# Efficient Bi-Level Optimization for Recommendation Denoising

Zongwei Wang
zongwei@cqu.edu.cn
Chongqing University
China

Min Gao*
gaomin@cqu.edu.cn
Chongqing University
China

Wentao Li*
wentaoli@hkust-gz.edu.cn
The Hong Kong University of Science
and Technology (Guangzhou)
China

Junliang Yu
jl.yu@uq.edu.au
The University of Queensland
Australia

Linxin Guo
guolinxin@cqu.edu.cn
Chongqing University
China

Hongzhi Yin
h.yin1@uq.edu.au
The University of Queensland
Australia

## ABSTRACT

The acquisition of explicit user feedback (e.g., ratings) in real-world recommender systems is often hindered by the need for active user involvement. To mitigate this issue, implicit feedback (e.g., clicks) generated during user browsing is exploited as a viable substitute. However, implicit feedback possesses a high degree of noise, which significantly undermines recommendation quality. While many methods have been proposed to address this issue by assigning varying weights to implicit feedback, two shortcomings persist: (1) the weight calculation in these methods is iteration-independent, without considering the influence of weights in previous iterations, and (2) the weight calculation often relies on prior knowledge, which may not always be readily available or universally applicable.

To overcome these two limitations, we model recommendation denoising as a bi-level optimization problem. The inner optimization aims to derive an effective model for the recommendation, as well as guiding the weight determination, thereby eliminating the need for prior knowledge. The outer optimization leverages gradients of the inner optimization and adjusts the weights in a manner considering the impact of previous weights. To efficiently solve this bi-level optimization problem, we employ a weight generator to avoid the storage of weights and a one-step gradient-matching-based loss to significantly reduce computational time. The experimental results on three benchmark datasets demonstrate that our proposed approach outperforms both state-of-the-art general and denoising recommendation models. The code is available at https://github.com/CoderWZW/BOD.

## CCS CONCEPTS

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

*Corresponding author

## KEYWORDS

recommendation, denoising, bi-level optimization, implicit feedback

## 1 INTRODUCTION

Recommender systems [9, 20, 28, 38], which leverage historical behavioral data of users to discover their latent interests, have been highly successful in domains such as E-commerce [40], for improving user experience and driving incremental revenue. While the explicit user feedback (e.g., ratings) is the best fuel for recommender systems, its acquisition is often impeded by the need for active user participation. Hence, implicit feedback (e.g., clicks) generated during user browsing is exploited as a viable substitute [7, 27]. For implicit feedback, an observed user-item interaction is generally regarded as a *positive sample*, while an unobserved interaction is deemed as a *negative sample* [39]. However, implicit feedback is plagued by a significant level of noise [3, 33, 35] as evidenced by the following factors: (1) users may inadvertently interact with certain items, giving rise to *false positive samples* [1, 22] due to curiosity or misclicks; (2) users may encounter situations where they lack exposure to an item, leading to *false negative samples* [7], despite having a positive preference for it. The presence of such noise exacerbates the overestimation and underestimation of some certain user preferences. Hence, it is imperative to mitigate this noise in order to enhance the accuracy of recommendations [33].

A prevalent approach to mitigating the impact of noise in implicit feedback involves utilizing auxiliary information such as attribute features [46] and external data (e.g., social networks) [40]. However, the acquisition of such information may be hindered by privacy concerns [30]. In the absence of auxiliary information, current denoising methods often employ a weighting strategy in which the samples are assigned varying weights iteratively [30]. Specifically, this is achieved through the repetitive execution of two phases: first training the recommendation model using the initial weight assignments, and then recomputing the sample weights based on the output of the recommendation model. The weights reflect the level of noise in each sample, as well as their contribution to the recommendation [7], thus enabling effective denoising.

Despite the demonstrated efficacy, current research still faces two limitations: (1) the weight assignment strategy is performed in an iteration-independent manner [30], leading to the random re-initialization of weights. This implies that the previously computed weights, which possess the ability to express the confidence level regarding the presence of noise in the samples, are disregarded during the calculation of sample weights in the current iteration. As a result, the initialization weights in every iteration is mere blind guesses, while the valuable information generated during previous iterations is not utilized to its full potential; (2) existing methods often resort to *prior knowledge* for determining sample weights. For example, Wang *et al.* [30] claim that a sample with high loss is more likely to be noisy, whereas [32] argues that a noisy sample varies greatly in loss across multiple recommendation models. Nevertheless, prior knowledge may not be readily available, and its applicability may vary in different situations.

To address the limitations inherent in current methods, we model recommendation denoising as a bi-level optimization problem, which consists of two interdependent optimization sub-tasks: an inner optimization for deriving an effective recommendation model and an outer optimization for learning the weights while considering the impact of previous weights. The logic behind the bi-level optimization can be succinctly explained by two considerations. Firstly, by retaining the weights learned by the outer optimization as global variables, previous adjustments to the weights can be persistent, enabling the sharing of weight information throughout model training. Secondly, it inherently unveils the fundamental goal of recommendation denoising: augmenting recommendation accuracy to guide the denoising procedure. The inner optimization phase primarily focuses on enhancing the accuracy of the recommendation model, while the outer optimization phase completes the denoising process guided by the accuracy of the inner recommendation model. Consequently, this approach primarily emphasizes the improvement of accuracy and thus releases the dependence on prior knowledge.

While the bi-level optimization approach to denoising has shown promise in addressing the limitations of existing methods, it has introduced new challenges as well. (1) It is storage-intensive because the size of the weights equals to the product of the user number and the item number in recommender systems if all samples are considered. As recommender systems expand, the number of users/items can become substantial, leading to a requirement for a large amount of memory to store the weights. (2) It is time-consuming to solve the bi-level optimization problem as conventional solutions require full training of the inner optimization to produce suitable gradients for the outer optimization. To tackle these challenges, we provide an efficient <u>bi</u>-level <u>o</u>ptimization framework for <u>d</u>enoising (**BOD**). For alleviating the storage demand, BOD employs the autoencoder network [18] based generator to generate the weights on the fly. This approach helps circumvent the high memory cost by storing the generator, which possesses fewer parameters in comparison to the weight variables. For reducing the computational time, BOD adopts a one-step update strategy to optimize the outer task instead of accumulating gradients from the inner optimization. Particularly, we employ the gradient matching technique [16] to stabilize the training and ensure the validity of the gradient information in the one-step update.

Our contributions are summarized as follows.
- We propose a bi-level optimization framework to tackle the noise in implicit feedback, which is prior knowledge-free and fully utilizes the varying sample weights throughout model training.
- We provide an efficient solution for the proposed bi-level optimization framework, which significantly reduces the storage demand and computational time.
- The experimental results on three benchmark datasets demonstrate that our proposed approach outperforms both state-of-the-art general and denoising recommendation models.

The rest of this paper is structured as follows. Section 2 provides the background on the relevant preliminaries. Section 3 presents the bi-level optimization framework for recommendation denoising. Section 4 provides the efficient solution to the proposed bi-level optimization problem. The experimental results and analysis are presented in Section 5. Section 6 reviews the related work for recommendation denoising. The paper is concluded in Section 7.

## 2 PRELIMINARY

### 2.1 Implicit Feedback Based Recommendation

Given a user-item interaction data $\mathcal{D} = \{u, i, r_{u,i} \mid u \in \mathcal{U}, i \in \mathcal{I}\}$, where $\mathcal{U} \in \mathbb{R}^{|\mathcal{U}|}$ denotes the set of users, $\mathcal{I} \in \mathbb{R}^{|\mathcal{I}|}$ denotes the set of items, and $r_{u,i} = \{0, 1\}$ indicates whether user $u$ has interacted with item $i$. In general, recommendation methods based on implicit feedback are trained on interaction data $\mathcal{D}$, to learn user representations $\mathbf{Z}_{\mathcal{U}} \in \mathbb{R}^{|\mathcal{U}| \times d}$ ($d$ is the dimension of representations), item representations $\mathbf{Z}_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}| \times d}$ and a model $f$ with parameters $\theta_f$ to make predictions.

The training of recommendation model is formulated as follows:

$$\theta_f^* = \underset{\theta_f}{\arg\min} \, \mathcal{L}_{rec}(\mathcal{D}), \tag{1}$$

where $\mathcal{L}_{rec}$ is the recommendation loss, such as BPR loss [22, 23, 48] and AU loss [26], and $\theta_f^*$ is the optimal parameters of $f$. Here, we use the BPR loss as an instantiation of $L_{rec}$:

$$\mathcal{L}_{rec} = \underset{(u,i,j) \sim P_{\mathcal{D}}}{\mathbb{E}} log(\sigma(f(\mathbf{z}_u)^T f(\mathbf{z}_i) - f(\mathbf{z}_u)^T f(\mathbf{z}_j))), \tag{2}$$

where $P_{\mathcal{D}}(\cdot)$ refers to the distribution defined on the interaction data, the tuple $(u, i, j)$ denotes user $u$, a positive sample item $i$ with observed interactions, and a negative sample item $j$ without observed interactions with $u$. This triple is obtained through the pairwise sampling of positive and negative samples of $u \in \mathcal{U}$ following $P_{\mathcal{D}}(\cdot)$, and $\sigma$ is the sigmoid function.

### 2.2 Recommendation Denoising

Implicit feedback possesses a high degree of noise. The common way to denoise implicit feedback is to search a weight set $\mathcal{W} = \{u, i, w_{u,i} \mid u \in \mathcal{U}, i \in \mathcal{I}\}$, where $0 \leq w_{u,i} \leq 1$ indicates the weight for a sample (corresponds to the interaction between a user $u$ and an item $i$), which can reflect the probability of a sample being truly positive. If the weights are known, the recommendation model is trained by considering the weights on the samples. Taking the weighted BRP loss as an instance:

$$\mathcal{L}_{rec} = \underset{(u,i,j) \sim P_{\mathcal{D}}}{\mathbb{E}} log(\sigma(w_{u,i} f(\mathbf{z}_u)^T f(\mathbf{z}_i) \\ - w_{u,j} f(\mathbf{z}_u)^T f(\mathbf{z}_j))). \tag{3}$$
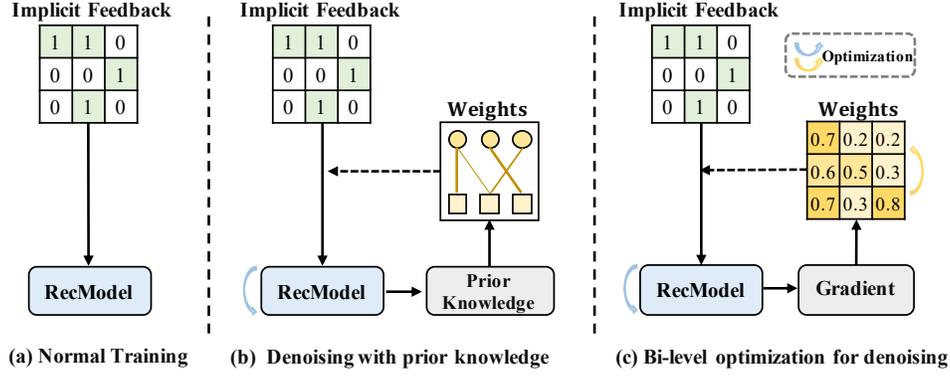
**Figure 1: The comparison among (a) Normal recommendation model training without denoising, (b) Denoising training using prior knowledge, and (c) Bi-level optimization for denoising. (b) and (c) are for recommendation denoising without extra information, where the weights in (b) are initialized for recommendation model training in each iteration, but the weights in (c) are stored as a global variable set.**

Existing denoising methods typically follow an iterative process of training a recommendation model with current weights and Equation 3, followed by weight updates based on the recommendation model's loss. As discussed in Section 1, this process has two limitations: the neglect of previous weight information and the reliance on prior knowledge to identify potential noisy samples. To address these limitations, we introduce a novel denoising method based on bi-level optimization.

## 3   BI-LEVEL OPTIMIZATION FOR DENOISING

Under the bi-level optimization setting, the inner optimization is to derive the recommendation model $f$, and the outer optimization is to refine the sample weight variable $\mathcal{W}$. The bi-level optimization is formulated as follows:

$$\min_{\mathcal{W}} \mathcal{L}_{\mathcal{W}}(f_{\theta_f^*}), \quad s.t., \theta_f^* = \arg\min_{\theta_f} \mathcal{L}_{rec}(\mathcal{W}, \mathcal{D}), \quad (4)$$

where $\mathcal{L}_{\mathcal{W}}$ is a defined loss for optimizing $\mathcal{W}$, and $\mathcal{L}_{rec}$ is a recommendation loss over interaction data $\mathcal{D}$ with weights $\mathcal{W}$. When optimizing $\mathcal{W}$ with $\mathcal{L}_{\mathcal{W}}$, we fix the recommendation model's current parameters $\theta_f$. Analogously, we fix the current weights when optimizing the recommendation model $f$ with $\mathcal{L}_{rec}$. In this framework, the weight set $\mathcal{W}$ is a learnable global variable. Therefore, the adjustments of $\mathcal{W}$ in previous iterations affect the computation of the current $\mathcal{W}$, which avoids the limitation that existing methods cannot share weight information throughout model training. In addition, the inner optimization guides the outer optimization for learning appropriate weight variable $\mathcal{W}$, eliminating the need for prior knowledge.

Despite the benefits of bi-level optimization, it also comes with new challenges. As stated in Section 1, in our setting it is infeasible to store $\mathcal{W}$ due to its sheer size, which is equal to $|\mathcal{U}| \times |\mathcal{I}|$ when all possible interactions are considered. Our empirical findings indicate that explicitly storing $\mathcal{W}$ for even moderate-sized datasets can cause out-of-memory issues. Besides, an intuitive approach to solving the bi-level optimization is to use the gradient information generated in the inner optimization to guide the learning of $\mathcal{W}$

in the outer optimization [49], which requires the full training of the inner optimization to obtain suitable gradients for the outer optimization, leading to a significant time cost.

## 4   SOLVING BI-LEVEL OPTIMIZATION

In this section, we provide an efficient solution to the bi-level optimization problem in Section 3.

### 4.1   Generator-Based Parameter Generation

To reduce the storage demand of $\mathcal{W}$, we apply a generative network to generate each entry $w_{u,i}$ in $\mathcal{W}$. Equation 5 shows the generator as follows:

$$w_{u,i} = g(f_{\theta_f^*}(\mathbf{z}_u), f_{\theta_f^*}(\mathbf{z}_i)), \quad (5)$$

where $g$ is the function of the generator whose parameters are $\theta_g$. The input of $g$ is the user embedding $\mathbf{z}_u$ and the item embedding $\mathbf{z}_i$ (transformed by the recommendation model $f_{\theta_f^*}$), and the output of $g$ is the corresponding $w_{u,i}$.

The network architecture of $g$ is a simplified version of the AE (AutoEncoder) network [24]. As per the conventional AE network, it comprises two essential components, namely the encoder layer and the decoder layer. The formulation of $g$ is outlined as follows:

$$\mathbf{z} = g_e(concat(f_{\theta_f^*}(\mathbf{z}_u), f_{\theta_f^*}(\mathbf{z}_i))), \quad (6)$$

$$w_{u,i} = g_d(\mathbf{z}), \quad (7)$$

where $g_e$ and $g_d$ refer to the fully connected layers of the encoder and decoder, respectively. In an effort to reduce the number of network parameters, we exclusively employ 1-lay fully connected layer as the encoder and decoder layers.

As demonstrated by Equation 5, the weight $w_{u,i}$ is not solely dependent on an update to the generator $g$, but is contingent upon the interplay between the generator $g$, the optimized recommendation model $f_{\theta_f^*}$, the user representation $\mathbf{z}_u$, and the item representation $\mathbf{z}_i$. This design is efficient because of two reasons: (1) Great scalability and flexibility. The generator unifies all $w_{u,i}$ in

$\mathcal{W}$. The input $\mathbf{z}_u$ and $\mathbf{z}_i$ encoded by $f_{\theta_f^*}$ can guarantee that $w_{u,i}$ becomes close if two user-item pairs are similar in the latent space, and the generator $g$ provides $w_{u,i}$ more variation. (2) Low space cost. Despite the additional training time of the generator, this design can avoid a significant space cost. More analysis of the model size and time complexity can be seen in Section 4.4.

After using the generator, the update of $\mathcal{W}$ in outer optimization changes to the training of generator $g$, and the bi-level optimization changes as follows:

$$\min_{\theta_g} \mathcal{L}_{\mathcal{W}}(f_{\theta_f^*}), \quad s.t., \quad \theta_f^* = \arg\min_{\theta_f} \mathcal{L}_{rec} \quad (8)$$

## 4.2 One-Step Gradient-Matching-Based Solving

To reduce the high time cost, in this part, we propose a one-step gradient matching scheme. Specifically, we perform only a single update to the recommendation model within the inner optimization. To ensure that the gradient information from a single update can effectively contribute to the outer optimization, we introduce a gradient-matching mechanism that matches the gradient information of the two models during the inner optimization.

This idea builds upon a prior work [49] that addresses the bi-level optimization using gradient information from multiple updates. During the inner optimization, the recommendation model is optimized using gradient descent with a learning rate $\eta$, shown as:

$$\theta_{f_{t+1}} \leftarrow \theta_{f_t} - \eta \nabla_{\theta_{f_t}} \mathcal{L}_{rec}, \quad (9)$$

where $\nabla_{\theta_{f_t}} \mathcal{L}_{rec}$ represents the $t$-step gradient of recommendation model $f$ with regard to $\mathcal{L}_{rec}$, and multiple updates are allowed. In the outer optimization, it is required to unroll the whole training trajectory of the inner problem, which is computationally challenging and hinders scalability, leading us to develop an efficient optimization technique.

With the aim to reduce the high time cost associated with the optimization process, it is logical to think about updating the parameters in the inner optimization only once: in Equation 9, only one gradient $\nabla_{\theta_f} \mathcal{L}_{rec}$ is generated, and the parameter $\theta_f$ is updated immediately. We thus remove the step notation $t$ from Equation 9 since we only perform the one-step calculation. After obtaining the (one-step) gradient $\nabla_{\theta_f} \mathcal{L}_{rec}$ and update $\theta_f$, we proceed to update $\theta_g$ based on gradient $\nabla_{\theta_f} \mathcal{L}_{rec}$. This leads to an efficient optimization, as it avoids the need for multiple updates *w.r.t* parameters in the inner optimization, as required by the outer optimization.

However, while the one-step gradient leading training accelerates the computational process, it may negatively impact the training of the generator $g$, as it reduces the amount of gradient information available. Multiple gradients contain more valuable information, such as the dynamic trend of the gradient, which is necessary for guiding the training of $g$ in the correct direction. To tackle this issue, we propose to use additional recommendation losses to train recommendation model simultaneously, which provide more gradient information. It is important to note that blindly incorporating all gradient information into the training of $g$ would not result in a positive change. Instead, an elegant combination of the gradient information is necessary to effectively train $g$.

Motivated by [15, 16, 21, 47], we adopt the gradient matching scheme to guide the optimization of the generator $g$. Specifically,

gradient matching involves defining two types of losses and optimizing the model by minimizing the difference between them, i.e., the models trained on these two losses converge to similar parameters. With the gradient matching scheme, the ultimate form of the bi-level optimization is as follows:

$$\min_{\theta_g} D(\nabla_{\theta_{f_t}} \mathcal{L}_{rec1}, \nabla_{\theta_{f_t}} \mathcal{L}_{rec2}),$$
$$s.t., \theta_f^* = \arg\min_{\theta_f}(\mathcal{L}_{rec1}(\mathcal{W}, \mathcal{D}) + \alpha \mathcal{L}_{rec2}(\mathcal{W}, \mathcal{D})), \quad (10)$$

where weight $\alpha$ control the balance of two recommendation loss, and $\mathcal{L}_{rec1}$ and $\mathcal{L}_{rec2}$ can be any two common recommendation loss functions. $\nabla_{\theta_{f_t}} \mathcal{L}_{rec1}$ and $\nabla_{\theta_{f_t}} \mathcal{L}_{rec2}$ are gradients of $\mathcal{L}_{rec1}$ and $\mathcal{L}_{rec2}$ respectively. In the gradient matching, our purpose is to optimize $g$ such that two different one-step gradients are pulled closer, and the $D(\cdot)$ is the distance function [16], which is shown as follows:

$$D(G^1, G^2) = \sum_{c=1}^{d_2}(1 - \frac{G_c^1 \cdot G_c^2}{\|G_c^1\| \|G_c^2\|}), \quad (11)$$

where $G^1 \in \mathbb{R}^{d_1 \times d_2}$ and $G^2 \in \mathbb{R}^{d_1 \times d_2}$ are gradients at a specific layer, $d_1, d_2$ are the number of rows and columns of the gradient matrices, and $G_c^1 \in \mathbb{R}^{d_1}$ and $G_c^2 \in \mathbb{R}^{d_1}$ refer to the $c$-th column vectors of the gradient matrices.

The advantages of using the gradient matching scheme include: (1) Diverse perspectives. Gradient matching entails utilizing two distinct gradients, enabling the model to update and explore the data from various angles, thereby enhancing its performance. (2) Improved stability. By aligning or matching the different gradients, the model's training process can potentially remain stable. This alignment ensures reasonable exploitation of the data. In clean samples, the gradients corresponding to the two losses exhibit minimal differences and remain stable. Conversely, noisy samples may result in substantial differences between the gradients. By aligning the two gradients, the scheme effectively emphasizes the differences between clean and noisy samples, further highlighting the disparity between them. (3) Efficiency in computational cost. The use of one-step gradient information during the optimization of the generator $g$ reduces the computational overhead as it is generated during the training of the recommendation model;

Note that, since $\nabla_{\theta_{f_t}} \mathcal{L}_{rec1}$ and $\nabla_{\theta_{f_t}} \mathcal{L}_{rec2}$ are unstable in the initial stage of recommendation model training, the one-step gradient matching might fail if directly using to train generator $g$ because of the large difference between $\nabla_{\theta_{f_t}} \mathcal{L}_{rec1}$ and $\nabla_{\theta_{f_t}} \mathcal{L}_{rec2}$. To get a relatively suitable $\nabla_{\theta_{f_t}} \mathcal{L}_{rec1}$ and $\nabla_{\theta_{f_t}} \mathcal{L}_{rec2}$, we will pre-train the recommendation model $f$ to a stable state, i.e., the warm-up stage, and then use one-step gradient matching to train $g$.

## 4.3 The Framework BOD

By incorporating the concept of generator-based parameter generation and one-step gradient matching scheme, we effectively minimize the computational demands of bi-level optimization for denoising. The resulting framework, referred to as BOD (as depicted in Fig. 2), serves as a general solution for recommendation denoising. Algorithm 1 shows the learning process of BOD. It is described as follows: (1) Towards inner optimization solving (Line 3-6): the inner optimization is to train the recommendation model over
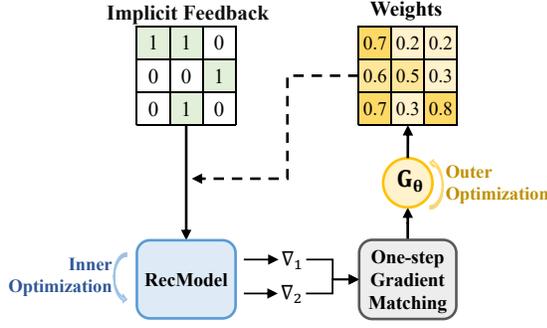
**Figure 2: The framework** BOD**.**

implicit feedback data with weights. Using two losses, we can generate corresponding one-step gradients $\nabla_{\theta_{ft}^{rec1}}$ and $\nabla_{\theta_{ft}^{rec2}}$. (2) Towards outer optimization solving (Line 7-8): the outer optimization uses gradient matching to guide the generator training, which can generate weights for the next step of bi-level optimization.

---

**Algorithm 1:** The proposed BOD framework

**Input:** interactions data $\mathcal{D}$, recommendation model $f$, generator $g$

**Output:** recommendation model with optimal para. $f_{\theta_f^*}$

1   warm-up the recommendation model $f$ and randomly initialize parameters of $g$;

2   **while** *not converged* **do**

     // inner optimization: fix $g$'s parameters $\theta_g$

3      sample a batch tuple $(u, i, j)$, and generate corresponding $w_{u,i}$ and $w_{u,j}$ by $g$;

4      calculate $\mathcal{L}_{rec1}$ and $\mathcal{L}_{rec2}$;

5      one update $f_{\theta_f}$ to $f_{\theta_f^*}$ based on $\mathcal{L}_{rec1} + \alpha\mathcal{L}_{rec2}$;

6      record $\nabla_{\theta_{ft}^{rec1}}$ and $\nabla_{\theta_{ft}^{rec2}}$;

     // outer optimization:fix $f$'s optimal parameters $\theta_f^*$

7      calculate matching distance loss $D(\nabla_{\theta_{ft}^{rec1}}, \nabla_{\theta_{ft}^{rec2}})$;

8      optimize $\theta_g$ based on $D(\nabla_{\theta_{ft}^{rec1}}, \nabla_{\theta_{ft}^{rec2}})$;

---

The framework BOD is characterized by its generality, as demonstrated in Equation 10, where the specific recommendation model is not specified. To apply the framework to a specific recommendation model $f$, the two recommendation losses for gradient matching must be specified, as described in Algorithm 1. The choice of the recommendation losses is flexible and not subject to any constraints. As a demonstration, we choose the BRP loss [22, 23, 48], which is a classical method used in recommendation denoising (referred to in Section 2), and the AU loss [26] (presented in Appendix A.1), which is a recent development. Through this example, the denoising method for the specific recommendation model $f$ can be generated based on the BOD framework and the two loss functions specified. Our contribution to the field is not the specific choice of recommendation losses, but rather the design of a general denoising

framework BOD that can incorporate any two recommendation losses for recommendation denoising.

## 4.4 Complexity Analysis

We summarize the complexity of the base model, existing denoising methods, BOD without the generator, and BOD in Table 1.

**Model Size.** The parameters of BOD come from two parts: (1) the parameters of recommendation models, which we assume is $M$; and (2) the parameters of the generator. For the generator, we need additional parameters for the encoder layer $EN \in \mathbb{R}^{2d \times d_g}$, and the decoder layer $DE \in \mathbb{R}^{d_g \times 1}$, where $d$ is the embedding size of user and item in recommendation model, and $d_g$ is the hidden layer size of the generator. Overall, the space cost of BOD is $M + EN + DE$, which is negligible compared with full weights matrix $\mathcal{W} \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{I}|}$. This shows that the space cost can be greatly reduced by using the generator.

**Time Complexity.** The complexity of BOD consists of two parts. (1) Inner optimization: the update of recommendation model parameters, which we assume is $O(M)$; and (2) Outer optimization: the update of the generator, which generates weights. The time complexity of the generator includes encoder layer $O(d(d_g))$ and the decoder layer $O(d_g)$, respectively. It is obvious that our scheme can save a lot of time costs because $O(M|\mathcal{U}||\mathcal{I}|)$ appears in other denoising methods are much larger.

**Table 1: Model complexity comparison.**

| Model | Model Size | Time Complexity |
|---|---|---|
| Base Model | $M$ | $O(M)$ |
| Denoising methods | $M + |\mathcal{U}||\mathcal{I}|$ | $O(M|\mathcal{U}||\mathcal{I}|)$ |
| BOD *w/o gen* | $M + |\mathcal{U}||\mathcal{I}|$ | $O(M + |\mathcal{U}||\mathcal{I}|)$ |
| BOD | $M + EN + DE$ | $O(M + d(d_g) + (d_g))$ |

## 5 EXPERIMENTS

This section tests the effectiveness of our proposed BOD. Specifically, we aim to answer the following questions.(RQ1): How does the performance and robustness of BOD against cutting-edge denoising methods? (RQ2): What is the running time of BOD? (RQ3): How does the inner optimization affect BOD? (RQ4): How does the generator affect BOD? (RQ5): What is the stability of BOD? (RQ6): How about the generalizability of BOD?

**Table 2: Statistics of datasets.**

| Dataset | #Users | #Items | #Interactions | Density |
|---|---|---|---|---|
| Beauty | 22,363 | 12,099 | 198,503 | 0.073% |
| iFashion | 300,000 | 81,614 | 1,607,813 | 0.006% |
| Yelp2018 | 31,668 | 38,048 | 1,561,406 | 0.130% |

## 5.1 Performance and Robustness (RQ1)

## 5.2 Experimental Settings

**Datasets.** Three commonly used public benchmark datasets: Beauty [26], iFashion [36], and Yelp2018 [11], are used in our experiments. The dataset statistics are shown in Table 2.

**Table 3: Performance comparison of different denoising methods on the robust recommendation. The highest scores are in bold, and the second best are with underlines. R and N refer to Recall and NDCG, respectively.**

| Dataset | | Beauty | | | | iFashion | | | | Yelp2018 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Base Model | Method | R@10 | R@20 | N@10 | N@20 | R@10 | R@20 | N@10 | N@20 | R@10 | R@20 | N@10 | N@20 |
| GMF | Normal | 0.0639 | 0.0922 | 0.0429 | 0.0518 | 0.0286 | 0.0462 | 0.0157 | 0.0205 | 0.0290 | 0.0500 | 0.0331 | 0.0411 |
| | WBPR | 0.0640 | 0.0919 | 0.0430 | 0.0516 | 0.0284 | 0.0461 | 0.0158 | 0.0205 | 0.0289 | 0.0498 | 0.0330 | 0.0409 |
| | WRMF | 0.0714 | 0.1039 | 0.0480 | 0.0582 | 0.0397 | 0.0589 | 0.0211 | 0.0265 | 0.0304 | 0.0539 | 0.0338 | 0.0422 |
| | T-CE | 0.0785 | 0.1085 | 0.0532 | 0.0654 | 0.0385 | 0.0583 | 0.0193 | 0.0256 | 0.0302 | 0.0538 | 0.0335 | 0.0421 |
| | DeCA | 0.0767 | 0.1081 | 0.0535 | 0.0669 | 0.0426 | 0.0626 | <u>0.0234</u> | <u>0.0297</u> | 0.0298 | 0.0532 | 0.0336 | 0.0415 |
| | SGDL | <u>0.0804</u> | <u>0.1092</u> | <u>0.0546</u> | <u>0.0675</u> | <u>0.0435</u> | <u>0.0646</u> | 0.0226 | 0.0289 | <u>0.0305</u> | <u>0.0541</u> | <u>0.0351</u> | <u>0.0436</u> |
| | BOD | **0.0843** | **0.1193** | **0.0579** | **0.0688** | **0.0596** | **0.0860** | **0.0355** | **0.0426** | **0.0357** | **0.0621** | **0.0408** | **0.0513** |
| NCF | Normal | 0.0738 | 0.1075 | 0.0488 | 0.0593 | 0.0404 | 0.0632 | 0.0227 | 0.0288 | 0.0281 | 0.0493 | 0.0324 | 0.0406 |
| | WBPR | 0.0741 | 0.1082 | 0.0494 | 0.0600 | 0.0410 | 0.0638 | 0.0231 | 0.0292 | 0.0288 | 0.0499 | 0.0331 | 0.0410 |
| | WRMF | 0.0745 | 0.1097 | 0.0498 | 0.0611 | 0.0455 | 0.0667 | 0.0247 | 0.0306 | 0.0292 | 0.0503 | 0.0335 | 0.0413 |
| | T-CE | 0.0798 | 0.1123 | 0.0506 | 0.0632 | 0.0546 | 0.0684 | 0.0268 | 0.0368 | 0.0301 | 0.0521 | 0.0335 | 0.0415 |
| | DeCA | <u>0.0886</u> | <u>0.1265</u> | 0.0588 | 0.0646 | 0.0547 | 0.0756 | <u>0.0315</u> | <u>0.0421</u> | 0.0297 | 0.0535 | 0.0351 | 0.0459 |
| | SGDL | 0.0864 | 0.1235 | <u>0.0610</u> | <u>0.0698</u> | <u>0.0588</u> | <u>0.0846</u> | 0.0311 | 0.0329 | <u>0.0325</u> | <u>0.0598</u> | <u>0.0388</u> | <u>0.0465</u> |
| | BOD | **0.0959** | **0.1352** | **0.0652** | **0.0774** | **0.0624** | **0.0909** | **0.0366** | **0.0442** | **0.0370** | **0.0633** | **0.0434** | **0.0531** |
| NGCF | Normal | 0.0726 | 0.1080 | 0.0465 | 0.0575 | 0.0351 | 0.0579 | 0.0188 | 0.0250 | 0.0231 | 0.0418 | 0.0263 | 0.0336 |
| | WBPR | 0.0731 | 0.1092 | 0.0476 | 0.0589 | 0.0364 | 0.0598 | 0.0197 | 0.0266 | 0.0240 | 0.0422 | 0.0271 | 0.0343 |
| | DeCA | 0.0834 | 0.1275 | 0.0587 | 0.0665 | <u>0.0587</u> | <u>0.0855</u> | <u>0.0247</u> | <u>0.0356</u> | 0.0287 | 0.0486 | 0.0299 | 0.0435 |
| | SGDL | <u>0.0935</u> | <u>0.1389</u> | <u>0.0635</u> | <u>0.0758</u> | 0.0566 | 0.0848 | 0.0225 | 0.0345 | <u>0.0297</u> | <u>0.0503</u> | <u>0.0348</u> | <u>0.0445</u> |
| | BOD | **0.1015** | **0.1415** | **0.0702** | **0.0827** | **0.0716** | **0.1053** | **0.0390** | **0.0480** | **0.0331** | **0.0565** | **0.0377** | **0.0466** |
| LightGCN | Normal | 0.0855 | 0.1221 | 0.0561 | 0.0675 | 0.0429 | 0.0661 | 0.0247 | 0.0310 | 0.0308 | 0.0532 | 0.0361 | 0.0445 |
| | WBPR | 0.0864 | 0.1261 | 0.0588 | 0.0685 | 0.0431 | 0.0662 | 0.0253 | 0.0316 | 0.0317 | 0.0529 | 0.0365 | 0.0448 |
| | DeCA | 0.0967 | 0.1345 | 0.0665 | 0.0754 | 0.0540 | 0.0865 | 0.0354 | 0.0435 | 0.0337 | 0.0511 | 0.0432 | 0.0524 |
| | SGDL | 0.0946 | 0.1365 | 0.0677 | 0.0769 | 0.0591 | 0.0908 | 0.0342 | 0.0415 | 0.0339 | 0.0541 | 0.0441 | 0.0575 |
| | SGL | <u>0.1005</u> | <u>0.1422</u> | <u>0.0692</u> | <u>0.0821</u> | 0.0665 | 0.0973 | 0.0392 | 0.0475 | 0.0374 | 0.0639 | 0.0436 | 0.0534 |
| | SimGCL | 0.0960 | 0.1316 | 0.0671 | 0.0781 | <u>0.0738</u> | <u>0.1070</u> | <u>0.0434</u> | <u>0.0523</u> | <u>0.0414</u> | **0.0711** | <u>0.0485</u> | **0.0593** |
| | BOD | **0.1095** | **0.1548** | **0.0755** | **0.0895** | **0.0811** | **0.1180** | **0.0477** | **0.0576** | **0.0416** | <u>0.0700</u> | **0.0490** | <u>0.0588</u> |

**Evaluation Metrics.** We split the datasets into three parts (training set, validation set, and test set) with a ratio of 7:1:2. Two common evaluation metrics are used, Recall@$K$ and NDCG@$K$. We set $K$=10 and $K$=20. Each metric is conducted 10 times, and then we report the average results.

**Baselines.** The main objective of this paper is to denoise the feedback to improve the performance of recommender systems. For this purpose, four commonly used (implicit feedback-based) recommendation models are chosen as **base models** for recommendation denoising:
- GMF [10] is a generalized version of matrix factorization based recommendation model.
- NCF [13] generalizes collaborative filtering with a Multi-Layer Perceptron.
- NGCF [31] applies graph convolution network (GCN) to encode user-item bipartite graph.
- LightGCN [11] is a state-of-the-art graph model, which discards the nonlinear feature transformations to simplify the design of GCN for the recommendation.

To compare the denoising effect, we first choose four existing denoising methods for the above recommendation models:
- WBPR [6] considers popular but uninteracted items as true negative samples. We take the number of interactions of items as the popularity.

- WRMF [14] uses weighted matrix factorization whose weights are fixed to denoise recommendation.
- T-CE [30] is the denoising method, which uses the binary cross-entropy (BCE) loss [19, 44] to assign weights to large loss samples with a dynamic threshold. Because the BCE loss only applies to some recommendation models, we can only get (and thus report) the results of T-CE in GMF and NCF, but not in the other two base models.
- DeCA [32] combines predictions of two different models to consider the disagreement of noisy samples.
- SGDL [7] is the state-of-the-art denoising model, which collects clean interactions in the initial training stages, and uses them to distinguish noisy samples based on the similarity of collected clean samples.

To further confirm the effectiveness of our model, we also compare BOD with the state-of-the-art robust recommendation models SGL and SimGCL:
- SGL [36] applies edge dropout to modify discrete graph structure randomly and trains different graphs based on contrastive learning to enhance the representation of nodes.
- SimGCL [41] simplifies the contrastive loss and optimizes representations' uniformity.

Since SGL and SimGCL use LightGCN as the base model in their work, so we only compare them when LightGCN is used and ignore the results on the other base models.

**Our Method BOD.** If not stressed, the two loss functions of BOD are BPR and AU losses, respectively.

**Parameter Settings.** We implement BOD in Pytorch. If not stressed, we use recommended parameter settings for all models, in which batch size, learning rate, embedding size, and the number of Light-GCN layers are 128, 0.001, 64, and 2, respectively. For SGL, the edge dropout rate is set to 0.1. We optimize them with Adam [17] optimizer and use the Xavier initializer [8] to initialize the model parameters.

**Performance Comparison.** We first compare BOD with exist-ing denoising methods and robust recommendation methods on three different datasets. Table 3 shows the results, where figures in bold represent the best performing indices, and the runner-ups are presented with underlines. According to Table 3, we can draw the following observations and conclusions:

- The proposed BOD can effectively improve the performance of four base models and show the best/second performance in all cases. We attribute these improvements to the storage and uti-lization of weights. In the bi-level process, BOD dynamically up-dates and stores the weights for all samples, including observed and unobserved samples. While other baselines (e.g., DeCA and SGDL) are insufficient to provide dynamically updated weights.
- All denoising approaches have better results than normal train-ing, which indicates the validity of denoising in recommenda-tions. The results are consistent with prior studies [7, 32].
- The improvement on Yelp2018 dataset is less significant than that on other datasets. One possible reason is that Yelp2018 dataset has a higher density than the others. Thus there are enough pure informative interactions for identifying user be-havior, compensating for the effect of noisy negative samples.

**Robustness Comparison.** We also conduct experiments to check BOD's robustness. Following previous work [36], we add a certain proportion of adversarial examples (i.e., 5%, 10%, 15%, 20% negative user-item interactions) into the training set, while keeping the testing set unchanged. Figure. 3 shows the experimental results on Beauty dataset. We can see the results of BOD maintain the best in all cases, despite performances reduced as noise data adding. Besides, the trend of BOD's drop percent curve keeps the minor change, which illustrates that BOD is the least affected by noise. This suggests that the denoising process of BOD can better figure out useful patterns under noisy conditions.

## 5.3 Time Complexity Analysis (RQ2)

In this part, we set LightGCN as the base model and report the real running time of compared methods for one epoch. The results in Table 4 are collected on an Intel(R) Core(TM) i9-10900X CPU and a GeForce RTX 3090Ti GPU. As shown in Table 4, we can see the running time increases with the volume of the datasets. Besides, the running speed of BOD is much quicker than denoising methods (T-CE, DeCA, and SGDL), even though they only deal with positive samples. Furthermore, SGL and SimGCL do not need additional time on denoising training, but the time cost of BOD is still less
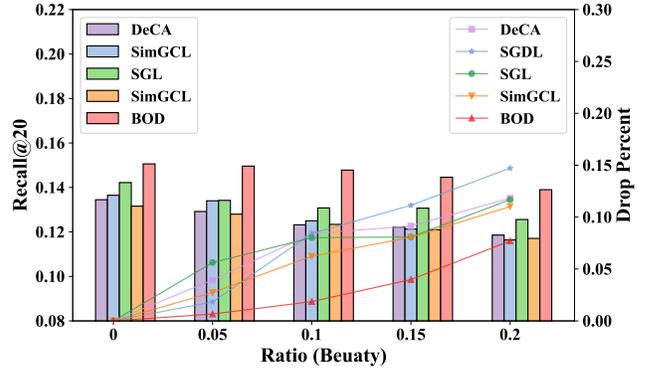


**Figure 3: Model performance w.r.t. noise ratio.**

than those of SGL and SimGCL, which proves that the generator-based one-step gradient matching scheme can save much time from the computation.

## 5.4 The Effect of Inner Optimization on BOD (RQ3)

In Equation 10 of Section 4.2, we design the objective function of the inner optimization using the weighted sum of the two recommen-dation losses. However, in order to solve for the outer optimization, we only need the inner optimization to provide (two gradients). It may appear that solving the inner optimization with the weighted sum of losses is unnecessary. Yet, our analysis demonstrates the importance of including both losses. For this purpose, we discuss the following three cases:

- $BOD_{BPR}$: We calculate BPR and AU loss both to get correspond-ing gradients, but only use the BPR's gradient to optimize the recommendation model in the inner optimization.
- $BOD_{AU}$: Similar to the previous case, the difference is that we use the gradient generated from AU loss to optimize the recom-mendation model.
- $BOD_{BPR+AU}$: In this (default) case, we use BPR's gradient and AU's gradients with the weight $\alpha$ to control the balance, to optimize recommendation model.

It is worth reminding that we hardly execute BOD without the generator due to the expensive space cost ("out of memory" will be displayed). As Table 5 shows, we find that if we only use one loss to optimize the recommendation model in the inner optimization, its performance is not as good as the case when two losses are used.

**Table 4: Running time in seconds per epoch.**

| method | Beauty | iFashion | Yelp2018 |
|---|---|---|---|
| Normal | 3.47(+0.08) | 89.22(+3.86) | 63.39(+3.23) |
| T-CE | 41.5(+2.06) | 765.64(+26.87) | 631.45(+21.33) |
| DeCA | 25.10(+1.28) | 556.50(+18.78) | 478.68(+19.01) |
| SGDL | 32.12(+1.82) | 688.21(+20.43) | 564.46(+16.92) |
| SGL | 8.12(+0.58) | 302.32(+12.26) | 270.17(+11.81) |
| SimGCL | 7.26(+0.40) | 273.66(+11.31) | 159.24(+8.01) |
| BOD | **6.82(+0.16)** | **97.38(+4.11)** | **75.98(+3.44)** |

**Table 5: The Effect of Inner optimization on** BOD.

| Dataset | | Beauty | | iFashion | | Yelp2018 | |
|---------|----------|--------|--------|----------|--------|----------|--------|
| Base Model | Component | R@20 | N@20 | R@20 | N@20 | R@20 | N@20 |
| GMF | $BOD_{BPR}$ | 0.1029 | 0.0599 | 0.0664 | 0.0306 | 0.0531 | 0.0422 |
| | $BOD_{AU}$ | 0.1133 | 0.0659 | 0.0705 | 0.0313 | 0.0547 | 0.0461 |
| | $BOD_{BPR+AU}$ | **0.1193** | **0.0688** | **0.0860** | **0.0426** | **0.0621** | **0.0513** |
| NCF | $BOD_{BPR}$ | 0.1208 | 0.0683 | 0.0806 | 0.0385 | 0.0564 | 0.0432 |
| | $BOD_{AU}$ | 0.1289 | 0.0739 | 0.0824 | 0.0386 | 0.0511 | 0.0411 |
| | $BOD_{BPR+AU}$ | **0.1352** | **0.0774** | **0.0909** | **0.0442** | **0.0633** | **0.0531** |
| NGCF | $BOD_{BPR}$ | 0.1295 | 0.0715 | 0.0976 | 0.0354 | 0.0476 | 0.0406 |
| | $BOD_{AU}$ | 0.1297 | 0.0717 | 0.0983 | 0.0385 | 0.0545 | 0.0454 |
| | $BOD_{AU}$ | **0.1415** | **0.0827** | **0.1053** | **0.0480** | **0.0565** | **0.0466** |
| LightGCN | $BOD_{BPR}$ | 0.1386 | 0.0868 | 0.0984 | 0.0415 | 0.0598 | 0.0523 |
| | $BOD_{AU}$ | 0.1400 | 0.0876 | 0.1054 | 0.0443 | 0.0625 | 0.0546 |
| | $BOD_{BPR+AU}$ | **0.1548** | **0.0895** | **0.1180** | **0.0576** | **0.7000** | **0.0588** |

## 5.5 The Effect of Generator on BOD (RQ4)

We delve deeper into the impact of generator design on performance. For this reason, we compare the applied AE method (using 2-layer Multi-Layer Perceptron, denoted as 2-MLP) with the classical variational autoencoder method (denoted as VAE) [18]. We also change the number of MLP layers to 1 and 3 (denoted as 1-MLP and 3-MLP) to further illustrate the difference. We perform experiments on the Beauty dataset and use LightGCN as the recommendation encoder. The experimental in Table 6 outcomes demonstrate that employing the 2-MLP method as a generator can yield improvements to a certain extent compared to the classical VAE method. In addition, the 2-MLP increases with the number of MLP layers and reaches stability at 2 layers. Therefore, we used the 2-layer MLP based AE method in the paper.
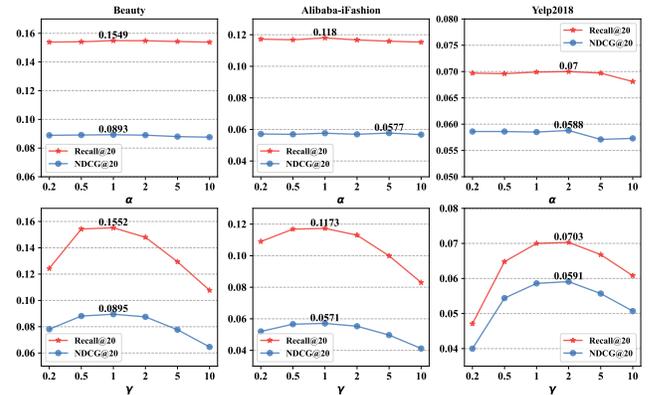
**Table 6: The Effect of Generator on** BOD.

| Method | Recall@10 | Recall@20 | NDCG@10 | NDCG@20 |
|--------|-----------|-----------|---------|---------|
| 1-MLP | 0.1011 | 0.1498 | 0.0718 | 0.0855 |
| 3-MLP | 0.1095 | 0.1548 | 0.0755 | 0.0894 |
| VAE | 0.1074 | 0.1532 | 0.0731 | 0.0888 |
| 2-MLP (ours) | **0.1095** | **0.1548** | **0.0755** | **0.0895** |

## 5.6 Parameters Sensitivity Analysis (RQ5)

We investigate the model's sensitivity to the hyperparameters $\alpha$ in Equation 10 and $\gamma$ in Equation 12 (see Appendix). We fix $\alpha$ as 1 and try different values of $\gamma$ with the set [0.2, 0.5, 1, 2, 5, 10], then we fix $\gamma$ as 1 and try different values of $\alpha$ with the set [0.2, 0.5, 1, 2, 5, 10]. As shown in Fig. 4, we can observe the performance keep stable as $\alpha$ changes. Besides, a similar trend in the performance of $\gamma$ increases first and then decreases on all the datasets. BOD reaches its best performance when $\gamma$ is in [0.5, 1, 2], which suggests that the weights of alignment and uniformity should be balanced at the same level.

## 5.7 Generalizability Analysis (RQ6)

We emphasize that our framework BOD is generic, which means that it can apply to any recommendation model in a plug-and-play fashion. One thing of interest is whether we can further improve



**Figure 4: Parameter sensitivity with regard to** $\alpha$ **and** $\gamma$.

the performance of the latest robust recommendation models SGL and SimGCL. To this end, we use SGL and SimGCL as base models for BOD, and the results are shown in Fig. 5. From the experimental results, we conclude that BOD does further improve the performance of these two methods. This confirms the generalizability of our method.



**Figure 5: The generalizability of** BOD **to SGL and SimGCL.**

## 6 RELATED WORK

Our study endeavors to denoise recommendations by employing a bi-level optimization framework. To begin with, we present various bi-level optimization methods applied in recommendation. Concurrently, researchers attempt to denoise implicit feedback through many avenues, and we mainly focus on denoising methods without extra information in this paper.

**Bi-level Optimization in Recommendation.** Bi-level optimization techniques have been extensively utilized in the domain of recommender systems [2, 5, 12]. For instance, APR introduces adversarial perturbations to model parameters as a means to improve recommendations through adversarial training. IRGAN employs a game-theoretical minimax game to iteratively optimize both a generative model and a discriminative model, ultimately generating refined recommendation results. Adv-MultVAE incorporates adversarial training into MultVAE architecture to eliminate implicit information pertaining to protected attributes while preserving recommendation performance. These approaches showcase the effectiveness of bi-level optimization in enhancing recommendation

performance, yet there is currently a dearth of research addressing denoising within this context.

**Denoising Methods.** A straightforward way approach [6, 34] to denoising implicit feedback is to select clean and informative samples and use them to train a recommendation model. For example, WBPR [6] notices that popular items are more likely real negative instances if they have missing actions, and then samples popular items with higher probabilities. IR [34] produces synthetic labels for user preferences depending on the distinction between labels and predictions, to uncover false-positive and false-negative samples. However, the variability of their performance is wide due to their reliance on the sampling distribution [43]. Thus, some methods [7, 30, 32] provide contribution weight for each sample during recommendation model training. T-CE [30] proposes to assign lower weights to large loss samples with a dynamic threshold because they find that a sample with high loss is more likely to be noisy. DeCA [32] assumes that different models make relatively similar predictions on clean examples, so it incorporates the training process of two recommendation models and distinguishes clean examples from positive samples based on two different predictions. SGDL [7] claims that models tend to dig out clean samples in the initial training stages, and then adaptively assign weights for samples based on the similarity of collected clean samples. We can notice that existing methods often resort to prior knowledge for determining weights, but prior knowledge has its own scope of application since we cannot be sure of the validity of prior knowledge itself.

**Other Directions.** There are also some robust recommendation methods [42] that do not belong to the above methods. For example, SGL [36] creates two additional views of graphs by adding and dropping edges and leverages contrastive learning to distinguish hard noisy samples. Recently, several works have considered augmentation from the perspective of embeddings. For example, SimGCL [41] considers embeddings distribution, and thus trains models to force the distribution of embeddings to fit uniformity.

It should be noted that previous research often uses additional information (attribute features [46], social information [40], etc.) to process weights. This paper discusses denoising tasks in extreme situations where only implicit feedback information can be obtained, so we do not introduce some denoising methods with extra information [4, 19, 25, 37, 45].

## 7 CONCLUSION

In this paper, we model recommendation denoising as bi-level optimization and propose an efficient solution for the proposed bi-level optimization framework, called BOD. In BOD, a generator-based parameter generation technique to save space, and a gradient-matching-based parameter-solving technique to save time. Extensive experiments on three real-world datasets demonstrate the superiority of BOD. Moving forward, our future research will focus on investigating more efficient approaches within the bi-level optimization framework to address specific challenges associated with recommendation denoising. One existing challenge pertains to the presence of noisy data in the test sets, which cannot be effectively identified during the bi-level optimization process.

## REFERENCES

[1] Zhi Bian, Shaojun Zhou, Hao Fu, Qihong Yang, Zhenqi Sun, Junjie Tang, Guiquan Liu, Kaikui Liu, and Xiaolong Li. 2021. Denoising user-aware memory network for recommendation. In *Proceedings of ACM Conference on Recommender Systems 2021*. 400–410.

[2] Dong-Kyu Chae, Jin-Soo Kang, Sang-Wook Kim, and Jung-Tae Lee. 2018. Cfgan: A generic collaborative filtering framework based on generative adversarial networks. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 137–146.

[3] Jiawei Chen, Hande Dong, Yang Qiu, Xiangnan He, Xin Xin, Liang Chen, Guli Lin, and Keping Yang. 2021. AutoDebias: Learning to debias for recommendation. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval 2021*. 21–30.

[4] 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.

[5] Christian Ganhör, David Penz, Navid Rekabsaz, Oleg Lesota, and Markus Schedl. 2022. Unlearning Protected User Attributes in Recommendations with Adversarial Training. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2142–2147.

[6] 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.

[7] 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 International ACM SIGIR Conference on Research and Development in Information Retrieval 2022*. ACM, 1412–1422.

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

[9] Lei Guo, Hongzhi Yin, Qinyong Wang, Tong Chen, Alexander Zhou, and Nguyen Quoc Viet Hung. 2019. Streaming session-based recommendation. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining 2019*. 1569–1577.

[10] Xiangnan He and Tat-Seng Chua. 2017. Neural factorization machines for sparse predictive analytics. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval 2017*. 355–364.

[11] 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 International ACM SIGIR Conference on Research and Development in Information Retrieval 2020*. 639–648.

[12] Xiangnan He, Zhankui He, Xiaoyu Du, and Tat-Seng Chua. 2018. Adversarial personalized ranking for recommendation. In *The 41st International ACM SIGIR conference on research & development in information retrieval*. 355–364.

[13] 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.

[14] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *IEEE International Conference on Data Mining 2008*. 263–272.

[15] Wei Jin, Xianfeng Tang, Haoming Jiang, Zheng Li, Danqing Zhang, Jiliang Tang, and Bing Yin. 2022. Condensing graphs via one-step gradient matching. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining 2022*. 720–730.

[16] Wei Jin, Lingxiao Zhao, Shichang Zhang, Yozen Liu, Jiliang Tang, and Neil Shah. 2021. Graph Condensation for Graph Neural Networks. In *International Conference on Learning Representations 2021*.

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

[18] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

[19] Defu Lian, Yongji Wu, Yong Ge, Xing Xie, and Enhong Chen. 2020. Geography-aware sequential location recommendation. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining 2020*. 2009–2019.

[20] Hongyu Lu, Min Zhang, and Shaoping Ma. 2018. Between clicks and satisfaction: Study on multi-phase user preferences and satisfaction for online news reading. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval 2018*. 435–444.

[21] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. 2021. Dataset distillation with infinitely wide convolutional networks. *Advances in Neural Information Processing Systems* 34 (2021), 5186–5198.

[22] Weike Pan and Li Chen. 2013. Gbpr: Group preference based bayesian personalized ranking for one-class collaborative filtering. In *Proceedings of International Joint Conference on Artificial Intelligence 2013*.

[23] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars BPR Schmidt-Thieme. 2014. Bayesian personalized ranking from implicit feedback. In *Proc. of Uncertainty in Artificial Intelligence 2014*. 452–461.

[24] Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio. 2011. Contractive auto-encoders: Explicit invariance during feature extraction. In *Proceedings of the 28th international conference on international conference on machine learning 2011*. 833–840.

[25] Jiliang Tang, Xia Hu, Huiji Gao, and Huan Liu. 2013. Exploiting local and global social context for recommendation.. In *Proceedings of International Joint Conference on Artificial Intelligence 2013*, Vol. 13. 2712–2718.

[26] Chenyang Wang, Yuanqing Yu, Weizhi Ma, Min Zhang, Chong Chen, Yiqun Liu, and Shaoping Ma. 2022. Towards Representation Alignment and Uniformity in Collaborative Filtering. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining 2022*. 1816–1825.

[27] J Wang, L Yu, W Zhang, Y Gong, Y Xu, B Wang, P Zhang, and D Zhang. 2018. A minimax game for unifying generative and discriminative information retrieval models. *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval 2018* (2018).

[28] Qinyong Wang, Hongzhi Yin, Hao Wang, Quoc Viet Hung Nguyen, Zi Huang, and Lizhen Cui. 2019. Enhancing collaborative filtering with generative augmentation. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining 2019*. 548–556.

[29] Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Proceedings of International Conference on Machine Learning 2020*. PMLR, 9929–9939.

[30] Wenjie Wang, Fuli Feng, Xiangnan He, Liqiang Nie, and Tat-Seng Chua. 2021. Denoising implicit feedback for recommendation. In *Proceedings of the ACM International Conference on Web Search and Data Mining 2021*. 373–381.

[31] 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.

[32] 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.

[33] 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 2021*. 3070–3078.

[34] 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 2021*. 3070–3078.

[35] Fan Wu, Min Gao, Junliang Yu, Zongwei Wang, Kecheng Liu, and Xu Wang. 2021. Ready for emerging threats to recommender systems? A graph convolution-based generative shilling attack. *Information Sciences 2021* 578 (2021), 683–701.

[36] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval 2021*. 726–735.

[37] Hongzhi Yin, Bin Cui, Zi Huang, Weiqing Wang, Xian Wu, and Xiaofang Zhou. 2015. Joint modeling of users' interests and mobility patterns for point-of-interest recommendation. In *Proceedings of the 23rd ACM International Conference on Multimedia*. 819–822.

[38] Hongzhi Yin, Qinyong Wang, Kai Zheng, Zhixu Li, Jiali Yang, and Xiaofang Zhou. 2019. Social influence-based group representation learning for group recommendation. In *IEEE 35th International Conference on Data Engineering 2019*. 566–577.

[39] Junliang Yu, Hongzhi Yin, Min Gao, Xin Xia, Xiangliang Zhang, and Nguyen Quoc Viet Hung. 2021. Socially-aware self-supervised tri-training for recommendation. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery & Data Mining 2021*. 2084–2092.

[40] Junliang Yu, Hongzhi Yin, Jundong Li, Qinyong Wang, Nguyen Quoc Viet Hung, and Xiangliang Zhang. 2021. Self-supervised multi-channel hypergraph convolutional network for social recommendation. In *Proceedings of the Web Conference 2021*. 413–424.

[41] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen. 2022. Are graph augmentations necessary? simple graph contrastive learning for recommendation. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval 2022*. 1294–1303.

[42] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Jundong Li, and Zi Huang. 2022. Self-supervised learning for recommender systems: A survey. *arXiv preprint arXiv:2203.15876* (2022).

[43] Fajie Yuan, Xin Xin, Xiangnan He, Guibing Guo, Weinan Zhang, Chua Tat-Seng, and Joemon M Jose. 2018. fBGD: Learning embeddings from positive unlabeled data with BGD. (2018).

[44] Junwei Zhang, Min Gao, Junliang Yu, Lei Guo, Jundong Li, and Hongzhi Yin. 2021. Double-scale self-supervised hypergraph learning for group recommendation. In *Proceedings of the ACM International Conference on Information & Knowledge Management 2021*. 2557–2567.

[45] Junwei Zhang, Min Gao, Junliang Yu, Xinyi Wang, Yuqi Song, and Qingyu Xiong. 2019. Nonlinear Transformation for Multiple Auxiliary Information in Music Recommendation. In *Proceedings of International Joint Conference on Neural Networks 2019*. IEEE, 1–8.

[46] Wei Zhang, Junbing Yan, Zhuo Wang, and Jianyong Wang. 2022. Neuro-Symbolic Interpretable Collaborative Filtering for Attribute-based Recommendation. In *Proceedings of the ACM Web Conference 2022*. 3229–3238.

[47] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. 2021. Dataset Condensation with Gradient Matching. *Proceedings of International Conference on Learning Representations 2021* 1, 2 (2021), 3.

[48] Tong Zhao, Julian McAuley, and Irwin King. 2014. Leveraging social connections to improve personalized ranking for collaborative filtering. In *Proceedings of the ACM international conference on conference on information and knowledge management 2014*. 261–270.

[49] Daniel Zügner and Stephan Günnemann. 2018. Adversarial Attacks on Graph Neural Networks via Meta Learning. In *Proceedings of International Conference on Learning Representations 2018*.

# A APPENDIX

## A.1 AU Loss

Recent studies [29] have proved alignment and uniformity of representations are related to prediction performance in representation learning, which is also applicable to the recommendation model [26, 41]. The alignment is defined as the expected distance between normalized embeddings of positive pairs, and on the other hand, uniformity is defined as the logarithm of the average pairwise Gaussian potential. AU (Alignment and Uniformity) loss [26] quantified alignment and uniformity in recommendation as follows:

$$L_{rec} = L_{alignment} + \gamma L_{uniformity}, \tag{12}$$

$$L_{alignment} = \mathop{\mathbb{E}}_{(u,i,j) \sim P_{\mathcal{D}}} ||\widetilde{f}(z_u) - \widetilde{f}(z_i)||^2, \tag{13}$$

$$L_{uniformity} = (log \mathop{\mathbb{E}}_{(u,u') \sim P_U} e^{-2||\widetilde{f}(z_u) - \widetilde{f}(z_{u'})||^2}$$
$$+ log \mathop{\mathbb{E}}_{(i,i') \sim P_I} e^{-2||\widetilde{f}(z_i) - \widetilde{f}(z_{i'})||^2})/2, \tag{14}$$

where $P_U(\cdot)$ and $P_I(\cdot)$ are the distributions of the user set and item set, respectively. $|| \cdot ||$ is the $l_1$ distance, and $\widetilde{f}(\cdot)$ indicates the $l_2$ normalized representations of $f(\cdot)$. The weight $\gamma$ controls the desired degree of uniformity.