DCF (Ours)</b>	$M$	$O(n \log n)$

**In-depth Comparison.** The T-CE method overlooks the potential of hard samples, possibly hindering model performance. Our cautious hard sample search, however, recognizes the value in these samples. By focusing on the lower bound of loss rather than just mean loss, we can better differentiate between noise and hard challenges. This approach ensures we capture the full learning spectrum, addressing T-CE's limitation.

## 4 EXPERIMENTS

To demonstrate the effectiveness and generalization of our proposed method, we compare the results of our DCF with state-of-the-art recommendation models on four backbones and three real-world datasets. We aim to answer the following research questions:

- **RQ1** : How does our proposed method perform compared to normal training and other state-of-the-art denoising methods?

- **RQ2** : How does each component in the DCF contribute to the overall performance and the impact of hyperparameters?
- **RQ3** : Do the hard samples we search for improve other model performance compared to random sampling?
- **RQ4** : Does our progressive strategy provide better results compared to the fixed flip ratio strategy?

**Table 2: Statistics of datasets.**

Dataset	#Users	#Items	#Interactions	Sparsity
Adressa	212,231	6,596	419,491	0.99970
MovieLens	943	1,682	100,000	0.93695
Yelp	45,548	57,396	1,672,520	0.99936

## 4.1 Experimental Settings

**4.1.1 Datasets.** We conduct extensive comparative experiments on three public popular datasets: Adressa<sup>2</sup> [11, 32, 36], MovieLens<sup>3</sup> [11, 36], Yelp<sup>2</sup> [11, 32]. Detailed statistics of the datasets are in Table 2. For getting clean test datasets, we only include interactions with a dwell time of at least 10 seconds, based on [32, 36] in Adressa. For MovieLens, we create a test set that only includes interactions with a rating of five, following the setting from [36]. Similarly, for Yelp, we refer to [32] and only include interactions with a rating higher than three in our clean test set.

Note that to maintain consistency with the settings of previous denoising works [32, 36, 37] and ensure fair performance comparisons. All ratings and dwell times filter the clean test set, while the training and validation sets are not filtered.

**4.1.2 Evaluation protocols.** We follow [32, 37] to split datasets into the training set, validation set, and clean test set with the ratio 8:1:1. Following existing works on denoising recommendations [11, 32, 36], we report the results *w.r.t.* two widely used metrics: NDCG@K and Recall@K, where higher scores indicate better performance. For a comprehensive comparison of different models, we set  $K=5$  and  $K=20$  for all datasets. Each experiment is repeated five times, as well as we conduct a significance test on the results of the experiments. In addition, limited by the page length, we place the detailed parameter setting at A.2.

**4.1.3 Baselines.** Our goal in this paper is to weaken the adverse impact of noisy interactions on model performance. For this purpose, we choose four **backbones** based on implicit feedback:

- **GMF** [16]: This generalized version of matrix factorization captures the latent factors of users and items.
- **NeuMF** [16]: NeuMF combines matrix factorization and neural networks to enhance the accuracy of collaborative filtering.
- **NGCF** [33]: NGCF is a graph neural network model for enhanced and personalized recommendation systems.
- **LightGCN** [15]: LightGCN is a simplified GCN designed specifically for recommendation, leveraging the power of graph-based methods to enhance interaction predictions.

The following approaches are used to train each model:

- **Normal** [16]: The model is being trained using the simple binary-cross-entropy (BCE) loss function.

- **WBPR** [10]: This method considers popular but non-interactive items as true negatives based on interaction count.
- **WRMF** [18]: WRMF design weighted matrix factorization where the weights remain fixed to remove noise from recommendation.
- **T-CE** [32]: This is a generalized version of BCE to truncate large-loss examples with a dynamic threshold in each iteration.
- **DeCA** [36]: DeCA utilizes two different models to predict and account for any disagreement that may arise from noisy samples.
- **BOD** [37]: BOD automatically learns samples weights using bi-level optimization.

## 4.2 Performance Comparison (RQ1)

To validate the effectiveness and generalization of our framework, we conduct extensive experiments on four popular backbones and three datasets. Table 3 shows the results of our DCF with existing denoising methods. Indices that perform the **best** are in bold, while **runner-ups** are underlined. According to Table 3, we can draw the following observations and conclusions:

- DCF achieves good performance on all backbones and datasets. We attribute these improvements to the extended observation of sample loss values, searching for valuable hard samples, and relabeling some highly discriminated noisy samples. After relabeling, these samples can be utilized to avoid performance degradation due to sparse sample space. However, baseline models such as BOD and DeCA lack these capabilities, making them less effective than ours.
- We observe that T-CE is second only to us on most datasets and backbone, whereas the other DeCA and BOD are not as good as T-CE. T-CE directly drops the high loss samples, whereas DeCA and BOD do not utilize the high and low loss phenomena. The results are consistent with previous studies [32, 43]. Furthermore, DeCA is based on the prediction agreement of the two models, and BOD is based on bi-level optimization, and we speculate that the suboptimal performance because they are both unstable in training. Hence, we argue that using loss values is a straightforward and efficient method.
- In other baselines, the BOD is second only to us in Yelp. We think this is because the bi-level optimization of BOD requires more density data in order to better learn the weighting matrix involved. Whereas WRMF sometimes achieves good results, we speculate it is because matrix factorization has better effects on sparse data. We also note that on the MovieLens dataset, other denoising methods do not outperform normal training. We argue this is because on sparse datasets, these reweighting or dropping strategies may be less effective due to inadequate representation learning, which leads to true-positive samples being mistakenly down-weighted or dropped.

## 4.3 Model Investigation (RQ2)

**4.3.1 Ablation Study.** DCF consists of three components, Confirmed Loss Calculation (CL), Hard Sample Search (HS), and Progressive Label Correction (LC). We are eager to validate the effectiveness of each component and the combination between them on the performance. Therefore, we conduct the following seven sets of ablation experiments (as shown in Tabel 4). However, due to the length of the paper, we only show results on the MoiveLens dataset

<sup>2</sup><https://github.com/WenjieWWJ/DenoisingRec>

<sup>3</sup><https://github.com/wangyu-ustc/DeCA>



**Table 3: Performance comparison of different denoising methods on the robust recommendation. The highest scores are in bold, and the runner-ups are with underlines. R and N refer to Recall and NDCG, respectively. A significant improvement over the runner-up is marked with \* (i.e., two-sided t-test with  $p < 0.05$ ) and \*\* (i.e., two-sided t-test with  $0.05 \leq p < 0.1$ ).**

Dataset		Adressa				MovieLens				Yelp			
Base Model	Method	R@5	R@20	N@5	N@20	R@5	R@20	N@5	N@20	R@5	R@20	N@5	N@20
GMF	Normal	0.1257	0.2126	0.0908	0.1210	0.0355	0.1196	0.0482	0.0715	0.0149	0.0431	0.0152	0.0243
	WBPR	0.1258	0.2131	0.0912	0.1212	0.0357	0.1199	0.0486	0.0718	0.0148	0.0434	0.0155	0.0246
	WRMF	<u>0.1263</u>	0.2132	0.0911	0.1214	0.0363	<u>0.1201</u>	0.0491	0.0722	0.0144	0.0446	0.0140	0.0242
	T-CE	0.1251	<u>0.2155</u>	<u>0.0913</u>	<u>0.1228</u>	<u>0.0374</u>	<b>0.1202</b>	<u>0.0509</u>	<u>0.0738</u>	0.0143	<u>0.0447</u>	0.0143	0.0244
	DeCA	0.1232	0.2147	0.0866	0.1181	0.0364	0.1075	0.0423	0.0631	0.0141	0.0442	0.0139	0.0236
	BOD	0.1237	0.2153	0.0892	0.1193	0.0366	0.1083	0.0474	0.0685	<u>0.0151</u>	0.0436	<u>0.0157</u>	<u>0.0247</u>
	DCF (Ours)	<b>0.1296*</b>	<b>0.2183*</b>	<b>0.0938*</b>	<b>0.1254*</b>	<b>0.0427*</b>	0.1175	<b>0.0543*</b>	<b>0.0743*</b>	<b>0.0155**</b>	<b>0.0458*</b>	<b>0.0158</b>	<b>0.0257*</b>
NeuMF	Normal	0.1909	0.3078	0.1427	0.1851	<u>0.0439</u>	0.1084	<u>0.0516</u>	0.0724	0.0123	0.0386	<u>0.0123</u>	0.0210
	WBPR	0.1903	0.3082	<u>0.1428</u>	0.1848	<u>0.0426</u>	0.1132	<u>0.0504</u>	<u>0.0735</u>	0.0104	0.0384	<u>0.0108</u>	0.0193
	WRMF	<u>0.1922</u>	<u>0.3084</u>	0.1424	<u>0.1852</u>	0.0418	<u>0.1180</u>	0.0512	0.0729	0.0119	0.0378	0.0121	0.0198
	T-CE	0.1880	0.3080	0.1410	0.1847	0.0366	<u>0.1065</u>	0.0482	0.0680	0.0108	0.0383	0.0105	0.0190
	DeCA	0.1870	0.3076	0.1402	0.1804	<u>0.0327</u>	0.0990	0.0388	0.0590	0.0103	0.0381	0.0101	0.0182
	BOD	0.1890	0.3082	0.1414	0.1828	0.0375	0.0134	0.0489	0.0703	<u>0.0126</u>	<u>0.0389</u>	0.0119	<u>0.0215</u>
	DCF (Ours)	<b>0.1979*</b>	<b>0.3134*</b>	<b>0.1439*</b>	<b>0.1853</b>	<b>0.0513*</b>	<b>0.1210*</b>	<b>0.0642*</b>	<b>0.0816*</b>	<b>0.0132*</b>	<b>0.0411*</b>	<b>0.0132*</b>	<b>0.0223*</b>
NGCF	Normal	0.1235	0.2257	0.0934	0.1291	0.0335	<u>0.1015</u>	<u>0.0452</u>	0.0634	0.0172	0.0495	<u>0.0174</u>	0.0273
	WBPR	0.1237	0.2252	0.0936	0.1289	0.0331	0.1013	0.0446	0.0637	0.0166	0.0481	0.0164	0.0267
	WRMF	0.1244	0.2255	0.0942	0.1304	0.0334	0.1011	0.0449	0.0639	0.0169	0.0486	0.0167	0.0270
	T-CE	<u>0.1260</u>	<u>0.2270</u>	<u>0.0959</u>	<u>0.1313</u>	<u>0.0335</u>	0.1014	0.0450	0.0635	0.0173	<u>0.0497</u>	0.0173	0.0270
	DeCA	0.1172	0.2235	0.0846	0.1037	<u>0.0318</u>	0.0973	0.0436	0.0627	0.0169	0.0464	0.0166	0.0268
	BOD	0.1212	0.2246	0.0901	0.1265	0.0321	0.1008	0.0437	0.0633	<u>0.0174</u>	0.0492	0.0169	<u>0.0274</u>
	DCF (Ours)	<b>0.1267*</b>	<b>0.2275**</b>	<b>0.0970*</b>	<b>0.1321*</b>	<b>0.0353*</b>	<b>0.1037*</b>	<b>0.0468*</b>	<b>0.0647*</b>	<b>0.0180*</b>	<b>0.0503**</b>	<b>0.0179**</b>	<b>0.0277**</b>
LightGCN	Normal	0.1236	0.2257	0.0933	0.1289	0.0347	0.1029	0.0457	0.0642	0.0185	0.0514	0.0183	<u>0.0291</u>
	WBPR	0.1239	0.2253	0.0914	0.1295	0.0353	0.1043	0.0460	0.0638	0.0182	0.0514	0.0178	0.0282
	WRMF	0.1242	0.2260	0.0928	0.1298	<u>0.0357</u>	<u>0.1046</u>	<u>0.0464</u>	<u>0.0650</u>	0.0181	0.0515	0.0180	0.0283
	T-CE	<b>0.1261</b>	<u>0.2261</u>	<u>0.0956</u>	<u>0.1306</u>	0.0290	0.1020	0.0409	0.0612	<u>0.0185</u>	<u>0.0516</u>	<u>0.0183</u>	0.0290
	DeCA	0.1185	0.2251	0.0859	0.1038	0.0347	0.0987	0.0440	0.0640	0.0176	0.0503	0.0172	0.0273
	BOD	0.1225	0.2254	0.0902	0.1287	0.0325	0.1014	0.0443	0.0638	0.0182	0.0513	0.0177	0.0282
	DCF (Ours)	<u>0.1258</u>	<b>0.2274*</b>	<b>0.0961**</b>	<b>0.1315*</b>	<b>0.0365*</b>	<b>0.1050</b>	<b>0.0472*</b>	<b>0.0659*</b>	<b>0.0192*</b>	<b>0.0523*</b>	<b>0.0187**</b>	<b>0.0296**</b>

and GMF, and other datasets are similar and omitted. In addition, we also conduct ablation experiments on damping function in A.3.1.

We observe better results from more components, which aligns with our expectations. We analyze this trend in effectiveness as being attributable to each component playing its desired role. In the single component section, we observe the best results for DCF<sub>HS</sub>. This proves that we obtain hard samples by deriving a confidence lower bound on the loss. In addition, DCF<sub>LC</sub> also achieves good performance, which is attributed to the fact that we let the corrected samples be retrained to avoid sample wasting, as well as misalignment of the training and testing sample spaces. Also, in the combination of components, HS and LC achieved excellent results, proving the effectiveness of the two-by-two combination. It is worth noting that the combination of CL and HS does not achieve more significant results. We hypothesize that the reason for this is that the search for HS is a game of risk and benefit, and therefore, we have to tune the hyperparameters carefully.

**4.3.2 Parameters Sensitivity Analysis.** We are eager to know the sensitivity of different critical parameters: (i) relabeling ratio  $R$ ; (ii) discretion level  $\sigma^2$  of searching hard samples; (iii) calculating time interval  $v$  of mean losses. Experiments are conducted on the MovieLens dataset. As shown in Fig. 4, we can find that:

**Table 4: The effect of each components on DCF.**

Method	R@5	R@10	N@5	N@10
T-CE	0.0374	0.0734	0.0509	0.0591
DCF <sub>CL</sub>	0.0432(15.5%)	0.0751(2.3%)	0.0527(3.5%)	0.0600(1.5%)
DCF <sub>HS</sub>	0.0435(16.3%)	0.0772(5.2%)	0.0535(5.1%)	0.0610(3.2%)
DCF <sub>LC</sub>	0.0430(15.0%)	0.0756(3.0%)	0.0529(3.9%)	0.0604(2.2%)
DCF <sub>CL+HS</sub>	0.0436(16.6%)	0.0743(1.2%)	0.0531(4.3%)	0.0596(0.8%)
DCF <sub>CL+LC</sub>	0.0459(22.7%)	0.0763(4.0%)	0.0549(7.9%)	0.0612(3.6%)
DCF <sub>HS+LC</sub>	0.0454(21.4%)	0.0779(6.1%)	0.0545(7.1%)	0.0613(3.7%)
DCF <sub>ALL</sub>	<b>0.0471(25.9%)</b>	<b>0.0789(7.5%)</b>	<b>0.0553(8.6%)</b>	<b>0.0621(5.1%)</b>

- NDCG@5 and Recall@5 are affected by the hyperparameters, and the overall trend is up and then down. We observe that the performance is not getting better as the flip rate  $R$  increases. We speculate that this is because only a tiny fraction of the highest loss samples are strongly correlated with the noisy samples, and thus more flips may lead to more errors.
- $\sigma^2$  controls the intensity of the search for hard samples, which is a risk versus reward game; the larger the  $\sigma^2$ , the higher the

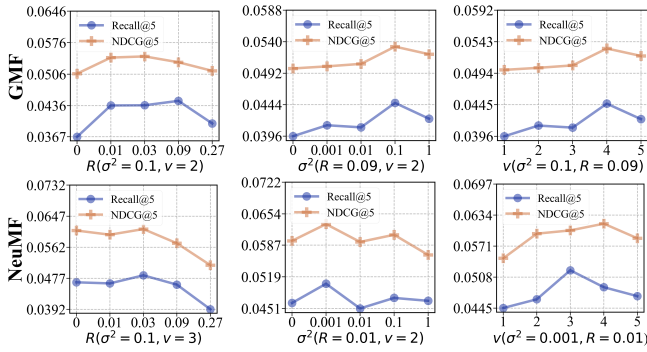


Figure 4: Impact of the relabel ratio  $R$ , search discretion level  $\sigma^2$  and time interval  $v$ .

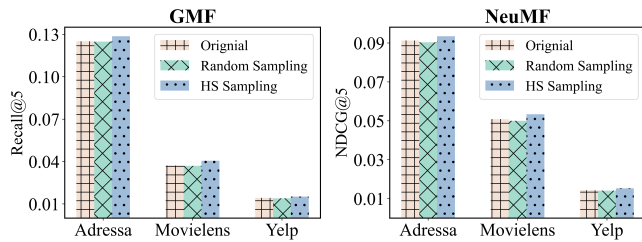


Figure 5: Comparative experiments on three datasets with two backbones validate the effectiveness of our hard sample search strategy to improve model performance.

probability of performance degradation, but at the same time it is possible to improve the effectiveness of the model.

- We observe that extending the time interval  $v$  for mean loss computation has apparent benefits, but this is limited to a specific range. If the range is too wide, the calculation of the mean may be negatively affected by the instability of early model training.

#### 4.4 Hard Sample Search Strategy (RQ3)

Despite confirming the effectiveness of hard sample search in our ablation experiments, we aim to ascertain whether the identified samples are indeed hard samples rather than noisy samples. Therefore, we design a comparative experiment. The first group is to use the searched samples for T-CE training. The second group is to randomly sample the same size of samples from the dropped samples for training, and the third group is not to use the original training of both samples to compare the performance of these three. Experiments are conducted on three datasets and two backbones.

From Fig. 5, we observe that adding hard samples improves model performance, which is consistent with our expectations. This is due to the efficient identification of noisy samples through our lower bound on loss. In contrast, adding random samples from the dropped samples to the model does not enhance the recommendation model performance but has a negative impact.

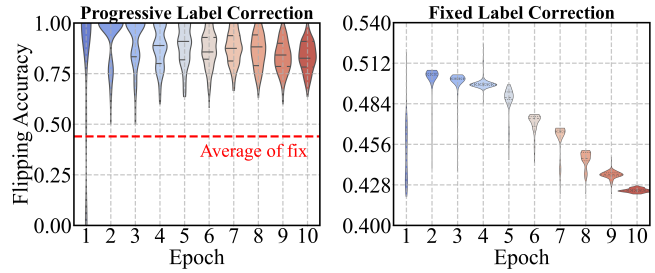


Figure 6: Comparison of flip accuracy between progressive label correction and fixed. For a clear presentation, we use a violin plot here. Additionally, we mark the average flip accuracy of fixed with a red line to clearly highlight the superiority of our progressive strategy.

#### 4.5 Progressive Strategy in Label Correction (RQ4)

We have demonstrated the effectiveness of progressive label correction in the ablation section. However, we are interested in gaining a more detailed understanding of our strategy’s performance, particularly regarding flip accuracy. We also wish to determine if our progressive relabeling approach is superior to fixed relabeling methods. The results of the experiment are shown in Fig. 6.

- From the two figures, it can be observed that the flip accuracy of progressive label correction is significantly higher than that of the fixed strategy (red lines). We attribute this to our dynamic relabeling ratio strategy, better suited for transitioning from instability to stability training characteristics. In the early stages of training, the model may be more susceptible to the interference of noisy labels, which could impact its final performance. However, in the later stages, the accuracy of the progressive relabeling decreases slightly. We presume that this is because the model already has enough information, and there is a slight overfitting phenomenon.
- We also observe that the stability of the flip accuracy in the progressive strategy increases as training iterations progress and ultimately maintain a high level. Additionally, in the beginning, even if we flip a small proportion, the volatility is still high. If we were to flip a larger proportion, there is a high likelihood of introducing incorrect information. This further validates the effectiveness of our strategy of gradually increasing the flip ratio from low to high.

### 5 RELATED WORK

Recommendation systems based on implicit feedback have attracted a lot of attention. However, recent studies point out implicit feedback is easily vulnerable to users’ unconscious behaviors and various biases (e.g., popularity bias, position bias, etc.), which degrade the generalization ability. Thus to weaken the problem caused by noisy implicit feedback, some denoising methods [9, 18, 23, 49, 52] have been proposed, and they can be categorized into sample drop [9, 17, 20, 21, 30, 32, 38, 51], sample reweight [10, 12, 32, 37], and other methods [36] designed with the help of other information.



**Drop-based Methods.** An intuitive idea is to pick out the clean interactions and send them for training. Then, distinguishing between the characteristics of clean and noisy samples is critical. T-CE [32] observe experimentally loss values of noisy samples are higher than clean samples, and experiments using this observation find it effective. Subsequent studies [11, 21] incorporate this observation into the design of denoising recommendation models. For example, AutoDenoise [21] considers the denoising work to consist of two behaviors: searching and deciding. So they use reinforcement learning to automate these two behaviors. Additionally, [51] proposes to augment recommendation algorithms with noisy examples of user preferences and mitigate the challenge of data sparsity.

**Reweight-based Methods.** In addition to selecting clean samples, some work [12, 32, 37] reduce the weights on the noisy samples, which reduces the impact on model parameter updates and avoids the reduction of generalization ability. For example, R-CE [32] utilizes the loss value as a denoising signal, thus assigning a small weight to the noisy sample. AutoDenoise [12] also utilizes the loss value as a denoising signal and fuses it with the user representation and item representation in the weight calculation with good results. BOD [37] converts the learning of weights into a bi-level optimization problem and uses a simple and elegant solution to learn the weight parameters with relatively good results.

**Other Methods.** First, methods inspired by the phenomenon of noisy samples during training. For example, DeCA [36] observes that different models have greater predictive divergence for noisy samples and lesser divergence for clean samples. SGDL [11] observe the memorization effect of noise, thus removing noise in the pre-training period. Thus based on this observation, an efficient denoising recommendation model is designed and performs well. Second, there is also some work to denoise based on graph collaborative filtering, e.g., RGCF [29], RocSE [46]. Their method's main design idea lies in removing noisy interaction edges. RocSE [46] removes the noisy interactions from both graph structure denoising and contrastive learning and performs well on graph collaborative filtering. Third, some methods to denoise in specific recommendation scenarios, design denoising modules for scenarios such as movies [50], music [7], and social [31, 44]. As well as exploring the removal of noisy samples in sequential recommendation [5, 28, 47].

**Comparison with Our Method.** We argue that the loss-based method is a simple and effective approach but suffers from predictive instability due to random initialization of parameters. Furthermore, both the dropped sample and reweighted sample strategies based on loss values design have inherent flaws. For example, it leads to misalignment of the training and testing spaces or the need for complex methods to learn the weights. Therefore, we consider a more robust loss value calculation as well as relabeling.

## 6 CONCLUSION

In this paper, we present a novel framework named DCF, designed to mitigate the adverse effects of noisy samples on the representation learning of users and items. The DCF contains two modules. The first module, sample dropping correction, achieves more stable loss estimations by calculating the mean loss value of samples

and focuses on identifying and retraining hard samples. The second module, progressive label correction, we relabel samples with a higher likelihood of being noisy and reintegrate them into the learning process. Extensive experiments on widely used benchmarks and datasets demonstrate the effectiveness and generalization of our proposed framework.

## ACKNOWLEDGMENTS

This work is supported by grants from the National Key Research and Development Program of China (Grant No.2021ZD0111802) and National Science Foundation of China (Grant No.U23B2031, No.721881011, No.U21B2026).

## REFERENCES

- [1] Haoyue Bai, Min Hou, Le Wu, Yonghui Yang, Kun Zhang, Richang Hong, and Meng Wang. 2023. GoRec: A Generative Cold-start Recommendation Framework. (2023).
- [2] Haoyue Bai, Min Hou, Le Wu, Yonghui Yang, Kun Zhang, Richang Hong, and Meng Wang. 2024. Unified Representation Learning for Discrete Attribute Enhanced Completely Cold-Start Recommendation. *IEEE Transactions on Big Data* (2024), 1–12.
- [3] Miaomiao Cai, Min Hou, Lei Chen, Le Wu, Haoyue Bai, Yong Li, and Meng Wang. 2024. Mitigating Recommendation Biases via Group-Alignment and Global-Uniformity in Representation Learning. *ACM Transactions on Intelligent Systems and Technology* (2024).
- [4] Olivier Catoni. 2012. Challenging the empirical mean and empirical variance: a deviation study. In *Annales de l'IHP Probabilités et statistiques*, Vol. 48. 1148–1185.
- [5] Huiyuan Chen, Yusan Lin, Menghai Pan, Lan Wang, Chin-Chia Michael Yeh, Xiaoting Li, Yan Zheng, Fei Wang, and Hao Yang. 2022. Denoising self-attentive sequential recommendation. In *Proceedings of the 16th ACM Conference on Recommender Systems*. 92–101.
- [6] Peng Chen, Xinghu Jin, Xiang Li, and Lihu Xu. 2021. A generalized catoni's m-estimator under finite  $\alpha$ -th moment assumption with  $\alpha \in (1, 2)$ . *Electronic Journal of Statistics* 15, 2 (2021), 5523–5544.
- [7] Quanyu Dai, Yalei Lv, Jieming Zhu, Junjie Ye, Zhenhua Dong, Rui Zhang, Shu-Tao Xia, and Ruiming Tang. 2022. LCD: Adaptive Label Correction for Denoising Music Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3903–3907.
- [8] Jingtao Ding, Yuhao Quan, Quanming Yao, Yong Li, and Depeng Jin. 2020. Simplify and robustify negative sampling for implicit collaborative filtering. *Advances in Neural Information Processing Systems* 33 (2020), 1094–1105.
- [9] Jingtao Ding, Guanghui Yu, Xiangnan He, Fuli Feng, Yong Li, and Depeng Jin. 2019. Sampler design for bayesian personalized ranking by leveraging view data. *IEEE transactions on knowledge and data engineering* 33, 2 (2019), 667–681.
- [10] Zeno Gantner, Lucas Drumond, Christoph Freudenthaler, and Lars Schmidt-Thieme. 2012. Personalized ranking for non-uniformly sampled items. In *Proceedings of KDD Cup 2011*. PMLR, 231–247.
- [11] Yunjun Gao, Yuntao Du, Yujia Hu, Lu Chen, Xinjun Zhu, Ziquan Fang, and Baihua Zheng. 2022. Self-guided learning to denoise for robust recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1412–1422.
- [12] Yingqiang Ge, Mostafa Rahmani, Athirai Irissappane, Jose Sepulveda, Fei Wang, James Caverlee, and Yongfeng Zhang. 2023. Automated Data Denoising for Recommendation. *arXiv preprint arXiv:2305.07070* (2023).
- [13] Justin Gilmer, Behrooz Ghorbani, Ankush Garg, Sneha Kudugunta, Behnam Neyshabur, David Cardoze, George Dahl, Zachary Nado, and Orhan Firat. 2021. A loss curvature perspective on training instability in deep learning. *arXiv preprint arXiv:2110.04369* (2021).
- [14] Evarist Giné, Rafal Latała, and Joel Zinn. 2000. Exponential and moment inequalities for U-statistics. In *High Dimensional Probability II*. Springer, 13–38.
- [15] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.
- [16] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.
- [17] Cho-Jui Hsieh, Nagarajan Natarajan, and Inderjit Dhillon. 2015. PU learning for matrix completion. In *International conference on machine learning*. PMLR, 2445–2453.
- [18] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *2008 Eighth IEEE international conference on data*

- mining. *Ieee*, 263–272.
- [19] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436–444.
- [20] Dawen Liang, Laurent Charlin, James McInerney, and David M Blei. 2016. Modeling user exposure in recommendation. In *Proceedings of the 25th international conference on World Wide Web*. 951–961.
- [21] Weilin Lin, Xiangyu Zhao, Yejing Wang, Yuanshao Zhu, and Wanyu Wang. 2023. AutoDenoise: Automatic Data Instance Denoising for Recommendations. In *Proceedings of the ACM Web Conference 2023*. 1003–1011.
- [22] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. 2012. Foundations of Machine Learning. (Aug 2012).
- [23] Rong Pan, Yunhong Zhou, Bin Cao, Nathan N Liu, Rajan Lukose, Martin Scholz, and Qiang Yang. 2008. One-class collaborative filtering. In *2008 Eighth IEEE international conference on data mining*. IEEE, 502–511.
- [24] Wentao Shi, Jiawei Chen, Fuli Feng, Jizhi Zhang, Junkang Wu, Chongming Gao, and Xiangnan He. 2023. On the Theories Behind Hard Negative Sampling for Recommendation. In *Proceedings of the ACM Web Conference 2023*. 812–822.
- [25] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. 2016. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 761–769.
- [26] Peijie Sun, Yifan Wang, Min Zhang, Chuhan Wu, Yan Fang, Hong Zhu, Yuan Fang, and Meng Wang. 2024. Collaborative-Enhanced Prediction of Spending on Newly Downloaded Mobile Games under Consumption Uncertainty. *WWW2024, Industry Track* (2024).
- [27] Peijie Sun, Le Wu, Kun Zhang, Xiangzhi Chen, and Meng Wang. 2023. Neighborhood-Enhanced Supervised Contrastive Learning for Collaborative Filtering. *IEEE Transactions on Knowledge and Data Engineering* (2023).
- [28] Yatong Sun, Bin Wang, Zhu Sun, and Xiaochun Yang. 2021. Does Every Data Instance Matter? Enhancing Sequential Recommendation by Eliminating Unreliable Data. In *IJCAI*. 1579–1585.
- [29] Changxin Tian, Yuexiang Xie, Yaliang Li, Nan Yang, and Wayne Xin Zhao. 2022. Learning to denoise unreliable interactions for graph collaborative filtering. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 122–132.
- [30] Menghan Wang, Yujie Lin, Guli Lin, Keping Yang, and Xiao-ming Wu. 2020. M2GRL: A multi-task multi-view graph representation learning framework for web-scale recommender systems. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 2349–2358.
- [31] Tianle Wang, Lianghao Xia, and Chao Huang. 2023. Denoised Self-Augmented Learning for Social Recommendation. *arXiv preprint arXiv:2305.12685* (2023).
- [32] Wenjie Wang, Fuli Feng, Xiangnan He, Liqiang Nie, and Tat-Seng Chua. 2021. Denoising implicit feedback for recommendation. In *Proceedings of the 14th ACM international conference on web search and data mining*. 373–381.
- [33] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*. 165–174.
- [34] Yifan Wang, Peijie Sun, Weizhi Ma, Min Zhang, Yuan Zhang, Peng Jiang, and Shaoping Ma. 2024. Intersectional Two-sided Fairness in Recommendation. In *Proceedings of the ACM on Web Conference 2024 (WWW '24)*. Association for Computing Machinery, New York, NY, USA, 3609–3620. <https://doi.org/10.1145/3589334.3645518>
- [35] Yifan Wang, Peijie Sun, Min Zhang, Qinglin Jia, Jingjie Li, and Shaoping Ma. 2023. Unbiased Delayed Feedback Label Correction for Conversion Rate Prediction. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2456–2466.
- [36] Yu Wang, Xin Xin, Zaiqiao Meng, Joemon M Jose, Fuli Feng, and Xiangnan He. 2022. Learning robust recommenders through cross-model agreement. In *Proceedings of the ACM Web Conference 2022*. 2015–2025.
- [37] Zongwei Wang, Min Gao, Wentao Li, Junliang Yu, Linxin Guo, and Hongzhi Yin. 2023. Efficient Bi-Level Optimization for Recommendation Denoising. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2502–2511.
- [38] Zitai Wang, Qianqian Xu, Zhiyong Yang, Xiaochun Cao, and Qingming Huang. 2021. Implicit feedbacks are not always favorable: Iterative relabeled one-class collaborative filtering against noisy interactions. In *Proceedings of the 29th ACM International Conference on Multimedia*. 3070–3078.
- [39] Hongyi Wen, Longqi Yang, and Deborah Estrin. 2019. Leveraging post-click feedback for content recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 278–286.
- [40] Le Wu, Peijie Sun, Yanjie Fu, Richang Hong, Xiting Wang, and Meng Wang. 2019. A neural influence diffusion model for social recommendation. In *Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval*. 235–244.
- [41] Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama. 2021. Sample selection with uncertainty of losses for learning with noisy labels. *arXiv preprint arXiv:2106.00445* (2021).
- [42] Junyuan Xie, Ross Girshick, and Ali Farhadi. 2016. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*. PMLR, 478–487.
- [43] Xin Xin, Xiangyuan Liu, Hanbing Wang, Pengjie Ren, Zhumin Chen, Jiahuan Lei, Xinlei Shi, Hengliang Luo, Joemon M. Jose, Maarten de Rijke, and Zhaochun Ren. 2023. Improving Implicit Feedback-Based Recommendation through Multi-Behavior Alignment. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 932–941.
- [44] Dezhao Yang, Jianghong Ma, Shanshan Feng, Haijun Zhang, and Zhao Zhang. 2023. IDVT: Interest-aware Denoising and View-guided Tuning for Social Recommendation. *arXiv preprint arXiv:2308.15926* (2023).
- [45] Yonghui Yang, Le Wu, Kun Zhang, Richang Hong, Hailin Zhou, Zhiqiang Zhang, Jun Zhou, and Meng Wang. 2023. Hyperbolic Graph Learning for Social Recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2023), 1–14. <https://doi.org/10.1109/TKDE.2023.3343402>
- [46] Haibo Ye, Xinjie Li, Yuan Yao, and Hanghang Tong. 2023. Towards robust neural graph collaborative filtering via structure denoising and embedding perturbation. *ACM Transactions on Information Systems* 41, 3 (2023), 1–28.
- [47] Chi Zhang, Yantong Du, Xiangyu Zhao, Qilong Han, Rui Chen, and Li Li. 2022. Hierarchical item inconsistency signal learning for sequence denoising in sequential recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2508–2518.
- [48] Honglei Zhang, Fangyuan Luo, Jun Wu, Xiangnan He, and Yidong Li. 2023. LightFR: Lightweight Federated Recommendation with Privacy-preserving Matrix Factorization. *ACM Trans. Inf. Syst.* 41, 4, Article 90 (mar 2023), 28 pages.
- [49] Kaike Zhang, Qi Cao, Fei Sun, Yunfan Wu, Shuchang Tao, Huawei Shen, and Xueqi Cheng. 2023. Robust Recommender System: A Survey and Future Directions. *arXiv preprint arXiv:2309.02057* (2023).
- [50] Haiyuan Zhao, Lei Zhang, Jun Xu, Guohao Cai, Zhenhua Dong, and Ji-Rong Wen. 2023. Uncovering User Interest from Biased and Noised Watch Time in Video Recommendation. *arXiv preprint arXiv:2308.08120* (2023).
- [51] Xiangyu Zhao, Zhendong Niu, Kaiyi Wang, Ke Niu, Zhongqiang Liu, et al. 2015. Improving top-N recommendation performance using missing data. *Mathematical Problems in Engineering* 2015 (2015).
- [52] Xinjun Zhu, Yuntao Du, Yuren Mao, Lu Chen, Yujia Hu, and Yunjun Gao. 2023. Knowledge-refined Denoising Network for Robust Recommendation. *arXiv preprint arXiv:2304.14987* (2023).

## A APPENDIX

### A.1 Derivation OF THEOREM 1

According to [4], we explore the upper bound  $\tilde{\mu}_z^-$  and the lower bound  $\tilde{\mu}_z^+$  for the mean  $\tilde{\mu}_z = \frac{1}{n} \sum_{i=1}^n \phi(z_i)$  computation, as follows

$$\tilde{\mu}_z^- \leq \tilde{\mu}_z \leq \tilde{\mu}_z^+. \quad (8)$$

To derive a bound for  $\tilde{\mu}_z$ . We build an estimator of  $\tilde{\mu}_z$ , the formal definition as follow,

$$r(\tilde{\mu}_z) = \sum_{i=1}^n \phi[\alpha(z_i - \tilde{\mu}_z)] = 0. \quad (9)$$

where  $\alpha$  is positive real parameter. We then adjust  $r(\tilde{\mu}_z)$  to be in the form of quantity

$$r(\theta) = \frac{1}{\alpha n} \sum_{i=1}^n \phi[\alpha(z_i - \theta)], \theta \in \mathbb{R}. \quad (10)$$

With the exponential moment inequality [14] and the  $C_r$  inequality [22], we analysis the  $r(\tilde{\mu}_z)$  in the form of inequalities

$$\begin{aligned} \mathbb{E}\{\exp[\alpha n r(\theta)]\} &\leq \left\{1 + \alpha(\mu_z - \theta) + \frac{\alpha^2}{2} [\alpha^2 + (\mu_z - \theta)^2]\right\}^n \\ &\leq \exp\left\{n\alpha(\mu_z - \theta) + \frac{n\alpha^2}{2} [\alpha^2 + (\mu_z - \theta)^2]\right\}. \end{aligned} \quad (11)$$

In order to bound  $\mu_z$ , we will find two non-random values  $\theta_-$  and  $\theta_+$  of the parameter such that with large probability  $r(\theta_-) > 0 > r(\theta_+)$ , which will imply that  $\theta_- < \hat{\theta}_\alpha < \theta_+$ , since  $r(\hat{\theta}_\alpha) = 0$  by construction and  $\theta \mapsto r(\theta)$  is non-increasing. The new bounds are as follow

$$B_-(\theta) = \mu_z - \theta - \alpha \left[ \sigma^2 + (\mu_z - \theta)^2 \right] - \frac{\log(\epsilon^{-1})}{\alpha n} \quad (12)$$

and

$$B_+(\theta) = \mu_z - \theta + \alpha \left[ \sigma^2 + (\mu_z - \theta)^2 \right] + \frac{\log(\epsilon^{-1})}{\alpha n}. \quad (13)$$

From [6] (Lemma 2.2), we obtain that

$$P(r(\theta) > B_-(\theta)) \geq 1 - \epsilon \quad \text{and} \quad P(r(\theta) < B_+(\theta)) \geq 1 - \epsilon. \quad (14)$$

According to Chebyshev's inequality and the previous proposition [6]. We assume

$$4\alpha^2\sigma^2 + \frac{4\log(\epsilon^{-1})}{n} \leq 1. \quad (15)$$

From [6] (Theorem 2.6), we then have

$$\tilde{\mu}_z^- \geq \mu_z - \frac{\alpha\sigma^2 + \frac{\log(\epsilon^{-1})}{\alpha n}}{\alpha - 1} \quad (16)$$

and

$$\tilde{\mu}_z^+ \leq \mu_z + \frac{\alpha\sigma^2 + \frac{\log(\epsilon^{-1})}{\alpha n}}{\alpha - 1}. \quad (17)$$

With probability at least  $1 - 2\epsilon$ , we have  $\tilde{\mu}_z^- \leq \tilde{\mu}_z \leq \tilde{\mu}_z^+$ . We set  $\alpha = \frac{n}{\sigma^2}$ , and obtain

$$|\tilde{\mu}_z - \mu_z| \leq \frac{\sigma^2 \left( n + \frac{\sigma^2 \log(\epsilon^{-1})}{n^2} \right)}{n - \sigma^2}. \quad (18)$$

With probability of at least  $1 - 2\epsilon$ , we proceed to exploit the lower bound by setting  $\epsilon = \frac{1}{2i}$ .

$$t^* = \tilde{\mu} - \frac{\sigma^2 \left( i + \frac{\sigma^2 \log(2i)}{i^2} \right)}{d - \sigma^2},$$

where  $d$  represents the number of times a sample has not been dropped.

### A.2 Parameter Settings

All of the models use the Adam optimizer with a learning rate of 0.001, batch size set to 1024, and embedding dimension set to 32. The number of graph convolution layers for LightGCN is set to 3 without dropout. For every training example, we select one observed interaction and one random unobserved interaction to input into the model. The ratio of relabeled interactions is tuned in  $\{0, 0.01, 0.03, 0.09, 0.27\}$ .  $\sigma^2$  is tuned in  $\{0, 0.001, 0.01, 0.1\}$ . The time step is tuned in  $\{1, 2, 3, 4, 5\}$ . Note that the purpose of using hyper-parameters  $\sigma^2$  is to lessen the dependence on strong assumptions, helping our framework work better in practice. Our experiments are based on RTX 2080 Ti GPU and PyTorch.

### A.3 Additional Experiments

*A.3.1 Effect of Damping Function.* We verify whether damping function (DF) is effective through comparative experiments. By observing Table 5, we conclude that DF produces a positive effect. And we reason this is because DF removes the negative impact of outliers, thus making confirmed loss calculation more accurate.

**Table 5: The effect of damping function with GMF on DCF.**

Dataset	Method	R@5	R@10	N@5	N@10
Adressa	DCF	0.1296	0.2183	0.0938	0.1254
	DCF w/o DF	0.1292	0.2178	0.0933	0.1251
MovieLens	DCF	0.0427	0.1175	0.0543	0.0743
	DCF w/o DF	0.0423	0.1172	0.0540	0.0739
Yelp	DCF	0.0155	0.0458	0.0158	0.0257
	DCF w/o DF	0.0151	0.0453	0.0154	0.0252