

Improving Recommendation Fairness via Data Augmentation

Lei Chen
Hefei University of
Technology, Hefei, China
chenlei.hfut@gmail.com

Le Wu*
Hefei University of
Technology, Hefei, China
lewu.ustc@gmail.com

Kun Zhang
Hefei University of
Technology, Hefei, China
zhang1028kun@gmail.com

Richang Hong
Hefei University of
Technology, Hefei, China
hongrc.hfut@gmail.com

Defu Lian
University of Science and Technology
of China, Hefei, China
liandefu@ustc.edu.cn

Zhiqiang Zhang
Ant Group, Hangzhou, China
lingyao.zzq@antfin.com

Jun Zhou
Ant Group, Hangzhou, China
jun.zhoujun@antfin.com

Meng Wang
Hefei University of
Technology, Hefei, China
Institute of Artificial Intelligence,
Hefei Comprehensive National
Science Center, Hefei, China
eric.mengwang@gmail.com

ABSTRACT

Collaborative filtering based recommendation learns users' preferences from all users' historical behavior data, and has been popular to facilitate decision making. Recently, the fairness issue of recommendation has become more and more essential. A recommender system is considered unfair when it does not perform equally well for different user groups according to users' sensitive attributes (e.g., gender, race). Plenty of methods have been proposed to alleviate unfairness by optimizing a predefined fairness goal or changing the distribution of unbalanced training data. However, they either suffered from the specific fairness optimization metrics or relied on redesigning the current recommendation architecture. In this paper, we study how to improve recommendation fairness from the data augmentation perspective. The recommendation model amplifies the inherent unfairness of imbalanced training data. We augment imbalanced training data towards balanced data distribution to improve fairness. Given each real original user-item interaction record, we propose the following hypotheses for augmenting the training data: each user in one group has a similar item preference (click or non-click) as the item preference of any user in the remaining group. With these hypotheses, we generate "fake" interaction behaviors to complement the original training data. After that, we design a bi-level optimization target, with the inner optimization generates better fake data to augment training data with our hypotheses, and the outer one updates the recommendation model

parameters based on the augmented training data. The proposed framework is generally applicable to any embedding-based recommendation, and does not need to pre-define a fairness metric. Extensive experiments on two real-world datasets clearly demonstrate the superiority of our proposed framework. We publish the source code at <https://github.com/newlei/FDA>.

CCS CONCEPTS

• **Information systems** → *Collaborative filtering*; • **Human-centered computing** → User models.

KEYWORDS

user modeling, fairness, data augmentation, fair recommendation

ACM Reference Format:

Lei Chen, Le Wu, Kun Zhang, Richang Hong, Defu Lian, Zhiqiang Zhang, Jun Zhou, and Meng Wang. 2023. Improving Recommendation Fairness via Data Augmentation. In *Proceedings of the ACM Web Conference 2023 (WWW '23)*, May 1–5, 2023, Austin, TX, USA. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3543507.3583341>

1 INTRODUCTION

Recommender systems automatically help users to find items they may like, and have been widely deployed in our daily life for decision-making [14, 37, 40]. Given the original user-item historical behavior data, most recommendation models design sophisticated techniques to learn user and item embeddings, and try to accurately predict users' unknown preferences to items [27, 40, 45]. Recently, researchers argued that simply optimizing recommendation accuracy lead to unfairness issues. Researchers have found that current recommender systems show apparent demographic performance bias of different demographic groups [16, 17]. Career recommender systems tend to favour male candidates compared to females even though they are equally qualified [39]. Besides, recommendation accuracy performance shows significant differences between advantaged user group and disadvantaged user group [41, 43].

*Le Wu is the corresponding author.

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

WWW '23, May 1–5, 2023, Austin, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-9416-1/23/04...\$15.00

<https://doi.org/10.1145/3543507.3583341>

Since the unfairness phenomenon has been ubiquitous in recommender systems and user-centric AI applications, how to define fairness measures and improve fairness to benefit all users is a trending research topic [24, 46]. Among all fairness metrics, group based fairness is widely accepted and adopted, which argues that the prediction performance does not discriminate particular user groups based on users' sensitive attributes (e.g., gender, race) [15, 20]. Given the basic ideas of group fairness, researchers proposed various group fairness measures and debiasing models to improve fairness, such as fairness regularization based models [3, 53], sensitive attribute disentangle models [58], and adversarial techniques to remove sensitive attributes [52]. E.g., researchers proposed different recommendation fairness regularization terms and added these terms in the loss function for collaborative filtering [53]. In summary, these works alleviate unfairness in the modeling process with the fixed training data and achieve better fairness results.

Table 1: An illustration of unfairness of recommendation accuracy performance on MovieLens dataset. We divide users into two subgroups based on the gender attribute. For each group, we calculate the distribution of clicked items of the corresponding group based on the training data, as well as the hit items on Top-K recommendation from two widely used recommendation models. Then, we measure the distribution differences of the two groups with JS divergence. We consider two typically recommendation models (BPR [49] and GCCF [40]). We consider the corrected hit from Top-K ranking list, denoted as “Top-20-Hit” and “Top-50-Hit”. The recommendation accuracy performance has larger JS divergence compared to the training data, showing recommendation models exacerbate unfairness from training data.

BPR [49]	JS divergence	GCCF [40]	JS divergence
Training data	0.1303	Training data	0.1303
Top-20-Hit	0.4842	Top-20-Hit	0.4879
Top-50-Hit	0.4349	Top-50-Hit	0.4229

In fact, researchers agree that unfairness in machine learning mainly comes from two processes. Firstly, the collected historical data shows imbalanced distribution among different user groups or reflects the real-world discrimination [9, 30]. After that, algorithms that learn the typical patterns amplify imbalance or bias inherited from training data and hurt the minority group [48, 52]. To show whether current recommendation algorithms amplify unfairness, let us take the MovieLens dataset as an example (detailed data description is shown in Section 5). We divide users into two subgroups based on the sensitive attribute *gender*. Given the training data, for each subgroup, we calculate the distribution over all items based on clicked items of users in this group. Then, we measure the distribution differences of the two sub user groups with Jensen-Shannon (JS) divergence [44]. After employing recommendation algorithms for Top-K recommendation, we measure the distribution difference of hit items of the two user groups. As can be seen from Table 1, the training data shows the preference distributions of the two groups are different. Both of the two recommendation models (i.e., BPR [49] and GCCF [40]) show larger group divergences of recommendation accuracy results compared to the divergence value of the training data. As a result, these two groups receive

different benefits from the recommendation algorithms. Since all recommendation models rely on the training data for model optimization, comparing with the huge works on model-level fairness research, how to consider and improve recommendation fairness from the data perspective is equally important.

When considering fairness from data perspective, some previous works proposed data resampling or data reweighting techniques to change data distribution [30, 48]. As a particular group of users (e.g., females) are underrepresented in the training data, BN-SLIM is designed to resample a balanced neighborhood set of males and females for neighborhood-based recommendation [9]. Besides, given a specific fairness measure in recommendation, researchers made attempts to add virtual users via a unified optimization framework [48]. These previous works show the possibility of improving fairness from the data perspective. However, they either suffered from applying to specific recommendation models or relied on a specific predefined fairness metric, limiting the generality to transferring to current RS architecture.

In this paper, we study the problem of designing a model-agnostic framework to improve recommendation fairness from data augmentation perspective. Given the original training data of user-item implicit feedback, we argue that the augmented training data are better balanced among different user groups, such that RS algorithms could better learn preferences of different user groups and avoid neglecting preferences of the minority groups. As a result, we argue the augmented data should satisfy the following hypotheses: for any user's two kinds of preference (a click record or a non-click record) in one group, there is another user in the remaining group (users with opposite sensitive attribute value) that has a similar item preference. With these hypotheses, we propose a general framework of Fairness-aware Data Augmentation (FDA) to generate “fake” data that complement the original training data. We design a bi-level optimization function with the inner and outer loop to optimize FDA. The proposed FDA is model-agnostic and can be easily integrated to any embedding-based recommendation. Besides, FDA does not rely on any predefined specific fairness metrics. Finally, extensive experiments on two real-world datasets clearly demonstrate the superiority of our proposed framework.

2 RELATED WORK

Fairness Discovery and Measures. As machine learning technologies have become a vital part of our daily lives with a high social impact, fairness issues are concerned and raised [33, 46]. Fairness refers to not discriminating against individuals or user groups based on sensitive user attributes [19, 42]. One increasing requirement is how to define and measure fairness. Current fairness metrics can be categorized into individual fairness and group fairness. Individual fairness refers to producing similar predictions for similar individuals, and group fairness argues not discriminating a particular user group based on the sensitive attribute [5, 7, 52, 54]. Among all group fairness metrics, Demographic Parity (DP) and Equality of Opportunity (EO) are widely accepted [12, 25]. DP requires that each particular demographic group has the same proportion of receiving a particular outcome [19, 56]. According to specific tasks, researchers have proposed specific demographic parity measures [10, 35, 56]. A notable disadvantage of demographic

parity lies in ignoring the natural differences across groups. To this end, EO is proposed for an equal proportion of receiving a particular outcome conditioned on the real outcome [2, 25, 47]. In other words, a model is fair if the predicted results and sensitive attributes are independently conditioned on the real outcome [25].

Fairness aware Models. Building on the mathematical fairness metrics, some researchers design task-oriented model structures to meet fairness requirements. Current model based fairness approaches can be classified into regularization based approaches [1, 18, 21, 29], causal based approaches [36, 38, 50], adversarial learning based approaches [4, 52, 57], and so on. These methods ensure that the outcome of models can meet fairness requirements, and the modification of model structures heavily relies on specific fairness definitions. E.g., regularization-based approaches add one fairness constraint to achieve a specific fairness at a time. Researchers need to design various constraints for achieving different fairness requirements [3, 18, 55].

Different from alleviating unfairness in the modeling process, some researchers attempted to solve fairness problems from the data perspective. Early works tried to directly remove the sensitive attribute from the training data [11, 34]. Some researchers argue that unfairness comes from the imbalanced training data of the protected and unprotected groups, and employ bagging and balance groups in each bag to build stratified training samples [30]. Besides, perturbation approaches change the training data distribution with some prior assumptions of sensitive attributes, input features and labels. After that, a perturbed distribution of disadvantaged groups is used to mitigate performance disparities [22, 32, 51]. Most of these data modification approaches are designed for classification tasks with abundant data samples, in which each data sample is independent. In CF based recommender systems, users and items are correlated with sparse interaction data. Therefore, current approaches of modifying data distribution in the classification task could not be easily adapted for the recommendation task with limited observed user behavior data. Recently, the authors [48] propose two metrics that capture the polarization and unfairness in recommendation. The framework needs to take one of the proposed metrics as the optimization direction, in order to generate corresponding antidote new user profiles for the selected metric. By proposing the concept of a balanced neighborhood, the authors [9] borrow the idea of data sampling and design corresponding regularization to control neighbor distribution of each user. The regularization is applied to the sparse linear method (SLIM) to improve the outcome fairness [9]. The above two models explore the possibility of improving recommendation fairness by changing the training data distribution. However, they either need to define/introduce specific fairness metrics or are only suitable for a particular kind of recommendation model with well-designed heuristics. Therefore, the problem of how to design a general fairness framework from the data perspective that is suitable for different recommendation backbones and multiple fairness metrics is still under explored.

3 PRELIMINARY

In a recommender system, there are two sets of entities: a user set U ($|U| = M$) and an item set V ($|V| = N$), we denote the user-item interaction matrix as $R = [r_{uv}]_{M \times N}$. We consider the common

implicit feedback scenario. If user u has clicked item v , then $r_{uv} = 1$ indicates user u likes the item v , otherwise $r_{uv} = 0$. For each user u , her clicked item set is denoted as $R_u = \{v : r_{uv} = 1\}$.

As most modern recommender systems are built on embedding based architecture, we focus on embedding based recommendation models. Generally speaking, there are two key components for recommendation. First, a recommendation model Rec employs an encoder Enc to project users and items into the embedding space, formulated as $\mathbf{E} = Enc(U, V) = [\mathbf{e}_1, \dots, \mathbf{e}_u, \dots, \mathbf{e}_v, \dots, \mathbf{e}_{M+N}]$, where \mathbf{e}_u is user u 's embedding, \mathbf{e}_v is item v 's embedding. Then, the predicted preference \hat{r}_{uv} can be calculated with $\hat{r}_{uv} = \mathbf{e}_u^T \mathbf{e}_v$. Learning high-quality user and item embeddings has become the key of modern recommender systems. There are two typical classes of embedding approaches: the classical matrix factorization models [45, 49] and neural graph-based models [27, 40].

Given a binary sensitive attribute $a \in \{0, 1\}$, a_u denotes user u 's attribute value. We divide the user set into two subsets: G_0 and G_1 . If $a_u = 0$, then $u \in G_0$, otherwise $u \in G_1$. Please note that as we do not focus on any specific recommendation models, we assume the embedding based recommendation models are available, such as matrix factorization models [45] or neural graph-based models [40]. Our goal is to improve recommendation fairness with relatively high accuracy from the data augmentation perspective.

Table 2: Mathematical Notations

Notations	Description
U, V	userset $ U = M$, itemset $ V = N$
$a \in \{0, 1\}$	a binary sensitive attribute
G_0, G_1	user group $\begin{cases} u \in G_0 & \text{if } a_u = 0 \\ u \in G_1 & \text{if } a_u = 1 \end{cases}$
u_0, u_1	users, $u_0 \in G_0, u_1 \in G_1$
$i_0, i_1, \bar{j}_0, \bar{j}_1$	real items
$\tilde{i}_0, \tilde{i}_1, \tilde{j}_0, \tilde{j}_1$	fake items
r_{ui}, r_{uj}	real positive data, real negative data
$\tilde{r}_{ui}, \tilde{r}_{uj}$	fake positive data, fake negative data

4 THE PROPOSED FRAMEWORK

In this section, we first introduce two hypotheses in our proposed FDA framework. Then, we show how to optimize two hypotheses given a recommendation model.

4.1 Hypotheses for Generating Fake Data

We argue that the augmented data should be balanced between two user groups, such that RS could learn latent preferences of different groups without neglecting the minority group. Since users have two kinds of behaviors, i.e., click (positive behavior) and non-click (negative behavior), we corresponding propose two hypotheses to augment data towards balanced distribution. In other words, for each behavior data (u_0, i, r) of user u_0 ($u_0 \in G_0$), item i , and the implicit feedback value r , we hope there is a user in the remaining group G_1 that shows the same preference value r to a similar item. Under this assumption, we can improve data balance of different user groups by generating fake behavior data that complements the training data. Specifically, the first hypothesis focuses on the positive behavior among two groups (G_0 and G_1).

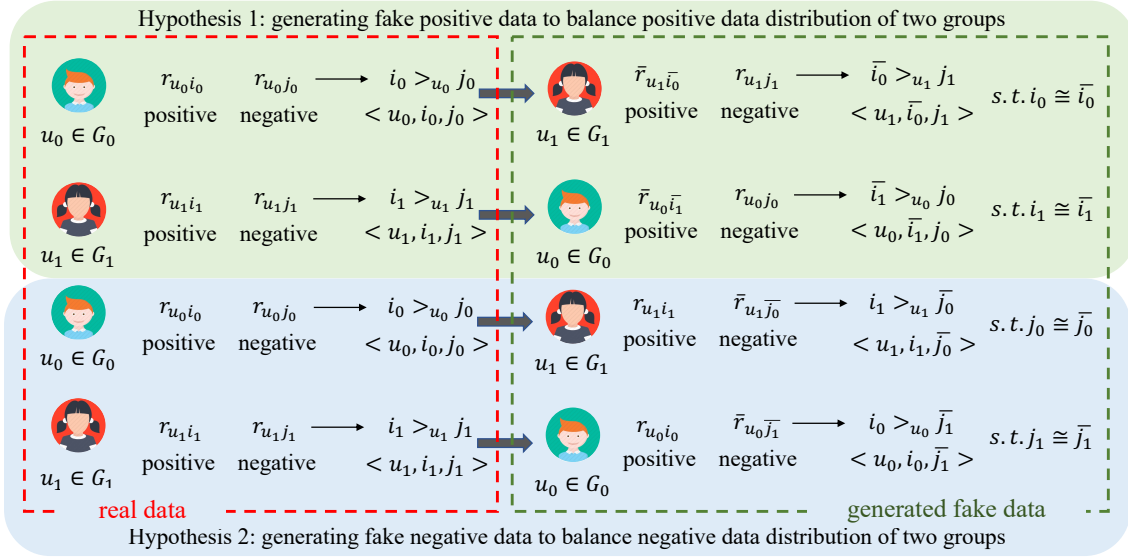


Figure 1: The overall architecture of the proposed FDA framework. By not changing the recommendation model, our key idea is to improve data balance between different user groups via data augmentation. For each real data (u_0, i, r) of user u_0 ($u_0 \in G_0$), item i , and the implicit feedback value r of this user-item pair from the training data, we hope there is a user in the remaining group G_1 that shows the same preference value r to a similar item \bar{i} . As the detailed implicit feedback value r can be positive or unknown, we turn the above idea into two hypotheses. In the upper part of the figure, according to Hypothesis 1, we generate two kinds of fake positive examples ($\bar{r}_{u_0 \bar{i}_1}$ and $\bar{r}_{u_1 \bar{i}_0}$). Correspondingly, based on the generated fake data, we can obtain fake records $(\langle u_0, \bar{i}_1, j_0 \rangle, \langle u_1, \bar{i}_0, j_1 \rangle)$. As shown in the lower part of the figure, according to Hypothesis 2, we generate two kinds of fake negative data ($\bar{r}_{u_0 \bar{j}_1}$ and $\bar{r}_{u_1 \bar{j}_0}$), and obtain fake records $(\langle u_0, i_0, \bar{j}_1 \rangle, \langle u_1, i_1, \bar{j}_0 \rangle)$.

Hypothesis 1. Assume that there is user $u_0 \in G_0$, and item i_0 is one of user u_0 's clicked items $i_0 \in R_{u_0}$. There should exist user $u_1 \in G_1$ that also clicks a similar item $\bar{i}_0 \approx i_0$. Similarly, if user $u_1 \in G_1$ has clicked item $i_1 \in R_{u_1}$, user $u_0 \in G_0$ should also click a similar item $\bar{i}_1 \approx i_1$. This hypothesis can be formulated as follows:

$$\forall u_0 \in G_0, u_1 \in G_1, \text{ if } r_{u_0 i_0} = 1, \text{ then } \bar{r}_{u_1 \bar{i}_0} = 1; \text{ where } i_0 \approx \bar{i}_0. \quad (1)$$

$$\forall u_1 \in G_1, u_0 \in G_0, \text{ if } r_{u_1 i_1} = 1, \text{ then } \bar{r}_{u_0 \bar{i}_1} = 1; \text{ where } i_1 \approx \bar{i}_1. \quad (2)$$

In the above equations, $\bar{r}_{u_1 \bar{i}_0}$ and $\bar{r}_{u_0 \bar{i}_1}$ are fake interactions data that does not appear in the training data.

The second hypothesis focuses on the negative behavior among two groups (G_0 and G_1):

Hypothesis 2. If user $u_0 \in G_0$ does not click item $j_0 \in V - R_{u_0}$, there should also exist user $u_1 \in G_1$ that does not click similar item $\bar{j}_0 \approx j_0$. Correspondingly, if user $u_1 \in G_1$ does not click item $j_1 \in R_{u_1}$, user $u_0 \in G_0$ should not click similar item $\bar{j}_1 \approx j_1$. This hypothesis can be formulated as follows:

$$\forall u_0 \in G_0, u_1 \in G_1, \text{ if } r_{u_0 j_0} = 0, \text{ then } \bar{r}_{u_1 \bar{j}_0} = 0; \text{ where } j_0 \approx \bar{j}_0. \quad (3)$$

$$\forall u_1 \in G_1, u_0 \in G_0, \text{ if } r_{u_1 j_1} = 0, \text{ then } \bar{r}_{u_0 \bar{j}_1} = 0; \text{ where } j_1 \approx \bar{j}_1. \quad (4)$$

In the above equations, $\bar{r}_{u_1 \bar{j}_0}$ and $\bar{r}_{u_0 \bar{j}_1}$ are fake interactions.

By employing Hypothesis 1 & 2, we can adjust the training data to make sure that the augmented training data are balanced for different groups.

4.2 Optimization For Fake Data

After generating fake data, we focus on augmenting the original training data for implicit feedback based recommendation. In implicit feedback based recommendation, pairwise learning has been widely used [27, 40, 49]. For each user u , if item i is clicked by user u , and item j is not clicked by user u , then this clicked item i is more relevant than a non-clicked item j , which can be formulated as $i >_u j$. As a result, the clicked item i should be assigned higher prediction value compared to the predicted preference of any non-clicked item j :

$$\forall i \in R_u, j \in V - R_u : \quad \hat{r}_{ui} > \hat{r}_{uj} \quad (5)$$

$$\text{i.e., } \mathbf{e}_u^T \mathbf{e}_i > \mathbf{e}_u^T \mathbf{e}_j. \quad (6)$$

Based on users' two sensitive attribute values and two kinds of implicit feedback behaviors, we can obtain four types of interactions (i.e., positive interactions for users from two groups: $r_{u_0 i_0}$ and $r_{u_1 i_1}$, negative interactions for users from two groups: $r_{u_0 j_0}$ and $r_{u_1 j_1}$) to support Hypothesis 1 and Hypothesis 2. The original training data contains two types of interactions as:

$$D_{real} = \{ \langle u_0, i_0, j_0 \rangle, \langle u_1, i_1, j_1 \rangle \}, \quad (7)$$

where user u_0 belongs to user group G_0 . Item i_0 is a positive feedback of user u_0 , and item j_0 is a negative feedback of user u_0 . Similarly, for user $u_1 \in G_1$, item i_1 is a positive feedback and item j_1 is a negative feedback. Next, we introduce how to construct the remaining two types of fake interactions and optimize them.

Optimization of Hypothesis 1. Hypothesis 1 focuses on clicked items and encourages the positive behavior distribution of two

groups are balanced. For each D_{real} , we can generate the corresponding fake data D_{fake1} according to Hypothesis 1 as:

$$D_{real} = \{ \langle u_0, i_0, j_0 \rangle, \langle u_1, i_1, j_1 \rangle \}, \quad (8)$$

$$D_{fake1} = \{ \langle u_1, \bar{i}_0, j_1 \rangle, \langle u_0, \bar{i}_1, j_0 \rangle \}, \quad (9)$$

where $i_0 \approx \bar{i}_0$ and $i_1 \approx \bar{i}_1$. D_{fake1} denotes fake positive interaction data. Note that, given any two real positive behavior data ($r_{u_0 i_0}$ and $\bar{r}_{u_1 i_1}$), D_{fake1} also contains two corresponding fake positive behavior data ($\bar{r}_{u_1 \bar{i}_0}$ and $\bar{r}_{u_0 \bar{i}_1}$). Therefore, we have the following expressions on D_{fake1} :

$$\bar{i}_0 >_{u_1} j_1 \text{ and } \bar{i}_1 >_{u_0} j_0. \quad (10)$$

Similar to Eq.(5), we can formulate Eq.(10) with the following optimization goal:

$$\hat{r}_{u_1 \bar{i}_0} > \hat{r}_{u_1 j_1} \text{ and } \hat{r}_{u_0 \bar{i}_1} > \hat{r}_{u_0 j_0}, \quad (11)$$

which can also be calculated as follows:

$$\mathbf{e}_{u_1}^T \bar{\mathbf{e}}_{i_0} > \mathbf{e}_{u_1}^T \mathbf{e}_{j_1} \text{ and } \mathbf{e}_{u_0}^T \bar{\mathbf{e}}_{i_1} > \mathbf{e}_{u_0}^T \mathbf{e}_{j_0}. \quad (12)$$

Optimization of Hypothesis 2. Hypothesis 2 focuses on non-clicked items and encourages the non-click behavior of two user groups to be balanced. For each triplet from D_{real} , we can generate the corresponding fake behavior data D_{fake2} according to Hypothesis 2:

$$D_{real} = \{ \langle u_0, i_0, j_0 \rangle, \langle u_1, i_1, j_1 \rangle \}, \quad (13)$$

$$D_{fake2} = \{ \langle u_1, i_1, \bar{j}_0 \rangle, \langle u_0, i_0, \bar{j}_1 \rangle \}, \quad (14)$$

where $j_0 \approx \bar{j}_0$ and $j_1 \approx \bar{j}_1$. Therefore, we have the following goal on D_{fake2} :

$$i_1 >_{u_1} \bar{j}_0 \text{ and } i_0 >_{u_0} \bar{j}_1, \quad (15)$$

Similarly, we can turn the above goal into optimization functions as:

$$\hat{r}_{u_1 i_1} > \hat{r}_{u_1 \bar{j}_0} \text{ and } \hat{r}_{u_0 i_0} > \hat{r}_{u_0 \bar{j}_1}, \quad (16)$$

which can be calculated as follows:

$$\mathbf{e}_{u_1}^T \mathbf{e}_{i_1} > \mathbf{e}_{u_1}^T \bar{\mathbf{e}}_{j_0} \text{ and } \mathbf{e}_{u_0}^T \mathbf{e}_{i_0} > \mathbf{e}_{u_0}^T \bar{\mathbf{e}}_{j_1}. \quad (17)$$

We construct corresponding fake data for each hypothesis, then we integrate all fake data from D_{fake1} and D_{fake2} :

$$\begin{aligned} D_{fake} &= \{D_{fake1}, D_{fake2}\} \\ &= \{ \langle u_1, \bar{i}_0, j_1 \rangle, \langle u_0, \bar{i}_1, j_0 \rangle, \langle u_1, i_1, \bar{j}_0 \rangle, \langle u_0, i_0, \bar{j}_1 \rangle \}. \end{aligned} \quad (18)$$

With the implicit feedback, Bayesian Personalized Ranking (BPR) is widely used for learning the pairwise based optimization function [45, 49]. We also adopt BPR loss to optimize the fake data generation process:

$$\begin{aligned} \min \mathcal{L}_{fake} &= - \sum_{\langle u, i, j \rangle \in D_{fake}} \ln(\sigma(\hat{r}_{ui} - \hat{r}_{uj})) \\ &= - \sum_{\langle u_1, \bar{i}_0, j_1 \rangle \in D_{fake1}} \ln(\sigma(\hat{r}_{u_1 \bar{i}_0} - \hat{r}_{u_1 j_1})) \\ &\quad - \sum_{\langle u_0, \bar{i}_1, j_0 \rangle \in D_{fake1}} \ln(\sigma(\hat{r}_{u_0 \bar{i}_1} - \hat{r}_{u_0 j_0})) \\ &\quad - \sum_{\langle u_1, i_1, \bar{j}_0 \rangle \in D_{fake2}} \ln(\sigma(\hat{r}_{u_1 i_1} - \hat{r}_{u_1 \bar{j}_0})) \\ &\quad - \sum_{\langle u_0, i_0, \bar{j}_1 \rangle \in D_{fake2}} \ln(\sigma(\hat{r}_{u_0 i_0} - \hat{r}_{u_0 \bar{j}_1})). \end{aligned} \quad (19)$$

The key challenge in Eq.(19) lies in estimating the fake data D_{fake} . A direct approach is to define a similarity function among items, and then find similar items within a predefined threshold. However, with sparse user-item click behavior data, directly computing item similarities from user-item interaction matrix is not only time-consuming, but also not accurate.

We propose to find similar items (\bar{i} and \bar{j}) from continuous embedding space. An intuitive idea is to utilize the well-trained embeddings from the recommendation model. Note that, recommender systems transform user/item ID to continuous embedding space: $\mathbf{E} = \text{Enc}(U, V)$. We therefore define and find similar items (\bar{i} and \bar{j}) based on continuous embedding space \mathbf{E} . In order to satisfy the similarity requirement, first of all, the similar items need to lie within the original embedding distribution. Otherwise, it will seriously affect the recommendation accuracy. Inspired by adversarial examples and poisoning attacks [6], we employ a non-random perturbation δ to generate similar items in continuous embedding \mathbf{E} .

For each item v , we add small random noise δ_v to the item original embedding \mathbf{e}_v , then we can construct the corresponding similar item embedding $\bar{\mathbf{e}}_v$. The similar item can be formulated as follows:

$$\bar{\mathbf{e}}_v = \mathbf{e}_v + \delta_v, \quad \text{where } \|\delta_v\| \leq \epsilon \quad (20)$$

The noise δ_v is bounded in a small range ϵ , and the operator \leq enforces the constraint for each dimension of δ_v . Since δ_v is a small “unseen” noise, it is natural that $\bar{\mathbf{e}}_v$ is similar to \mathbf{e}_v and $\bar{\mathbf{e}}_v$ also lies within the original embedding distribution. This method can meet the requirements of Hypothesis 1 and Hypothesis 2. By combining the similar item requirement in Eq.(20) and the optimization of the fake data in Eq.(19), the loss function on fake data can be changed into the embedding form as:

$$\begin{aligned} \min_{\Theta} \mathcal{L}_{fake} &= - \sum_{D_{fake}} \ln(\sigma(\mathbf{e}_{u_0}^T (\mathbf{e}_{i_1} + \delta_{i_1}) - \mathbf{e}_{u_0}^T \mathbf{e}_{j_0})) \\ &\quad - \sum_{D_{fake}} \ln(\sigma(\mathbf{e}_{u_1}^T (\mathbf{e}_{i_0} + \delta_{i_0}) - \mathbf{e}_{u_1}^T \mathbf{e}_{j_1})) \\ &\quad - \sum_{D_{fake}} \ln(\sigma(\mathbf{e}_{u_0}^T \mathbf{e}_{i_0} - \mathbf{e}_{u_0}^T (\mathbf{e}_{j_1} + \delta_{j_1}))) \\ &\quad - \sum_{D_{fake}} \ln(\sigma(\mathbf{e}_{u_1}^T \mathbf{e}_{i_1} - \mathbf{e}_{u_1}^T (\mathbf{e}_{j_0} + \delta_{j_0}))), \end{aligned} \quad (21)$$

where $\delta_{i_0}, \delta_{i_1}, \delta_{j_0}, \delta_{j_1}$ respectively denote the small noises adding to corresponding item embeddings, and Θ denotes the combination of all small noises in the fake data. With fixed embeddings \mathbf{E} from any recommender model, we optimize the Θ to generate fake data.

4.3 The Overall Bi-Level Optimization

After generating fake training data, we intend to integrate these fake data with original training data as augmented data. We could not use all the fake data for data augmentation as too many fake data records would dramatically modify the original data distribution, leading to decrease accuracy. To better trade-off, we develop a random mask operation to inject fake data. The mask operation can be formulated as follows:

$$m_v = \begin{cases} 1, & \text{selected} \\ 0, & \text{else} \end{cases} \quad \text{and} \quad \sum_{v=1}^N m_v \leq \text{Max}_{mask}, \quad (22)$$

where Max_{mask} denotes the maximum number of selected fake data in each update. Obviously, Max_{mask} should be less than the number of items N . Since only part of the fake data is selected for training, the impact on recommendation accuracy can be controlled.

After fake data records are selected by the mask operation, we combine the selected fake data D_{fake} and original training data D_{real} to construct the augmented data D_{data} . Next, we retrain the recommendation model on the augmented data. This training process of a recommender model is the same as the recommender model except that the data input is augmented. In short, the overall training process of FDA involves two iterative steps: generating fake data based on the previous recommendation embeddings and training recommendation models given updated fake data. The two iterative steps can be combined by solving the following bi-level optimization problem:

$$\min_E \min_{\Theta} \mathcal{L} = \sum_{D_{aug}} -\ln(\sigma([m_{i_0} * \bar{r}_{u_1 \bar{i}_0} + (1 - m_{i_0}) * \hat{r}_{u_1 i_0}] - \hat{r}_{u_1 j_1})) \\ -\ln(\sigma([m_{i_1} * \bar{r}_{u_0 \bar{i}_1} + (1 - m_{i_1}) * \hat{r}_{u_0 i_1}] - \hat{r}_{u_0 j_0})) \\ -\ln(\sigma(\hat{r}_{u_1 i_1} - [m_{j_0} * \bar{r}_{u_1 \bar{j}_0} + (1 - m_{j_0}) * \hat{r}_{u_1 j_0}])) \\ -\ln(\sigma(\hat{r}_{u_0 i_0} - [m_{j_1} * \bar{r}_{u_0 \bar{j}_1} + (1 - m_{j_1}) * \hat{r}_{u_0 j_1}]))), \quad (23)$$

where $D_{aug} = D_{real} \cup D_{fake}$ is the augmented data. From this formulation, we can observe that our proposed FDA involves two levels of optimization.

4.4 Discussion

Model Analysis. Given the above bi-level optimization process of FDA, we now analyze why FDA can achieve a better balance between recommendation accuracy and fairness without changing the recommendation architecture (i.e., the outer loop that updates recommendation parameters). In the inner minimization loop, our key idea is to encourage for each triple behavior of the user, there is a user in the remaining group that shows the same preference value to a similar item. As users have two kinds of preferences (i.e., click and non-click), the key idea turns to two hypotheses for generating fake data. Therefore, the augmented data that contains both the training data and the fake data are more balanced compared to the original training data. For the original training data, the recommendation results are unfairer as recommendation results would amplify unbalance inherited from the input data. With more balanced data, the recommender models can output better fairness metrics even though the recommendation part does not model fairness. Besides, as we constrain the fake data generated by adding a small random noise, and the fake data is controlled by a limited ratio with random mask operation, the recommendation accuracy can also be guaranteed. Compared to other data balance or data augmentation based fairness enhanced recommendation models, FDA shows the advantage of generally applicable to embedding based recommendation backbones and does not need to predefine a specific group fairness metric.

Extension to Multiple Sensitive Values. In FDA, similar as many fairness aware approaches, we start with the binary sensitive attribute values [9, 25, 53]. When dealing with a sensitive attribute with K ($K > 2$) values, a naive extension to multiple sensitive attribute values is to encourage that: for each user's one kind of

behavior to an item in one group, we encourage that there are other users in each remaining group that show the same behavior to a similar item. Therefore, for each hypothesis, we can generate $K - 1$ fake data, and use the bi-level optimization with both original and the selected fake data with mask operation.

5 EXPERIMENTS

5.1 Experimental Setup

Datasets. We conduct experiments on two publicly available datasets: MovieLens [26] and LastFM [13, 52]. For MovieLens¹, we adopt the same data splitting strategy as previous works for fair recommendation [8, 52]. We treat the items that a user's rating is larger than 3 as positive feedback. Moreover, we randomly select 80% of records for training and the remaining 20% records for the test. LastFM is a large music recommendation dataset². We treat the items that a user plays as the positive feedback. To ensure the quality of the dataset, we use the 10-core setting to ensure that users (items) have at least 10 interaction records. We split the historical records into training, validation, and test parts with the ratio of 7:1:2. Besides the user-item interaction records, these two datasets also have the user profile data, including gender (two classes) for users. Similar as previous works, we treat gender as the sensitive attribute and divide users into two subgroups. The statistics of these two datasets are shown in Table 3.

Table 3: Statistics of the two datasets.

Datasets	Users	Items	Traning Records	Density
MovieLens	6,040	3,952	513,112	2.150%
LastFM	139,371	60,081	4,017,311	0.048%

Evaluation Metrics. Since we focus on the trade-off between fairness and recommendation accuracy, we need to evaluate two aspects and report the trade-off results. On the one hand, we employ two widely used ranking metrics for recommendation accuracy: HR [23] and NDCG [31] to evaluate the Top-K recommendation. Larger values of HR and NDCG mean better recommendation accuracy performance. On the other hand, we adopt two group fairness metrics: *Demographic Parity (DP)* [56] and *Equality of Opportunity (EO)* [25] to evaluate fairness. DP and EO evaluate group fairness from different aspects. For both fairness metrics, the smaller values mean better fairness results.

The following equation calculates DP measure:

$$DP = 1/N \sum_{v \in V} \frac{|\sum_{u \in G_0} \mathbb{1}_{v \in TopK_u} - \sum_{u \in G_1} \mathbb{1}_{v \in TopK_u}|}{\sum_{u \in G_0} \mathbb{1}_{v \in TopK_u} + \sum_{u \in G_1} \mathbb{1}_{v \in TopK_u}},$$

where G_0 denotes user group with sensitive attribute $a_u = 0$, and G_1 denotes user group with sensitive attribute $a_u = 1$. $TopK_u$ is Top-K ranked items for user u .

Please note that as DP forcefully requires similarly predicted results across different groups, it naturally neglects the natural preference differences of user clicked patterns [53]. EO is proposed to solve the limitation of DP [12, 24]. Specifically, it requires similar

¹<https://grouplens.org/datasets/movielens/>

²<http://ocelma.net/MusicRecommendationDataset/lastfm-360K.html>

Table 4: Recommendation accuracy and fairness performance on MovieLens with varying Top-K values. We compare all fairness-aware models, the best results are presented in bold font and the second best results are presented in underline. The performance improvement of our model against the best baseline is significant under the paired-t test.

Model	K=10				K=20				K=30				K=40				K=50			
	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓
BPR	0.2478	0.2492	0.6541	0.7158	0.2770	0.2498	0.6198	0.6868	0.3147	0.2600	0.6088	0.6803	0.3519	0.2720	0.5960	0.6746	0.3849	0.2830	0.5861	0.6575
GCCF	0.2607	0.2602	0.6391	0.7025	0.2913	0.2617	0.6182	0.6869	0.3301	0.2722	0.6011	0.6842	0.3708	0.2849	0.5856	0.6611	0.4059	0.2967	0.5740	0.6462
BPR_DP	0.2209	0.2241	0.6524	0.7039	0.2446	0.2229	0.6063	0.6802	0.2798	0.2324	0.6132	0.6662	0.3145	0.2433	0.5939	0.6672	0.3427	0.2530	0.5921	0.6565
BPR_EO	0.2259	0.2225	0.6510	0.7053	0.2507	0.2279	0.6194	0.6712	0.2862	0.2372	0.6162	0.6740	0.3204	0.2480	0.5976	0.6631	0.3529	0.2589	0.5950	0.6608
GCCF_DP	0.2416	0.2420	0.6532	0.7155	0.2691	0.2434	0.6164	0.6865	0.3063	0.2535	0.6053	0.6767	0.3426	0.2650	0.5951	0.6665	0.3761	0.2763	0.5869	0.6607
GCCF_EO	0.2407	<u>0.2428</u>	0.6479	0.7060	0.2698	<u>0.2437</u>	0.6211	0.6784	0.3068	0.2537	0.6066	0.6773	0.3428	0.2652	0.5980	0.6713	0.3748	0.2759	0.5844	0.6579
FairUser	0.2306	0.2262	0.6502	0.7375	0.2656	0.2318	0.6167	0.7162	0.3088	0.2448	0.5986	0.6927	0.3494	0.2582	0.5864	0.6826	0.3847	0.2705	0.5761	0.6680
BN-SLIM	0.2305	0.2287	0.6502	0.7349	0.2671	0.2334	0.6078	0.7061	0.3096	0.2457	0.5906	0.6958	0.3500	0.2589	0.5743	0.6816	0.3868	0.2716	0.5661	0.6747
FDA_BPR	0.2307	0.2226	0.6132	0.6969	0.2716	0.2321	0.5730	0.6562	0.3157	0.2460	0.5635	0.6488	0.3566	0.2598	0.5537	0.6358	0.3941	0.2729	0.5450	0.6224
FDA_NCF	<u>0.2401</u>	0.2322	<u>0.6093</u>	<u>0.6946</u>	<u>0.2800</u>	0.2367	<u>0.5763</u>	0.6644	<u>0.3208</u>	<u>0.2544</u>	<u>0.5681</u>	0.6464	<u>0.3601</u>	<u>0.2711</u>	<u>0.5554</u>	<u>0.6306</u>	<u>0.3972</u>	<u>0.2824</u>	<u>0.5445</u>	<u>0.6205</u>
FDA_GCCF	0.2476	0.2430	0.6036	0.6773	0.2857	0.2487	0.5786	<u>0.6593</u>	0.3290	0.2614	0.5682	0.6462	0.3709	0.2748	0.5510	0.6231	0.4075	0.2873	0.5438	0.6130

Table 5: Recommendation accuracy and fairness performance on LastFM with varying Top-K values. The performance improvement of our model against the best baseline is significant under the paired-t test.

Model	K=10				K=20				K=30				K=40				K=50			
	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓	HR↑	NDCG↑	DP↓	EO↓
BPR	0.1323	0.1291	0.6372	0.6636	0.1991	0.1602	0.6235	0.6450	0.2480	0.1794	0.6178	0.6410	0.2869	0.1933	0.6126	0.6361	0.3195	0.2042	0.6109	0.6348
GCCF	0.1361	0.1324	0.6321	0.6580	0.2038	0.1640	0.6220	0.6477	0.2533	0.1834	0.6143	0.6419	0.2929	0.1976	0.6105	0.6377	0.3259	0.2087	0.6065	0.6307
BPR_DP	0.1272	0.1231	0.6064	0.6451	0.1927	0.1535	0.5956	0.6318	0.2407	0.1724	0.5899	0.6251	0.2789	0.1861	0.5839	0.6215	0.3108	0.1968	0.5821	0.6192
BPR_EO	0.1292	0.1247	0.5933	0.6440	0.1958	0.1557	0.5738	0.6271	0.2444	0.1748	0.5637	0.6218	0.2830	0.1886	0.5577	0.6193	0.3155	0.1995	0.5541	0.6150
GCCF_DP	0.1281	0.1234	0.6285	0.6636	0.1950	0.1544	0.6108	0.6473	0.2442	0.1738	0.6046	0.6400	0.2832	0.1877	0.5986	0.6355	0.3163	0.1988	0.5977	0.6377
GCCF_EO	0.1295	0.1245	0.5912	0.6522	0.1969	0.1558	0.5681	0.6308	0.2461	0.1752	0.5530	0.6160	0.2855	0.1892	0.5471	0.6135	0.3176	0.2003	<u>0.5427</u>	0.6125
BN-SLIM	0.0986	0.0943	0.6183	0.6630	0.1546	0.1204	0.6017	0.6557	0.1970	0.1371	0.5962	0.6517	0.2318	0.1496	0.5911	0.6467	0.2620	0.1597	0.5847	0.6411
FDA_BPR	<u>0.1301</u>	<u>0.1268</u>	0.5604	0.5937	0.1965	<u>0.1577</u>	0.5535	0.5825	0.2455	<u>0.1770</u>	0.5524	0.5788	0.2848	0.1911	0.5479	0.5761	0.3180	0.2022	0.5457	0.5707
FDA_NCF	0.1300	0.1248	<u>0.5599</u>	<u>0.5931</u>	<u>0.1970</u>	0.1575	<u>0.5505</u>	<u>0.5801</u>	<u>0.2464</u>	<u>0.1770</u>	<u>0.5470</u>	<u>0.5748</u>	<u>0.2856</u>	<u>0.1913</u>	<u>0.5449</u>	<u>0.5720</u>	<u>0.3194</u>	<u>0.2025</u>	0.5432	<u>0.5700</u>
FDA_GCCF	0.1304	0.1268	0.5477	0.5789	0.1976	0.1580	0.5450	0.5693	0.2470	0.1774	0.5412	0.5688	0.2868	0.1917	0.5417	0.5676	0.3200	0.2029	0.5403	0.5641

predicted results across different groups conditioned on user real preferences. We calculate EO as follows:

$$EO = 1/N \sum_{v \in V} \frac{|\sum_{u \in G_0} \mathbb{1}_{v \in TopK_u} \mathbb{1}_{v \in Test_u} - \sum_{u \in G_1} \mathbb{1}_{v \in TopK_u} \mathbb{1}_{v \in Test_u}|}{\sum_{u \in G_0} \mathbb{1}_{v \in TopK_u} \mathbb{1}_{v \in Test_u} + \sum_{u \in G_1} \mathbb{1}_{v \in TopK_u} \mathbb{1}_{v \in Test_u}},$$

where $Test_u$ is items that user u clicks on the test data. Because not all predicted results have corresponding ground true labels, we calculate EO metric only on the testing data. A smaller EO means there is less unfairness, as the two groups receive similar recommendation accuracy.

Baselines. The baseline models can be divided into two categories: recommendation based models *BPR* [49] and *GCCF* [40], and the fairness-oriented models *BN-SLIM* [9], data augmentation model (*FairUser*) [48] and the fairness regularization based model [53]. *FairUser* adds virtual users to achieve fairness. As its optimization process is very time-consuming and needs to compute the full user-item matrix at each iteration, we only test *FairUser* on MovieLens dataset. Similar as the fairness regularization based model [53], we add different group fairness metrics (DP and EO) as regularization terms into recommendation based models. E.g., *BPR_EO* denotes treating EO as the regularization term for the base recommendation model of *BPR*.

Our proposed framework *FDA* can be applied on different recommendation backbones. We select *BPR* [49], *NCF* [28] and *GCCF* [40] as recommendation backbones show state-of-the-art performance. We use *FDA_BPR*, *FDA_NCF* and *FDA_GCCF* to denote the variants of our proposed framework with different recommendation backbones.

Parameter Setting. Our implementations are based on Pytorch-GPU 1.6.0. The embedding size is set as 64 for *FDA*. All the parameters are differentiable in the objective function, and we use the

Adam optimizer to optimize the model. In Eq.(23), the initial learning rate of Adam is 0.001 for the outer minimization and the inner minimization.

5.2 Overall Performance on Two Datasets

Table 4 and 5 report the overall results on two datasets. We have several observations from these two tables. *FDA* achieves the best performance from the two aspects: 1) improving EO and DP concurrently; 2) the trade-off between recommender accuracy and fairness.

First, when comparing each fairness metric, our proposed *FDA* outperforms other models on both DP and EO group fairness metrics. *FairUser* and *BN-SLIM* show worse performance than *BPR* (not considering fairness) on EO metric. *BPR_EO* and *GCCF_EO* can improve EO and DP concurrently. But, *BPR_EO* and *GCCF_EO* is worse than *FDA*. DP and EO reflect the group fairness from two aspects. *FDA* achieves good results of multiple group metrics on all datasets. Thus, *FDA* can well alleviate the unfairness problem and achieve better fairness performance. Second, when comparing the trade-off between recommender accuracy and fairness, all fairness-aware baselines perform worse on recommendation accuracy performance than *FDA*. In other words, all fairness-aware baselines cause a larger decrease in recommender accuracy. Therefore, while achieving better fairness performance, *FDA* has the least damage to accuracy on different datasets. We conclude that *FDA* can reach the best balance between accuracy and fairness. Third, no matter the base backbone model is *BPR*, *NCF* or *GCCF*, *FDA* can improve fairness and show better recommendation performance. *FDA_BPR* shows worse performance than *FDA_GCCF*. This is due to the fact that the base model (i.e., *BPR*) in *FDA_BPR* does not

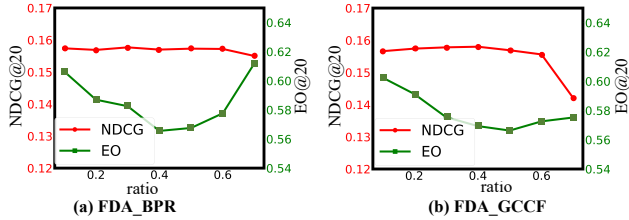
Table 6: Ablation study of the two modules of FDA: Hypothesis 1 and Hypothesis 2 .

Model	Hypothesis 1	Hypothesis 2	MovieLens				LastFM			
			HR@20	NDCG@20	DP@20	EO@20	HR@20	NDCG@20	DP@20	EO@20
FDA_BPR	✓	✗	0.2706	0.2462	0.6099	<i>0.6564</i>	0.1971	0.1613	0.6017	0.6202
FDA_BPR	✗	✓	0.2721	<i>0.2360</i>	<i>0.5841</i>	0.6682	0.1951	0.1552	<i>0.5585</i>	<i>0.6016</i>
FDA_BPR	✓	✓	<i>0.2716</i>	0.2321	0.5731	0.6562	<i>0.1965</i>	<i>0.1573</i>	0.5392	0.5656
FDA_GCCF	✓	✗	0.2650	0.2397	0.6078	<i>0.6610</i>	0.1850	<i>0.1511</i>	0.5905	0.6021
FDA_GCCF	✗	✓	0.2876	<i>0.2463</i>	<i>0.5824</i>	0.6696	<i>0.1906</i>	0.1499	<i>0.5647</i>	<i>0.5991</i>
FDA_GCCF	✓	✓	<i>0.2857</i>	0.2487	0.5786	0.6593	0.1976	0.1580	0.5450	0.5693

perform as well as the base graph embedding model GCCF. On the whole, for different recommendation backbone models and different experimental settings, FDA can effectively balance accuracy and fairness. This demonstrates the flexibility and effectiveness of FDA. Fourth, we observe the improvements of fairness on MovieLens are not so obvious as results on LastFM. In other words, eliminating unfairness on LastFM is easier than that on MovieLens. We guess a possible reason is that the sparsity and number of users are different on these two datasets.

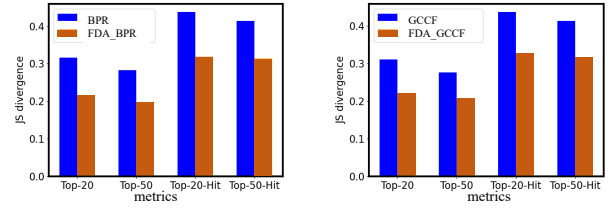
5.3 Model Analyses

Effects of fake data numbers. The setting of Max_{mask} plays an important role to control the maximum number of fake data. We conduct experiments on different Max_{mask} , as shown in Figure 2. Specifically, we show different ratios of maximum fake data and the number of all items (Max_{mask}/N). Since the fake data is similar to the real data, the fake data does not seriously affect the recommendation accuracy. When the ratio increases from 0.1 to 0.6, the recommended performance of FDA_BPR and FDA_GCCF does not decrease nearly. Among all ratios, we can find that FDA achieves a better balance on the fairness and recommender accuracy when the ratio equals 0.3 for FDA_BPR and 0.4 for FDA_GCCF on LastFM. The trend is similar on MovieLens. Due to the page length limit, the results on MovieLens are not reported.

**Figure 2: Performance under different ratios on LastFM.**

Measuring the distribution differences of recommendation results. As shown in Table 1, we employ the Jensen-Shannon (JS) divergence to measure differences of two user groups (i.e., G_0 and G_1) on the larger LastFM. We pay attention to Top-20 and Top-50 ranked items, denoted as “Top-20” and “Top-50”. We also consider the corrected Top-20 and Top-50 ranked items, denoted as “Top-20-Hit” and “Top-50-Hit” to measure the recommendation accuracy differences between the two groups. As shown in Figure 3, FDA can also improve fairness performance compared to its original recommendation backbone. Also, we observe that JS divergence of “Top-20” and “Top-50” is smaller than that of “Top-20-Hit” and “Top-50-Hit”. This is reasonable as “Top-K” divergence measures

ranking differences between the two groups and does not take recommendation accuracy into consideration. In contrast, “Top-K-Hit” measures the recommendation accuracy differences between the two groups. Nevertheless, FDA can improve these two metrics.

**Figure 3: The distribution differences of recommendation results with two recommendation backbones on LastFM.**

Ablation study. In this part, we investigate the effectiveness of each hypothesis: Hypothesis 1 and Hypothesis 2 of our proposed FDA framework. The results are illustrated in Table 6. From this table, we can obtain the following observations. First, each single hypothesis (Hypothesis 1 or Hypothesis 2) can help the model achieve comparable performance, indicating the usefulness of our proposed hypothesis. Second, compared with the performance of models with a single hypothesis, models with both of them (the entire FDA framework) have better performance, demonstrating the necessity of both hypotheses. As these two hypotheses consider different kinds of users’ click or non-click behavior, combining them together reaches the best performance.

6 CONCLUSION

In this paper, we studied the recommendation fairness issue from data augmentation perspective. Given the original training data, we proposed a FDA framework to generate fake user behavior data, in order to improve recommendation fairness. Specifically, given the overall idea of balanced data, we proposed two hypotheses to guide the generation of fake data. After that, we designed a bi-level optimization target, in which the inner optimization generates better fake data and the outer optimization finds recommendation parameters given the augmented data that comprises both the original training data and the generated fake data. Please note that, FDA can be applied to any embedding based recommendation backbones, and does not rely on any specific fairness metrics. Extensive experiments on two real-world datasets clearly showed FDA is effective to balance recommendation accuracy and fairness under different recommendation backbones.

7 ACKNOWLEDGEMENTS

This work is supported in part by grants from the National Key Research and Development Program of China (Grant No. 2021ZD0111802), the National Natural Science Foundation of China (Grant No. 72188101, 61932009, U1936219, U22A2094), Major Project of Anhui Province (Grant No. 202203a05020011), and the CCF-AFSG Research Fund (Grant No. CCF-AFSG RF20210006).

REFERENCES

- [1] Sina Aghaei, Mohammad Javad Azizi, and Phebe Vayanos. 2019. Learning optimal and fair decision trees for non-discriminative decision-making. In *AAAI*, Vol. 33. 1418–1426.
- [2] Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. 2021. Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research* 50, 1 (2021), 3–44.
- [3] Alex Beutel, Jilin Chen, Tulsee Doshi, Hai Qian, Li Wei, Yi Wu, Lukasz Heldt, Zhe Zhao, Lichan Hong, Ed H Chi, et al. 2019. Fairness in recommendation ranking through pairwise comparisons. In *SIGKDD*. 2212–2220.
- [4] Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi. 2017. Data decisions and theoretical implications when adversarially learning fair representations. In *FAT/ML*.
- [5] Asia J Biega, Krishna P Gummadi, and Gerhard Weikum. 2018. Equity of attention: Amortizing individual fairness in rankings. In *SIGIR*. 405–414.
- [6] Battista Biggio, Blaine Nelson, and Pavel Laskov. 2012. Poisoning attacks against support vector machines. In *ICML*. 1467–1474.
- [7] Reuben Binns. 2020. On the apparent conflict between individual and group fairness. In *FACCT*. 514–524.
- [8] Avishek Joey Bose and William Hamilton. 2019. Compositional fairness constraints for graph embeddings. In *ICML*. 715–724.
- [9] Robin Burke, Nasim Sonboli, and Aldo Ordonez-Gauger. 2018. Balanced neighborhoods for multi-sided fairness in recommendation. In *FACCT*. 202–214.
- [10] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. 2009. Building classifiers with independency constraints. In *ICDMW*. 13–18.
- [11] Toon Calders and Sicco Verwer. 2010. Three naive bayes approaches for discrimination-free classification. *DMKD* 21, 2 (2010), 277–292.
- [12] Simon Caton and Christian Haas. 2020. Fairness in machine learning: A survey. *arXiv preprint arXiv:2010.04053* (2020).
- [13] Óscar Celma Herrada et al. 2009. *Music recommendation and discovery in the long tail*. Universitat Pompeu Fabra.
- [14] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *RecSys*. 191–198.
- [15] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *ITCS*. 214–226.
- [16] Michael D Ekstrand, Mucun Tian, Ion Madrazo Azpiazu, Jennifer D Ekstrand, Oghenemaro Anuyah, David McNeill, and Maria Soledad Pera. 2018. All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. In *FAT*. 1–15.
- [17] Michael D Ekstrand, Mucun Tian, Mohammed R Imran Kazi, Hoda Mehrpouyan, and Daniel Kluver. 2018. Exploring author gender in book rating and recommendation. In *RecSys*. 242–250.
- [18] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In *SIGKDD*. 259–268.
- [19] Pratik Gajane and Mykola Pechenizkiy. 2017. On formalizing fairness in prediction with machine learning. *arXiv preprint arXiv:1710.03184* (2017).
- [20] Sahin Cem Geyik, Stuart Ambler, and Krishnamurthy Kenthapadi. 2019. Fairness-aware ranking in search & recommendation systems with application to linkedin talent search. In *SIGKDD*. 2221–2231.
- [21] Naman Goel, Mohammad Yaghini, and Boi Faltings. 2018. Non-discriminatory machine learning through convex fairness criteria. In *AAAI*.
- [22] Paula Gordaliza, Eustasio Del Barrio, Gamboa Fabrice, and Jean-Michel Loubes. 2019. Obtaining fairness using optimal transport theory. In *ICML*. 2357–2365.
- [23] Asela Gunawardana and Guy Shani. 2009. A survey of accuracy evaluation metrics of recommendation tasks. *JMLR* 10, 12 (2009).
- [24] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. In *NIPS*. 3315–3323.
- [25] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. *NIPS* 29 (2016), 3315–3323.
- [26] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *TIIS* 5, 4 (2015), 1–19.
- [27] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *SIGIR*. 639–648.
- [28] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the International Conference on World Wide Web*. 173–182.
- [29] Lingxiao Huang and Nisheeth Vishnoi. 2019. Stable and fair classification. In *ICML*. 2879–2890.
- [30] Vasileios Ioannidis, Besnik Fetahu, and Eirini Ntoutsi. 2019. Fae: A fairness-aware ensemble framework. In *Big Data*. 1375–1380.
- [31] Kalervo Järvelin and Jaana Kekäläinen. 2017. IR evaluation methods for retrieving highly relevant documents. In *SIGIR*, Vol. 51. ACM New York, NY, USA, 243–250.
- [32] Heinrich Jiang and Ofir Nachum. 2020. Identifying and correcting label bias in machine learning. In *AISTATS*. 702–712.
- [33] Anna Jobin, Marcello Lenca, and Effy Vayena. 2019. The global landscape of AI ethics guidelines. *Nature Machine Intelligence* 1, 9 (2019), 389–399.
- [34] Faisal Kamiran and Toon Calders. 2012. Data preprocessing techniques for classification without discrimination. *KAIS* 33, 1 (2012), 1–33.
- [35] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2012. Fairness-aware classifier with prejudice remover regularizer. In *ECML PKDD*. 35–50.
- [36] Niki Kilbertus, Mateo Rojas-Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. 2017. Avoiding discrimination through causal reasoning. In *NIPS*. 656–666.
- [37] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* 42, 8 (2009), 30–37.
- [38] Matt Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. In *NIPS*. 4069–4079.
- [39] Anja Lambrecht and Catherine Tucker. 2019. Algorithmic bias? An empirical study of apparent gender-based discrimination in the display of STEM career ads. *Management science* 65, 7 (2019), 2966–2981.
- [40] Chen Lei, Wu Le, Hong Richang, Zhang Kun, and Wang Meng. 2020. Revisiting Graph based Collaborative Filtering: A Linear Residual Graph Convolutional Network Approach. In *AAAI*, Vol. 34. 27–34.
- [41] Yunqi Li, Hanxiong Chen, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2021. User-oriented Fairness in Recommendation. In *WWW*. 624–632.
- [42] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *CSUR* 54, 6 (2021), 1–35.
- [43] Alessandro B Melchiorre, Navid Rekabsaz, Emilia Parada-Cabaleiro, Stefan Brandl, Oleg Lesota, and Markus Schedl. 2021. Investigating gender fairness of recommendation algorithms in the music domain. *IPM* 58, 5 (2021), 102666.
- [44] ML Menéndez, JA Pardo, L Pardo, and MC Pardo. 1997. The jensen-shannon divergence. *Journal of the Franklin Institute* 334, 2 (1997), 307–318.
- [45] Andriy Mnih and Ruslan R Salakhutdinov. 2008. Probabilistic matrix factorization. In *NIPS*. 1257–1264.
- [46] Dino Pedreshi, Salvatore Ruggieri, and Franco Turini. 2008. Discrimination-aware data mining. In *SIGKDD*. 560–568.
- [47] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. 2017. On Fairness and Calibration. *NIPS* 30 (2017), 5680–5689.
- [48] Bashir Rastegarpanah, Krishna P Gummadi, and Mark Crovella. 2019. Fighting fire with fire: Using antidote data to improve polarization and fairness of recommender systems. In *WSDM*. 231–239.
- [49] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *UAI*. 452–461.
- [50] Ashudeep Singh and Thorsten Joachims. 2018. Fairness of exposure in rankings. In *SIGKDD*. 2219–2228.
- [51] Hao Wang, Berk Ustun, and Flavio Calmon. 2019. Repairing without retraining: Avoiding disparate impact with counterfactual distributions. In *ICML*. 6618–6627.
- [52] Le Wu, Lei Chen, Pengyang Shao, Richang Hong, Xiting Wang, and Meng Wang. 2021. Learning Fair Representations for Recommendation: A Graph-based Perspective. In *WWW*. 2198–2208.
- [53] Sirui Yao and Bert Huang. 2017. Beyond parity: Fairness objectives for collaborative filtering. In *NIPS*. 2921–2930.
- [54] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. 2017. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *WWW*. 1171–1180.
- [55] Meike Zehlke, Francesco Bonchi, Carlos Castillo, Sara Hajian, Mohamed Megahed, and Ricardo Baeza-Yates. 2017. Fa* ir: A fair top-k ranking algorithm. In *CIKM*. 1569–1578.
- [56] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *ICML*. 325–333.
- [57] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. 2018. Mitigating unwanted biases with adversarial learning. In *AAAI*. 335–340.
- [58] Ziwei Zhu, Xia Hu, and James Caverlee. 2018. Fairness-aware tensor-based recommendation. In *CIKM*. 1153–1162.