

# User-controllable Recommendation Against Filter Bubbles

Wenjie Wang<sup>1</sup>, Fuli Feng<sup>2\*</sup>, Liqiang Nie<sup>3</sup>, and Tat-Seng Chua<sup>1</sup>

<sup>1</sup>Sea-NExT Joint Lab, National University of Singapore,

<sup>2</sup>University of Science and Technology of China, <sup>3</sup>Shandong University  
{wenjiewang96,fulifeng93,nieliqiang}@gmail.com,dcscts@nus.edu.sg

## ABSTRACT

Recommender systems usually face the issue of filter bubbles: over-recommending homogeneous items based on user features and historical interactions. Filter bubbles will grow along the feedback loop and inadvertently narrow user interests. Existing work usually mitigates filter bubbles by incorporating objectives apart from accuracy such as diversity and fairness. However, they typically sacrifice accuracy, hurting model fidelity and user experience. Worse still, users have to passively accept the recommendation strategy and influence the system in an inefficient manner with high latency, e.g., keeping providing feedback (e.g., like and dislike) until the system recognizes the user intention.

This work proposes a new recommender prototype called *User-Controllable Recommender System* (UCRS), which enables users to actively control the mitigation of filter bubbles. Functionally, 1) UCRS can alert users if they are deeply stuck in filter bubbles. 2) UCRS supports four kinds of control commands for users to mitigate the bubbles at different granularities. 3) UCRS can respond to the controls and adjust the recommendations on the fly. The key to adjusting lies in blocking the effect of out-of-date user representations on recommendations, which contains historical information inconsistent with the control commands. As such, we develop a causality-enhanced *User-Controllable Inference* (UCI) framework, which can quickly revise the recommendations based on user controls in the inference stage and utilize counterfactual inference to mitigate the effect of out-of-date user representations. Experiments on three datasets validate that the UCI framework can effectively recommend more desired items based on user controls, showing promising performance *w.r.t.* both accuracy and diversity.

## CCS CONCEPTS

• **Information systems** → **Recommender systems.**

## KEYWORDS

User-controllable Recommender Systems, Counterfactual Inference, Filter Bubbles, Causal Recommendation

\* Corresponding author: Fuli Feng. This research is supported by the Sea-NExT Joint Lab, the National Natural Science Foundation of China (No. U1936203), and the Major Basic Research Project of Natural Science Foundation of Shandong Province (No. ZR2021ZD15).

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

SIGIR '22, July 11–15, 2022, Madrid, Spain

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-8732-3/22/07...\$15.00

<https://doi.org/10.1145/3477495.3532075>

## ACM Reference Format:

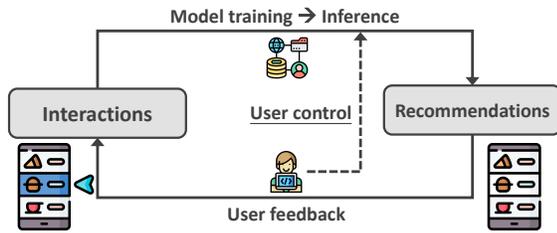
Wenjie Wang, Fuli Feng, Liqiang Nie, and Tat-Seng Chua. 2022. User-controllable Recommendation Against Filter Bubbles. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*, July 11–15, 2022, Madrid, Spain. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3477495.3532075>

## 1 INTRODUCTION

Recommender systems become increasingly important to provide personalized information filtering services in this information explosion era [35]. A *de facto* standard for building recommender systems is mining user interests from user features (e.g., gender and age) and historical interactions (e.g., click). Due to merely fitting the data, recommender systems typically face filter bubble issues: continually recommending many homogeneous items, isolating users from diverse contents [19, 33, 48]. For example, if a user has clicked many micro-videos to learn making coffee, the system may continuously recommend similar micro-videos from different uploaders, occupying the opportunities of other informative videos such as hot news. Worse still, due to the feedback loop as shown in Figure 1, filter bubbles might gradually become severer, narrow users' interests, and even intensify the segregation between users [33]. In the long-term, filter bubbles will decrease user activeness and item originality, hurting the ecosystem. Therefore, it is essential to mitigate filter bubbles.

Towards the goal, existing studies prevent the recommendations from merely fitting historical interactions by incorporating additional objectives. For instance, 1) diversity [6, 58], which pushes the recommendation list to cover more item categories; 2) fairness [3, 30], which pursues fair exposure opportunities over item categories; and 3) calibration [39, 48], which ensures that the recommendation list exhibits the same distribution over item categories as the user's history. However, these methods typically make the trade-off across multiple objectives, sacrificing accuracy and even degrading user experience [39, 58]. Moreover, in the feedback loop, users passively adjust the recommendations by user feedback (e.g., click, like, and dislike), which is inefficient and inadequate because users need to constantly provide user feedback until the system recognizes users' intention.

We argue that users have the right to decide whether to mitigate filter bubbles and choose which bubble to mitigate. To this end, we conceptually propose a new prototype called *User-Controllable Recommender System* (UCRS) with three main considerations: 1) the system has the responsibility to remind users if they are stuck in filter bubbles; 2) the system should provide various commands to fully support users' control intentions; and 3) the system should respond to the controls on the fly. UCRS achieves the three objectives with three additional functions beyond conventional recommender systems.



**Figure 1: Illustration of feedback loop and user controls. User feedback passively affects the recommendation strategy while user controls can directly adjust the strategy.**

- **Filter bubble alert.** We define several metrics to measure the strength of filter bubbles. With these metrics, *e.g.*, presented as a system notification, we aim to let users understand the status of filter bubbles and decide whether to mitigate the bubbles.
- **Control commands.** We suggest user controls at two levels regarding either a user or item feature. At the fine-grained level, UCRS supports the commands to increase the items *w.r.t.* a specified user or item feature, such as “more items liked by *young* users” and “more items in a target *category* (*e.g.*, action movies)”. Noticing that users may not intend to specify the target group, UCRS also supports commands at the coarse-grained level, *e.g.*, “no bubble *w.r.t.* my age” and “no bubble *w.r.t.* item category”.
- **Response to user controls.** Once receiving control commands, UCRS adjusts the recommendations by incorporating the commands into recommender inference<sup>1</sup>. This is however non-trivial because some out-of-date user representations learned from historical interactions have encoded the preference information leading to filter bubbles. Thus such user representations can still cause homogeneous recommendations.

To tackle the challenges, we propose a causality-enhanced *User-Controllable Inference* (UCI) framework, which inspects the generation procedure of recommendations from a causal view and leverages counterfactual inference to mitigate the effect of out-of-date user representations. Specifically, UCI imagines a counterfactual world where out-of-date user representations are discarded, and estimates their effects as the difference between factual and counterfactual worlds. After deducting such effects, UCI incorporates the control command into recommender inference. As to user-feature controls, UCI revises the user feature specified by the control command (*e.g.*, changing age from middle age to teenager) to conduct the final inference at the two levels. As to item-feature controls, UCI adopts a user-controllable ranking policy to control the recommendations *w.r.t.* item category. Extensive experiments on three datasets validate the superiority of UCI on mitigating filter bubbles without sacrificing recommendation accuracy. We release the code and data at: <https://github.com/WenjieWWJ/UCRS>.

To sum up, the contributions of this work are threefold:

- We study a new problem of using user controls to adjust filter bubbles, and propose a user-controllable recommender prototype, emphasizing the user rights of controlling recommender systems.
- We propose the UCI framework, which can mitigate the effect of out-of-date user representations via counterfactual inference and perform real-time adaptation to four kinds of user controls.

<sup>1</sup>UCRS cannot respond to user controls timely by model retraining or fine-tuning. The computation cost is also unaffordable.

- We define several metrics to measure filter bubbles and conduct extensive experiments on three datasets, validating the effectiveness of UCI in maintaining the accuracy and mitigating filter bubbles by following user controls.

## 2 RELATED WORK

• **Filter bubbles in recommendation.** Although recommender systems have achieved great success in the past years [11, 12, 22], the debate on filter bubbles has always attracted extensive attention [2, 31]. On one side, researchers claim that recommender systems provide users with satisfying items, and might expose some items that users would never see without recommendations [31]. On the other side, many studies [14, 33, 41] have stated that personalized recommendation would cause group polarization, where users are fragmented and the users with similar interests are grouped. Later, some work introduces filter bubbles in recommendation: users always receive similar content and gradually become isolated from diverse items [1, 18]. Moreover, due to the feedback loop [25, 29], the continual exposure of similar content will further intensify user interests over such items [7], leading to the issues of echo chamber [15, 19, 45] and ideological segregation [33]. In this work, we have found that filter bubbles do exist in the scenarios of content recommendation, such as movies and books. Different from previous work, our consideration is to let users decide whether to mitigate filter bubbles and directly control the recommendations. This is a big step for user engagement because it transfers the decision right from recommender platforms to users.

• **Diversity in recommendation.** Diversity has been widely used as one additional objective to alleviate filter bubbles [10], where recommender models are encouraged to generate dissimilar items in a recommendation list [6, 13]. Generally, item similarity can be compared by various distance functions (*e.g.*, cosine similarity) and item features (*e.g.*, item category and well-trained embeddings) [59]. Technically, the diversity-oriented recommendation can be divided into post-processing [5] and end-to-end methods [58]. The former diversifies the recommendation lists generated by some models via re-ranking [5, 59]. In contrast, the latter directly balances the objectives of accuracy and diversity during training and inference [9]. However, existing methods simply recommend diverse items to users, and then find new item categories liked by the users. This process does not only take lots of time and user feedback, but also brings many irrelevant items [59]. To solve these problems, our proposed UCRS utilizes user controls to indicate user intention and provide efficient diversification.

• **Fairness in recommendation.** Extensive fairness-oriented work has considered encouraging equal exposure across item groups, where groups can be partitioned by item features, such as producer and category. Previous studies [28, 30] usually focus on the definitions of fairness, spanning from amortized equity of attention [3], discounted cumulative fairness [54], to multi-sided fairness [53]. Besides, Steck *et al.* [39] proposed the objective of calibration, which forces the proportion of item categories in a recommendation list to follow that of user’s historical interactions [48]. Although fairness-related work is able to partly alleviate filter bubbles, they inevitably lie in the trade-off between accuracy and fairness, thus degrading the user experience.

• **User-controllable recommendation.** User controls can help users explicitly specify their interests and efficiently achieve recommendation adjustments [24, 43, 47]. Prior literature usually falls into two categories: user controls during the preference estimation [23, 27] and the controls over recommendation lists [38, 42]. In the stage of preference estimation, recommender systems can acquire user controls by multiple ways, such as preference forms [23] and interactive conversations [32, 40]. Once recommendations are presented, users can control by critiquing [8] and interactive explanations [44, 56]. Although these methods can perform user controls in different stages, they ignore the advantages of user controls on alleviating filter bubbles. Besides, previous studies never consider the possible inconsistency between out-of-date user representations and user controls. As such, existing methods can serve as the interface to acquire user controls, which are then used in UCRS to mitigate filter bubbles via the causal UCI framework.

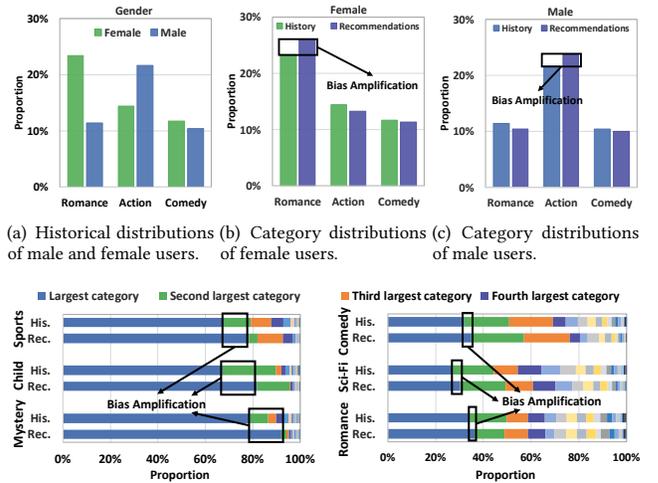
• **Causal recommendation.** We use counterfactual inference to mitigate the effect of out-of-date user representations, which is related to causal recommender models [4, 26, 60]. Generally, two causal frameworks have been applied to recommendation: *potential-outcome framework* [36] and *Structural Causal Models (SCMs)* [34]. The former mainly leverages inverse propensity scoring [37, 51] and doubly robust [21] to debias user feedback. SCMs typically abstract causal relationships into causal graph and estimate causal effects via intervention [48, 57] or counterfactual inference [52, 55], which are widely used for debiasing [48], explainable [44, 46], and out-of-distribution recommendations [50]. Nevertheless, using causality for diversity or alleviating filter bubbles receives little scrutiny.

### 3 PRELIMINARY ON FILTER BUBBLES

To intuitively understand filter bubbles, we conduct preliminary experiments to analyze their effects *w.r.t.* different user groups.

• **Experimental settings.** We train a representative recommender model, Factorization Machine (FM) [35], on three public datasets (*i.e.*, DIGIX-Video, Amazon-Book, and ML-1M), and then collect the top-10 recommended items for each user. Next, to study the phenomenon of filter bubbles, we split users into groups according to two factors: user features and user interactions. Specifically, we are able to divide user groups by available user features, such as gender and age. Besides, different users usually have interests in different item categories (*e.g.*, romance movies), and thus we can also distinguish user groups by user interactions over item categories. For each item category, we select the users whose interaction proportion over this category is larger than a threshold (*e.g.*, 0.5). Thereafter, we compare users' historical interactions and the recommendations generated by FM *w.r.t.* user groups.

• **Analysis.** For male and female users in DIGIX-Video, we visualize their historical distributions over top-3 item categories in Figure 2(a). From the figure, we can observe that male and female users express different interests in item categories. For example, as compared to females, male users prefer more action movies than romance movies. Consequently, the recommender models will inherit the biased distributions [39, 57]. As shown in Figure 2(b) and (c), the distribution of recommendations for male and female users is quite similar to that in the history, showing that the users will continually receive homogeneous items. Worse still, the



(a) Historical distributions of male and female users. (b) Category distributions of female users. (c) Category distributions of male users.

(d) Category proportion in history and recommendations on Amazon-Book. (e) Category proportion in history and recommendations on ML-1M

**Figure 2: Analysis of filter bubbles and bias amplification. In Figure (d) and (e), “His.” and “Rec.” denote the historical interactions and recommendations, respectively.**

models tend to amplify the bias and expose more historical majority categories [48] as shown in Figure 2(b) and (c), causing severer segregation between male and female users.

As to the user groups divided by user interactions, we present the results on Amazon-Book and ML-1M in Figure 2(d) and (e), respectively. The results on DIGIX-Video with similar observations are omitted to save space. From the figures, we have the following findings. 1) The largest categories in users' history are dominating the recommendation lists. Besides, as compared to Amazon-Book, the distributions on ML-1M are more diverse and the domination of majority categories is less severe. This is because most items in ML-1M have multiple categories. 2) The models usually have the bias amplification issue [48] and increase the proportions of recommended majority categories. Due to the bias amplification, filter bubbles will be gradually intensified, which inevitably narrow users' interests, fragment users, and lead to group segregation.

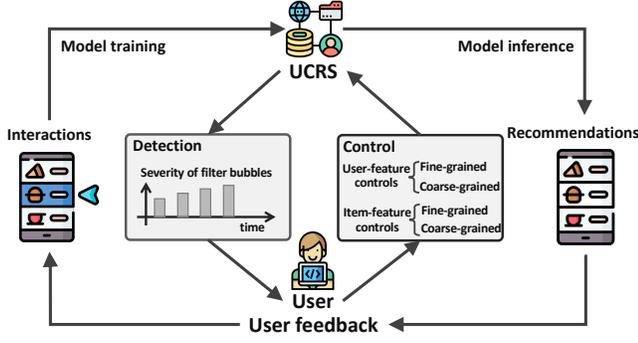
• **Summary.** We find that filter bubbles exist on the sides of user and item features. The bubbles *w.r.t.* item features are caused by the biased interactions over item categories. In this light, we propose the user-feature and item-feature controls correspondingly.

## 4 USER-CONTROLLABLE RECOMMENDATION

In this section, we first formulate the paradigm of user-controllable recommender systems, and then introduce the proposed causal UCI framework for the response to real-time user controls.

### 4.1 Formulation of UCRS

**4.1.1 User-controllable Recommender Systems.** As shown in Figure 3, UCRS introduces another loop between the users and recommender systems by incorporating two modules: *detection and control modules*. First, the detection module is used to measure



**Figure 3: Illustration of the proposed UCERS, which introduces another loop between users and recommender systems for the detection of filter bubbles and user controls.**

the severity of filter bubbles over time, and alert users if they are heavily stuck in filter bubbles. If users are willing to mitigate filter bubbles, they can utilize the control commands and perform real-time adjustments over the recommendations by the control module.

Formally, given users’ historical interactions  $D$ , traditional recommender models aim to predict the recommendations  $R$  via  $P(R|D)$ . In contrast, UCERS additionally considers user controls  $C$  and estimates  $P(R|D, do(C))$  with the user interventions  $do(C)$  [34], where the interventions formulate the four kinds of controls from a causal view. By the interventions, users are able to quickly adjust the recommendations, significantly decrease the items in historical majority categories, and freely jump out of filter bubbles. Specifically, we formulate four kinds of user controls based on user and item features at the fine-grained and coarse-grained levels.

**4.1.2 User-feature Controls.** We represent  $N$  features of user  $u$  as  $\mathbf{x}_u = [x_u^1, \dots, x_u^n, \dots, x_u^N]$ , where  $x_u^n \in \{0, 1\}$  denotes whether user  $u$  has the feature  $x^n$ . For instance, if  $[x^1, x^2]$  represents the features of male and female,  $\mathbf{x}_u = [0, 1]$  indicates that user  $u$  is female.

- **Fine-grained user-feature controls.** To alleviate the filter bubbles *w.r.t.* user features (e.g., gender and age), we design the fine-grained user-feature controls, which prompt UCERS to recommend more items liked by other user groups. For example, 30-year-old users might have interests in the videos liked by teenagers. Formally, to perform  $P(R|D, do(C))$  for user  $u$ , we formulate the control as  $do(C = c_u(+\hat{x}, \alpha))$ , where  $c_u(+\hat{x}, \alpha)$  is the control command to expose more items liked by other user group  $\hat{x}$  and  $c_u(+\hat{x}, \alpha)$  requires that user  $u$  does not have feature  $\hat{x}$ , i.e.,  $\hat{x}_u = 0$  for user  $u$ . Besides,  $\alpha \in [0, 1]$  is a coefficient to adjust the strength of user controls on recommendations.

- **Coarse-grained user-feature controls.** However, users might simply want to mitigate filter bubbles and do not enjoy the items liked by other user groups. In addition, some users possibly do not know which user group is attractive. As such, we propose the coarse-grained user-feature controls, which help to jump out the filter bubbles of users’ own groups. For example, 30-year-old users may not wish the recommendations to be restricted by the feature “age=30”. Formally, the control  $do(C)$  in  $P(R|D, do(C))$  is formulated as  $do(C = c_u(-\bar{x}, \alpha))$ , which reduces the items liked by user’s own group  $\bar{x}$ , i.e.,  $\bar{x}_u = 1$  for user  $u$ .

**4.1.3 Item-feature Controls.** Although user-feature controls are able to mitigate the filter bubbles *w.r.t.* user features, they ignore the filter bubbles caused by user interactions. As shown in Figure 2(d), recommender models typically expose more items in the historical majority categories. Therefore, to complement user-feature controls, we develop item-feature controls to adjust recommendations *w.r.t.* item features. Similar to user features, we represent  $M$  features of item  $i$  as  $\mathbf{h}_i = [h_i^1, \dots, h_i^m, \dots, h_i^M]$ , where  $h_i^m \in \{0, 1\}$  denotes whether item  $i$  has the feature  $h^m$ , e.g., action movies.

- **Fine-grained item-feature controls.** If users have the target item categories (e.g., more romance movies), fine-grained item-feature controls can be applied to increase their recommendations. Specifically, the intervention  $do(C)$  can be expressed as  $do(C = c_i(+\hat{h}, \beta))$ , where  $\hat{h}$  is the target item category and  $\beta \in [0, 1]$  is to modify the strength of user controls.

- **Coarse-grained item-feature controls.** Correspondingly, we suggest coarse-grained item-feature controls to alleviate the users’ burden of specifying target item categories. The goal of coarse-grained item-feature controls is to decrease the recommendations of the largest item category in user historical interactions. In particular, the intervention can be denoted as  $do(C = c_i(-\bar{h}, \beta))$ , where  $-\bar{h}$  means to reduce the largest category  $\bar{h}$  in the history.

## 4.2 Instantiation of UCERS

The key to instantiating UCERS lies in the implementation of the detection and control modules.

**4.2.1 Detection of Filter Bubbles.** We suggest several metrics to measure the severity of filter bubbles from different perspectives, such as diversity and isolation. At different time periods, we can calculate the metrics and obtain the severity level of filter bubbles (e.g., from 1 to 5) by some heuristic rules designed by the recommender platform. Subsequently, the severity level is presented to users and let users decide whether to control filter bubbles.

- **Coverage.** Filter bubbles usually decrease the diversity of recommended items, and thus we incorporate a widely adopted metric for diversity: coverage, which calculates the number of item categories in the recommendation list [58].

- **Isolation Index.** In addition to diversity-based metrics, we propose Isolation Index [20] to measure the segregation across different user groups, which is popular to estimate the ideological segregation in sociology [20]. Here, we revise it for the recommendation task. Formally, given two user groups  $a$  and  $b$ , we can calculate the Isolation Index of their recommendations by

$$s = \sum_{i \in \mathcal{I}} \left( \frac{a_i}{a_n} \cdot \frac{a_i}{a_i + b_i} \right) - \sum_{i \in \mathcal{I}} \left( \frac{b_i}{b_n} \cdot \frac{a_i}{a_i + b_i} \right), \quad (1)$$

where  $\mathcal{I}$  is the item set;  $a_i$  and  $b_i$  are the numbers of users in group  $a$  and  $b$  who receive the recommendation of item  $i$ . Besides,  $a_n = \sum_{i \in \mathcal{I}} a_i$ , which is the total frequency of items exposed to the users in group  $a$ . Meanwhile,  $b_n$  is similar to  $a_n$ . Finally,  $s \in [0, 1]$  is equal to the weighted average item exposure of group  $a$  minus that of group  $b$ , where the weights are  $\frac{a_i}{a_i + b_i}$  [20]. Intuitively,  $s$  captures the extent of recommendation segregation between two groups and higher values denote severer segregation. If there are multiple user groups, we take the average value of  $s$  between any pair of groups.

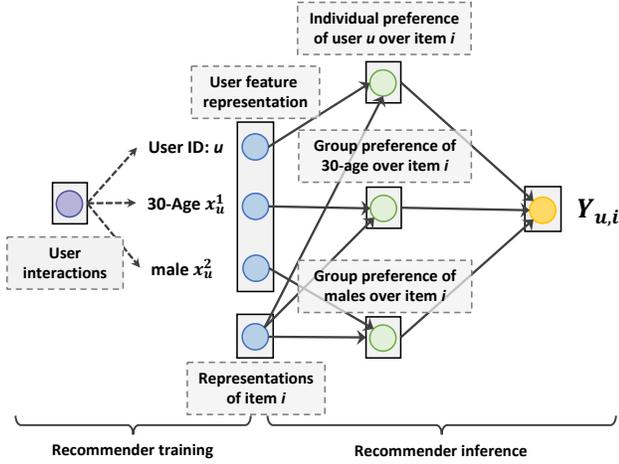


Figure 4: Causal graph behind the generation procedure of recommendations.

- **Majority Category Domination (MCD).** Isolation index is more suitable to measure the group segregation *w.r.t.* user features (e.g., age and gender). As to filter bubbles *w.r.t.* item features, we can utilize MCD to obtain the proportion of the historically largest item category in recommendation lists. The increase of MCD across different time periods reflects that the filter bubbles *w.r.t.* item category are becoming increasingly serious.

**4.2.2 Response to User-feature Controls.** If users aim to mitigate filter bubbles, UCRS needs real-time responses to user controls. For fine-grained user-feature controls  $do(C = c_u(+\hat{x}, \alpha))$ , users want more items liked by other user group  $\hat{x}$ . As such, UCRS is actually required to generate recommendations based on the changed user features, which are contrary to the facts. For example, age changes from 30 to 18. From a causal view, the objective of fine-grained user-feature controls is to answer a counterfactual question [34]: *what the recommendations would be if the user were in a counterfactual group  $\hat{x}$ ?* Similarly, coarse-grained user-feature controls are to answer the question: *what the recommendations would be if the user were not in the factual group  $\bar{x}$ ?* To answer the counterfactual questions, the UCI framework needs to inspect the causal relations between user features and recommendations, and then conduct counterfactual inference [34].

- **Causal view of generating recommendations.** As illustrated in Figure 4, we analyze the generation procedure of recommendations by a causal graph [34]. Specifically, for most models (e.g., FM [35]), recommender training learns the user representations from interactions, including the representations of ID, age, and gender. Thereafter, the representations of user  $u$  and item  $i$  are used to predict the probability of user  $u$  preferring item  $i$ , i.e.,  $Y_{u,i} \in [0, 1]$ . In detail,  $Y_{u,i}$  is fused from the preference scores of the individual ID and multiple group features, where the group preference is shared by the users in the corresponding user group.

To answer the counterfactual question of fine-grained user-feature controls, an intuitive solution is to change the user features for recommender inference, e.g., changing the age from 30 to 18. As to coarse-grained user-feature controls, we can directly discard the user feature  $\bar{x}$  (e.g., the age of 30) for inference. However,

as shown in Figure 4, user interactions are actually confounders, which affect the representations of user ID and other group features (e.g., age) during recommender training. Therefore, the correlations exist between the representations of user ID and group features. Although the group features are changed or discarded, user ID representations still encode the out-of-date interests of original features, which are inconsistent with user controls and hinder the recommendations of target user groups.

To remove the confounding effect, the popular choices are confounder balancing [36], back-door and front-door adjustments [34]. Nevertheless, confounder balancing and back-door adjustment require estimating the causal effect of confounders on representations. The estimation is infeasible because 1) user interactions are in a dynamic high-dimension space where new interactions are continually increasing; and 2) the effect of user interactions on representations is decided by the recommender training process, which differs across models and training manners (e.g., optimizer and learning rate). Besides, front-door adjustment needs to discover the mediator that blocks all back-door paths, which is not applicable in the causal graph of Figure 4. To avoid these challenges, we propose to directly reduce the causal effect of user ID representations on the prediction  $Y_{u,i}$  during inference, which can effectively decrease the influence of out-of-date representations without knowing the training process.

- **Implementation of counterfactual inference.** In particular, the UCI framework first estimates the effect of user ID representations via counterfactual inference, and then deducts it from the original prediction  $Y_{u,i}$  [16, 17, 49]. Formally, we image *what the prediction  $Y_{\hat{u},i}$  would be if user  $u$  had not the ID representations in a counterfactual world [49], where  $\hat{u}$  denotes the representations of user  $u$  without ID representations. By comparing  $Y_{u,i}$  with  $Y_{\hat{u},i}$ , we can measure the effect of user ID representations by  $Y_{u,i} - Y_{\hat{u},i}$ . Thereafter, we subtract it from the original prediction  $Y_{u,i}$  with the coefficient  $\alpha$ :*

$$\begin{aligned} Y_{u,i} - \alpha \cdot (Y_{u,i} - Y_{\hat{u},i}) \\ = f(u, i) - \alpha \cdot (f(u, i) - f(\hat{u}, i)) \\ = (1 - \alpha) \cdot f(u, i) + \alpha \cdot f(\hat{u}, i), \end{aligned} \quad (2)$$

where  $f(\cdot)$  can be any recommender function of using user and item representations to calculate the prediction  $Y$  (e.g., FM), and  $\alpha \in [0, 1]$  adjusts the strength of mitigating the effect of user ID representations.

- **Summary of UCI.** The UCI framework consists of two steps to answer the two questions of user-feature controls during inference: 1) changing specific user features to  $\hat{x}$  for fine-grained controls and discarding the user feature  $\bar{x}$  for coarse-grained controls; and 2) using counterfactual inference to mitigate the effect of out-of-date user ID representations via Equation (2).

**4.2.3 Response to Item-feature Controls.** As to item-feature controls, fine-grained controls aim to increase the target item category  $\hat{h}$  while coarse-grained ones are to decrease the largest item category  $\bar{h}$  of users' history. Indeed, they are asking two interventional questions [34]: *what the recommendations will be if users want more items in the target category  $\hat{h}$  or users do not want the largest category  $\bar{h}$ ?* To answer such questions, the UCI

framework utilizes a user-controllable ranking policy as:

$$Y'_{u,i} = Y_{u,i} + \beta \cdot r(i), \quad (3)$$

where  $Y'_{u,i}$  is the revised score for ranking, and  $\beta \in [0, 1]$  is a coefficient to adjust the strength of user controls. Besides,  $r(i)$  denotes a regularization term over item  $i$ . Specifically,

$$r(i) = \begin{cases} 2, & \text{if } \hat{h}_i = 1 \text{ for item } i \text{ with fine-grained controls} \\ 0, & \text{if } \hat{h}_i = 1 \text{ for item } i \text{ with coarse-grained controls} \\ 1, & \text{otherwise,} \end{cases} \quad (4)$$

where  $r(i)$  encourages more recommendations of the items in the target category  $\hat{h}$  under fine-grained controls, and decreases the largest category  $\hat{h}$  if coarse-grained controls are applied.

• **Target category prediction.** Due to the extensive item categories, it is a burden for users to specify target categories in fine-grained item-feature controls. Although coarse-grained controls partly alleviate the burden, we can further enhance it by predicting the possible target categories for users. Specifically, if users wish to decrease the largest category of the history, we can predict which item category users will prefer, and then improve the coarse-grained item-feature controls with fine-grained ones.

As shown in Figure 5, we sort users' interacted items by time, and then split the interaction sequence into two parts to obtain the distributions over item categories, respectively. Next, we predict the second category distribution based on the first one via multiple Multi-Layer Perceptrons (MLPs). During training, MLPs utilize the category distributions of all users to capture 1) the temporal interest transition (e.g., the increasing preference over some categories), and 2) the relationships between item categories (e.g., the users liking action movies probably prefer crime movies). In the inference stage, we leverage the second category distribution to predict top- $K$  target categories. Besides, we conduct intervention  $do(\hat{h} = 0)$  to indicate the user controls of reducing category  $\hat{h}$ . The top- $K$  item categories with the highest values are treated as the target ones in the fine-grained controls. Finally, UCI further enhances coarse-grained item-feature controls by using target category prediction.

• **Summary of UCI.** Under item-feature controls, user ID representations also encode the historical interests, which conflict with the objective of increasing target categories or decreasing the historical majority category. Therefore, 1) UCI first conducts counterfactual inference to mitigate the causal effect of user ID representations on  $Y_{u,i}$  as illustrated in Equation (2); 2) for coarse-grained item-feature controls, UCI leverages target category prediction to obtain the top- $K$  target categories; and 3) UCI adopts the ranking policy in Equation (3) for recommendations.

## 5 EXPERIMENTS

We conduct experiments to answer the following questions:

- **RQ1.** How does UCI perform to adjust recommendations for alleviating filter bubbles via four kinds of user controls?
- **RQ2.** How do users can control the recommendations by the coefficients (i.e.,  $\alpha$  and  $\beta$ )?
- **RQ3.** How does the proposed counterfactual inference affect the recommendations?

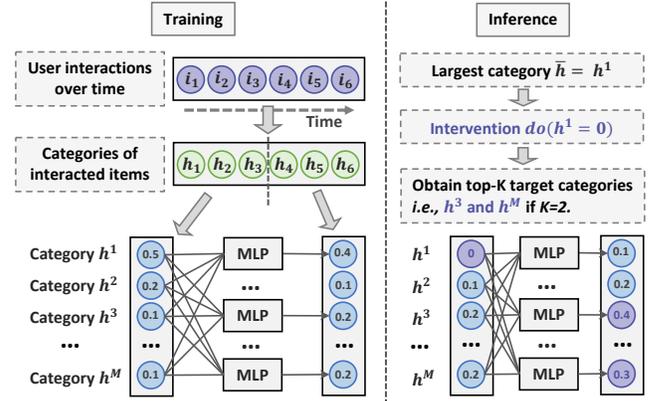


Figure 5: Illustration of the training and inference procedures to predict the target item categories.

Table 1: Statistics of the three datasets. “#IC” denotes the number of item categories.

Datasets	#Users	#Items	#Interactions	Density	#IC
DIGIX-Video	7,643	15,526	316,045	0.0027	135
ML-1M	6,040	3,883	575,276	0.0245	18
Amazon-Book	29,115	16,845	1,712,409	0.0035	29

### 5.1 Experimental Settings

• **Datasets.** We utilize three real-world datasets for experiments: DIGIX-Video, ML-1M, and Amazon-Book, which are publicly available and vary in terms of domain, user/item features, and sparsity. The statistics of the datasets are presented in Table 1. Specifically, 1) DIGIX-Video<sup>2</sup> is a video recommendation dataset, released by 2021 DIGIX AI Challenge. It covers rich user and item features, including age, gender, and item category. 2) ML-1M<sup>3</sup> is a widely-used movie dataset, in which each movie usually has multiple categories. 3) Amazon-Book<sup>4</sup> is a popular dataset for book recommendations, where users only have ID features and each item has a hierarchical taxonomy. In contrast to multi-label item categories in ML-1M, we only keep the largest category, and thus each book is assigned with only one category. For each dataset, we use the 10-core setting and treat the interactions with ratings  $\geq 4$  as positive samples. In addition, we sort interactions by timestamps, and then split 80%, 10%, and 10% of interactions as the training, validation, and test sets, respectively. For each interaction, we randomly sample an unobserved interaction as the negative sample for training.

• **Evaluation of user-feature controls.** Since online testing is expensive and infeasible for researchers, we design an offline evaluation setting: 1) assuming some users are willing to mitigate filter bubbles and provide the four kinds of controls; 2) generating the recommendations by different recommender methods according to user controls; and 3) evaluating the recommendations in terms of accuracy and the metrics on mitigating filter bubbles, such as Isolation Index, MCD, and Coverage.

<sup>2</sup><https://www.kaggle.com/voler2333/2021-digix-video-recommendation>.

<sup>3</sup><https://grouplens.org/datasets/movielens/1m/>.

<sup>4</sup><https://nijianmo.github.io/amazon/index.html>.

**Datasets.** For user-feature controls, we utilize DIGIX-Video for evaluation because it has rich user features (*i.e.*, gender and age) and video categories. In contrast, Amazon-Book only has user ID features; and the users in ML-1M are heavily affected by the dataset bias, where the favorite movie categories of over 77% of users are “drama” and “comedy”, and 10% popular movies occupy 52% interactions. Consequently, the users with different features (*e.g.*, age) show similar interaction distributions. As such, Amazon-Book and ML-1M are not well suitable for evaluating the filter bubbles *w.r.t.* user features. In this work, we test the fine-grained and coarse-grained user-feature controls over the gender groups and the age groups of DIGIX-Video, respectively. The users take the opposite gender group as the target under fine-grained controls and want to jump out of their own age groups under coarse-grained controls. This is because the number of age groups is larger, and thus users are more likely to utilize coarse-grained controls without the burden of specifying target age groups.

**Baselines.** All the baselines and our proposed UCI are model-agnostic, which are compared on two representative recommender models: FM and Neural Factorization Machine (NFM) [22].

- 1) **woUF** trains the models without user features (woUF) such as age and gender, which possibly alleviate the segregation across user groups during recommender training.
- 2) **changeUF** utilizes well-trained recommender models and only changes user features (UF) to the target  $\hat{x}$  for inference, *e.g.*, changing age from 30 to 18. ChangeUF is used for the fine-grained user-feature controls.
- 3) **maskUF** discards the original user features  $\bar{x}$  (*e.g.*, age=30) for the inference of coarse-grained user-feature controls.
- 4) **Fairco** [30] is a user-controllable ranking algorithm, which pursues fair exposure opportunities across item groups.
- 5) **Diversity** [59] incorporates a re-ranking method to diversify recommendations by minimizing the intra-list similarity.

**Metrics.** To measure the performance, we utilize the all-ranking protocol [49], and the top-10 items are returned as recommendations. We adopt **Recall** and **NDCG** to evaluate the accuracy. To quantify the severity of filter bubbles, we leverage **Isolation Index** (Iso-Index) and **Coverage** to estimate the group segregation and diversity. In addition, for the fine-grained user-feature controls with target user groups, we develop the metrics **DIS-EUC** to compare the distance between the recommendations of users and groups. Formally, we denote  $\bar{x}$  and  $\hat{x}$  as the original and target groups of user  $u$ , respectively;  $d_u \in \mathbb{R}^M$  is the distribution over item categories in the recommendations of user  $u$ ;  $\bar{g}_u \in \mathbb{R}^M$  denotes the same distribution by averaging over the users in the original group  $\bar{x}$  (*e.g.*, 30-year-old users); and  $\hat{g}_u \in \mathbb{R}^M$  represents the same distribution of the target group  $\hat{x}$ . Thereafter, we calculate  $\text{DIS-EUC} = \text{dis}(d_u, \hat{g}_u) - \text{dis}(d_u, \bar{g}_u)$  for user  $u$ , where  $\text{dis}(\cdot)$  uses Euclidean distance. DIS-EUC measures the distance difference from the user to two groups, where larger distances indicate severer group segregation and filter bubbles.

• **Evaluation of item-feature controls.** We conduct experiments on the users who have preference shifts from the training to test sets. Specifically, for each user, we obtain the largest item categories in the training and test sets, and then we select the users with different largest categories. This simulates the situation that users

**Table 2: Performance comparison between UCI and the baselines under the coarse-grained user-feature controls.**

	Recall $\uparrow$	NDCG $\uparrow$	Iso-Index $\downarrow$	Coverage $\uparrow$
Random	0.0008	0.0005	0.0008	11.6185
FM	0.0758	0.0584	0.1082	9.5191
FM-woUF	0.0757	0.0582	0.1195	9.9211
FM-maskUF	0.0756	0.0577	0.1048	9.6185
FM-Fairco	0.0728	0.0534	0.1050	<b>9.9241</b>
FM-Diversity	0.0756	0.0574	0.1025	9.8742
FM-UCI	<b>0.0767</b>	<b>0.0592</b>	<b>0.0777</b>	9.8802
NFM	<b>0.0774</b>	0.0585	0.1144	10.2670
NFM-woUF	0.0722	0.0542	0.1191	10.1445
NFM-maskUF	0.0759	0.0575	0.1378	9.9740
NFM-Fairco	0.0755	0.0571	0.1130	<b>10.3458</b>
NFM-Diversity	0.0741	0.0562	0.1026	10.3268
NFM-UCI	0.0767	<b>0.0596</b>	<b>0.0760</b>	9.9000

aim to mitigate historical filter bubbles and want more items in other categories. The numbers of selected users in DIGIX-Video, ML-1M, and Amazon-Book are 4320, 3806, and 5155, respectively.

**Baselines.** We generate recommendations for the selected users by using the following methods: 1) **woIF** trains FM and NFM without item features (woIF); 2) **Fairco**; 3) **Diversity**; 4) **Reranking** is one variant of UCI, which only uses the ranking policy in Equation (3) and discards counterfactual inference and target category prediction; 5) **C-UCI** denotes the UCI strategy with target category prediction under coarse-grained controls; and 6) **F-UCI** represents UCI under fine-grained controls, which knows the target category of each user, *i.e.*, the largest category in the test set.

**Metrics.** For performance comparison, we use Recall, NDCG, and Coverage. Besides, we introduce a new metric **Weighted NDCG** (W-NDCG), which assigns the NDCG relevance scores of the positive items in the target categories, the positive ones in non-target categories, and the negative ones as 2, 1, and 0, respectively. W-NDCG distinguishes the positive items in the target and non-target categories, and prefers the positive ones in the target categories. Furthermore, we employ MCD and **Target Category Domination** (TCD) to calculate the proportions of the historical majority category and users’ target category in the recommendations, respectively.

• **Hyper-parameter settings.** We train FM and NFM by following the settings in [22]: the sizes of user/item representations are 64; and Adagrad with the batch size of 1,024 is used for parameter optimization. The learning rate is searched in {0.001, 0.01, 0.05}. The hidden size of the MLP in NFM and target category prediction is tuned in {4, 8, 16, 32} and the normalization coefficient is searched from {0, 0.1, 0.2}.  $K$  in target category prediction is chosen from {1, 2, ..., 5}. Besides,  $\alpha$  and  $\beta$  in the controls of  $c_u(\cdot)$  and  $c_i(\cdot)$  are adjusted in {0, 0.1, ..., 0.5} and {0, 0.01, ..., 0.1}, respectively. We select the best model by Recall on the validation set.

## 5.2 Performance Comparison

**5.2.1 Performance under User-feature Controls (RQ1).** We present the results under the coarse-grained and fine-grained user-feature controls in Table 2 and Table 3, respectively. From the two tables, we have the following observations:

**Table 3: Performance comparison between UCI and the baseline under the fine-grained user-feature controls. The best results are highlighted in bold and the second best ones are underlined.**

	FM					NFM				
	Recall $\uparrow$	NDCG $\uparrow$	Iso-Index $\downarrow$	DIS-EUC $\downarrow$	Coverage $\uparrow$	Recall $\uparrow$	NDCG $\uparrow$	Iso-Index $\downarrow$	DIS-EUC $\downarrow$	Coverage $\uparrow$
Random	0.0008	0.0005	0.0008	0.0001	11.6185	0.0008	0.0005	0.0008	0.0001	11.6185
FM/NFM	<u>0.0861</u>	<u>0.0650</u>	0.1161	0.0568	8.6781	<b>0.0849</b>	<u>0.0630</u>	0.1046	0.0436	9.5126
woUF	0.0858	0.0648	0.1156	0.0561	8.7526	<u>0.0847</u>	0.0630	0.1048	0.0431	9.5193
changeUF	0.0858	0.0649	0.1152	0.0566	8.6859	0.0839	0.0626	0.1035	0.0432	9.5461
Fairco	0.0782	0.0550	0.1082	<u>0.0533</u>	<u>9.1206</u>	0.0619	0.0420	0.1011	<u>0.0357</u>	<u>9.7353</u>
Diversity	0.0750	0.0573	0.0995	0.0552	<b>9.4312</b>	0.0731	0.0552	0.0864	0.0399	<b>9.8614</b>
UCI	<b>0.0870</b>	<b>0.0661</b>	<b>0.0979</b>	<b>0.0516</b>	9.0304	0.0844	<b>0.0635</b>	<b>0.0844</b>	<b>0.0354</b>	9.6439

**Table 4: Results of item-feature controls on ML-1M and Amazon-Book. “UB” implies that F-UCI is the “upper bound” of C-UCI.**

Method	ML-1M						Amazon-Book					
	Recall $\uparrow$	NDCG $\uparrow$	W-NDCG $\uparrow$	MCD $\downarrow$	TCD $\uparrow$	Coverage $\uparrow$	Recall $\uparrow$	NDCG $\uparrow$	W-NDCG $\uparrow$	MCD $\downarrow$	TCD $\uparrow$	Coverage $\uparrow$
Random	0.0029	0.0031	0.0029	0.2859	0.0446	8.5024	0.0004	0.0004	0.0004	0.2414	0.0163	4.5892
FM	0.0659	0.0536	0.0485	0.5994	0.2222	8.6600	0.0118	0.0095	0.0085	0.5353	0.2305	3.0175
FM-woIF	0.0649	0.0529	0.0481	0.5952	0.2238	8.6755	0.0116	0.0097	0.0084	0.5310	0.2268	3.1792
FM-Fairco	0.0605	0.0506	0.0459	0.3812	0.2417	<u>9.3799</u>	0.0117	0.0095	0.0085	0.1559	0.2942	<b>3.2700</b>
FM-Diversity	0.0531	0.0473	0.0428	0.5597	0.2351	<b>9.5407</b>	0.0092	0.0080	0.0072	0.5199	0.2383	3.1899
FM-Reranking	0.0761	0.0637	0.0603	<b>0.1000</b>	0.3099	8.9409	<u>0.0178</u>	<u>0.0142</u>	0.0142	<u>0.0026</u>	0.5348	3.0196
FM-C-UCI	<u>0.0770</u>	<u>0.0665</u>	<u>0.0630</u>	<u>0.2466</u>	<u>0.3334</u>	9.1206	0.0173	0.0141	<u>0.0146</u>	0.0213	<u>0.6310</u>	2.0768
FM-F-UCI (UB)	<b>0.2095</b>	<b>0.1704</b>	<b>0.1792</b>	0.3544	<b>1.0000</b>	8.0712	<b>0.0334</b>	<b>0.0283</b>	<b>0.0337</b>	<b>0.0023</b>	<b>1.0000</b>	1.0002
NFM	0.0651	0.0556	0.0501	0.5748	0.2321	8.8854	0.0121	0.0102	0.0088	0.5488	0.2294	2.9818
NFM-woIF	0.0654	0.0551	0.0498	0.5732	0.2290	9.0110	0.0112	0.0092	0.0082	0.5330	0.2332	3.0900
NFM-Fairco	0.0626	0.0516	0.0470	0.4750	0.2441	<u>9.3715</u>	0.0114	0.0091	0.0085	0.1856	0.2891	<b>3.8497</b>
NFM-Diversity	0.0522	0.0481	0.0438	0.5391	0.2391	<b>9.7018</b>	0.0109	0.0092	0.0081	0.5146	0.2427	<u>3.1825</u>
NFM-Reranking	0.0752	0.0672	0.0631	<b>0.0296</b>	<u>0.3163</u>	9.0468	0.0180	0.0144	0.0144	<u>0.0025</u>	0.5421	3.0184
NFM-C-UCI	<u>0.0778</u>	<u>0.0687</u>	<u>0.0647</u>	<u>0.2753</u>	0.3119	9.0744	<u>0.0181</u>	<u>0.0148</u>	<u>0.0154</u>	0.0049	<u>0.6779</u>	1.4190
NFM-F-UCI (UB)	<b>0.2125</b>	<b>0.1729</b>	<b>0.1820</b>	0.3319	<b>1.0000</b>	8.2299	<b>0.0338</b>	<b>0.0276</b>	<b>0.0327</b>	<b>0.0023</b>	<b>1.0000</b>	1.0002

- The intuitive baselines (*i.e.*, woUF, changeUF, and maskUF) slightly decrease the recommendation accuracy and alleviate the group segregation in terms of Isolation Index and DIS-EUC. Meanwhile, diversity rises marginally in most cases. However, the overall performance is quite similar to FM or NFM. This is consistent with the analysis in Section 4.2.2: although the user features are discarded for training or changed/masked for inference, the user ID representations still encode the historical interactions, which are affected by the users’ original features and lead to similar recommendations with FM and NFM.
- Fairness and diversity methods can effectively mitigate the issue of filter bubbles and diversify the recommendation lists. Nevertheless, they bring sharp performance drop. For example, the accuracy of Fairco on FM declines by 15.38% *w.r.t.* NDCG in Table 3. It makes sense because pursuing the objectives of fairness and diversity will inevitably recommend many irrelevant items, occupying the opportunities of positive items.
- UCI significantly alleviates the group segregation in filter bubbles while achieving superior accuracy. Besides, the diversity also increases as compared to FM and NFM, which relieves the accuracy-diversity dilemma. We attribute the improvements to the effectiveness of counterfactual inference in reducing the effect of out-of-date user ID representations. The mitigation of such effect pushes the recommender model to expose fewer items similar to historical interactions, and makes the recommendations more

diverse. Meanwhile, the superior accuracy is because  $\alpha$  in the user control  $c_u(\cdot, \alpha)$  adjusts the influence of the representations of user ID and other group features (*e.g.*, age and gender), leading to a better balance between individual preference and group preference as illustrated in Figure 4.

**5.2.2 Performance under Item-feature Controls (RQ1).** Table 4 provides the results of item-feature controls on ML-1M and Amazon-Book. The observations on DIGIX-Video are similar to those on ML-1M. The consideration of presenting ML-1M and Amazon-Book is to compare the effects of single-label and multi-label item categories. From Table 4, we have the following findings:

- As compared with FM and NFM, woIF marginally decreases the historical majority categories *w.r.t.* MCD, where the marginal effect shows that woIF still recommends many items in the historical majority categories without knowing item features. Moreover, woIF slightly degrades the accuracy and improves the diversity. Such observations and the underlying reasons are analogous to woUF under user-feature controls.
- The performance of Fairco and Diversity is similar to that under user-feature controls: at the expense of sacrificing accuracy, they substantially alleviate the filter bubbles by recommending fewer historical majority categories and improving diversity.
- Reranking and C-UCI have significant performance improvements over the baselines *w.r.t.* the accuracy and mitigation

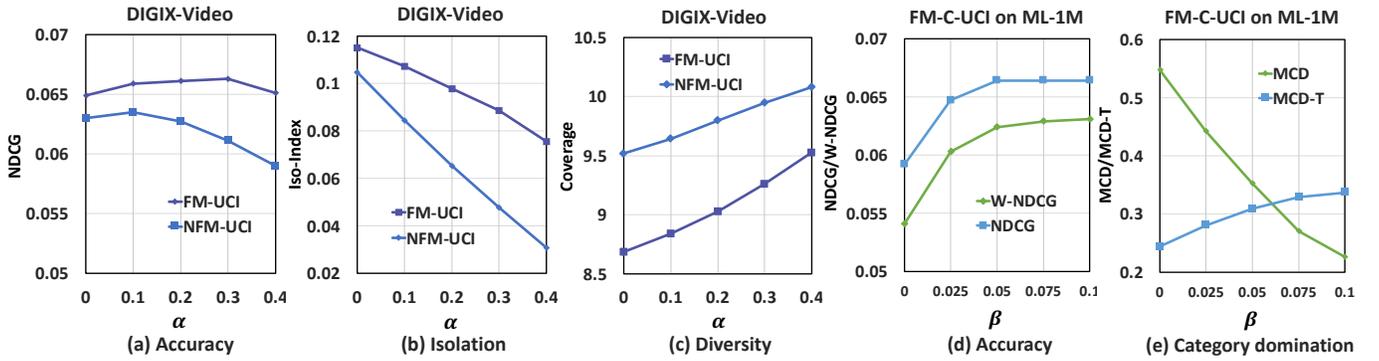


Figure 6: Effects of the control coefficients  $\alpha$  and  $\beta$  w.r.t. accuracy, isolation, diversity, and category domination.

of majority categories. This is due to using coarse-grained item-feature controls, which indicate the largest categories in the history to decrease. Besides, C-UCI performs better than Reranking, especially in terms of W-NDCG, validating the superiority of counterfactual inference and target category prediction. Counterfactual inference reduces the influence of out-of-date ID representations for these users with preference shifts; and meanwhile UCI recommends more target categories due to the target category prediction. Furthermore, as compared to ML-1M, the improvements of UCI over Reranking on Amazon-Book are only significant w.r.t. W-NDCG. It is reasonable because Amazon-Book only has user ID features and counterfactual inference is impractical: only ID representations are to represent users and reducing their effect usually does not change the ranking lists. Therefore, on Amazon-Book, purely target category prediction is helpful to enhance W-NDCG and TCD.

- F-UCI achieves the best accuracy by following users' fine-grained item-feature controls, showing that directly incorporating user controls into recommender inference is greatly effective to understand users' interests as compared to passively learning from user interactions. As the upper bound of C-UCI, it also implies the promising potential of target category prediction, where a more accurate prediction of target categories can lead to dramatic accuracy improvements.
- The performance of diversity varies from ML-1M to Amazon-Book, especially over Reranking, C-UCI, and F-UCI. Both C-UCI and F-UCI decrease the diversity while the relative drop on Amazon-Book is larger. The reasons are that 1) C-UCI and F-UCI emphasize the recommendations of target item categories; and 2) the items in Amazon-Book have single-label categories. Purely recommending the target categories will easily degrade the diversity. This indicates that users should adjust the control coefficients  $\alpha$  and  $\beta$  to balance the superior accuracy and diversity at their own will.

**5.2.3 Effect of Coefficients in User Controls (RQ2).** To study whether users can flexibly adjust the recommendations, we explore the effect of the control coefficients  $\alpha$  and  $\beta$ . The results of FM-UCI and NFM-UCI w.r.t. varying  $\alpha$  on DIGIX-Video are reported in Figure 6(a), (b), and (c). The performance of FM-C-UCI w.r.t.  $\beta$  on ML-1M is summarized in Figure 6(d) and (e). More results on other datasets have similar trends, which are omitted to save space.

Table 5: Performance comparison with (w/) and without (w/o) counterfactual inference (CI).

Method	Variants	Recall	NDCG	W-NDCG	MCD	Coverage
FM-F-UCI	w/o CI	0.2094	0.1689	0.1777	0.3587	7.9519
	w/ CI	<b>0.2095</b>	<b>0.1704</b>	<b>0.1792</b>	<b>0.3544</b>	<b>8.0712</b>
NFM-F-UCI	w/o CI	0.2094	0.1687	0.1775	0.3525	8.0155
	w/ CI	<b>0.2125</b>	<b>0.1729</b>	<b>0.1820</b>	<b>0.3319</b>	<b>8.2299</b>

From the figures, we can find: 1)  $\alpha$  controls the influence of user-feature controls, where a larger  $\alpha$  significantly alleviates the filter bubbles and enhances the diversity as shown in Figure 6(b) and (c); 2) the accuracy w.r.t. increasing  $\alpha$  first rises, and then gradually decreases. Besides, we can observe that the accuracy drop is more steady as compared to the isolation and diversity. Such findings verify that UCI can mitigate filter bubbles and improve diversity without sacrificing accuracy or with less accuracy decline; and 3) from the results in Figure 6(d) and (e) under item-feature controls,  $\beta$  is able to mitigate the domination of historical majority categories while improving the accuracy.

#### 5.2.4 Ablation Study of Counterfactual Inference (RQ3).

We conduct ablation study to further analyze the effect of counterfactual inference. In Figure 6(a), (b), and (c),  $\alpha = 0$  denotes the ablation of counterfactual inference under user-feature controls. Besides, we remove it from FM-F-UCI and NFM-F-UCI under item-feature controls, and summarize the results on ML-1M in Table 5. From Figure 6 and Table 5, we observe that using counterfactual inference alleviates the filter bubbles, improves the diversity, and even enhances the accuracy when  $\alpha$  is small. The higher accuracy is mainly due to the inconsistency between out-of-date user ID representations and the latest user interests.

## 6 CONCLUSION AND FUTURE WORK

In this work, we proposed a novel recommender prototype UCRS to flexibly alleviate filter bubbles, which provides users more choices to adjust recommendations. Functionally, the prototype can detect the severity of filter bubbles and allow users to adjust filter bubbles via user controls. In particular, we developed four kinds of user controls: the user-feature and item-feature controls at the fine-grained and coarse-grained levels. To implement the user controls, we designed a UCI framework for recommender inference, which leverages counterfactual inference to mitigate the effect of out-of-date user

ID representations on recommendations. Furthermore, UCI revises the user features for user-feature controls and adopts a ranking policy with target category prediction for item-feature controls. We proposed several metrics to measure filter bubbles and conducted experiments on three datasets, validating the effectiveness of UCI in alleviating filter bubbles and maintaining the accuracy.

The new UCRS prototype and the novel UCI framework can be widely deployed in the practical recommender systems. The recommender platform can design various interactive interfaces (e.g., conversational systems and control panels) to acquire the control commands, and then adopt UCRS to effectively adjust recommendations. This additional interaction paradigm between users and recommender systems will 1) ensure the user rights of controlling recommender strategies, 2) increase the user engagement in the recommendation ecosystem, and 3) significantly enhance the user satisfaction over the recommended items.

Nevertheless, this work takes the initial step to perform user-controllable recommendation against filter bubbles, leaving many potential directions to future work. In particular, 1) it is non-trivial to instantiate the proposed UCRS framework in the online testing platforms, which is costly and impractical for researchers but shows a better justification *w.r.t.* the efficiency and effectiveness of UCI; and 2) more user controls under the framework of UCRS can be designed by collecting users' opinions, which will help users to quickly adjust recommendations by more diverse interfaces.

## REFERENCES

- [1] Mahsa Badami, Olga Nasraoui, and Patrick Shafto. 2018. PrCP: Pre-recommendation Counter-Polarization. In *KDIR*. 280–287.
- [2] Eytan Bakshy, Solomon Messing, and Lada A Adamic. 2015. Exposure to ideologically diverse news and opinion on Facebook. *Science* 348, 6239 (2015), 1130–1132.
- [3] Asia J Biega, Krishna P Gummedi, and Gerhard Weikum. 2018. Equity of attention: Amortizing individual fairness in rankings. In *SIGIR*. ACM, 405–414.
- [4] Stephen Bonner and Flavian Vasile. 2018. Causal Embeddings for Recommendation. In *RecSys*. ACM, 104–112.
- [5] Jaime Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *SIGIR*. ACM, 335–336.
- [6] Praveen Chandar and Ben Carterette. 2013. Preference based evaluation measures for novelty and diversity. In *SIGIR*. ACM, 413–422.
- [7] Allison JB Chaney, Brandon M Stewart, and Barbara E Engelhardt. 2018. How Algorithmic Confounding in Recommendation Systems Increases Homogeneity and Decreases Utility. In *RecSys*. ACM, 224–232.
- [8] Li Chen and Pearl Pu. 2012. Critiquing-based recommenders: survey and emerging trends. *UMAUT* 22, 1 (2012), 125–150.
- [9] Laming Chen, Guoxin Zhang, and Hanning Zhou. 2018. Fast greedy map inference for determinantal point process to improve recommendation diversity. In *NeurIPS*. 5627–5638.
- [10] Peizhe Cheng, Shuaiqiang Wang, Jun Ma, Jiankai Sun, and Hui Xiong. 2017. Learning to Recommend Accurate and Diverse Items. In *WWW*. IW3C2, 183–192.
- [11] Zhiyong Cheng, Ying Ding, Lei Zhu, and Mohan Kankanhalli. 2018. Aspect-aware latent factor model: Rating prediction with ratings and reviews. In *WWW*. ACM, 639–648.
- [12] Zhiyong Cheng, Fan Liu, Shenghan Mei, Yangyang Guo, Lei Zhu, and Liqiang Nie. 2022. Feature-Level Attentive ICF for Recommendation. *TOIS* 40, 4 (2022), 1–24.
- [13] Charles LA Clarke, Maheedhar Kolla, Gordon V Cormack, Olga Vechtomova, Azin Ashkan, Stefan Bütcher, and Ian MacKinnon. 2008. Novelty and diversity in information retrieval evaluation. In *SIGIR*. ACM, 659–666.
- [14] Pranav Dandekar, Ashish Goel, and David T Lee. 2013. Biased assimilation, homophily, and the dynamics of polarization. *PNAS* 110, 15 (2013), 5791–5796.
- [15] Tim Donkers and Jürgen Ziegler. 2021. The Dual Echo Chamber: Modeling Social Media Polarization for Interventional Recommending. In *RecSys*. ACM, 12–22.
- [16] Fuli Feng, Weiran Huang, Xin Xin, Xiangnan He, and Tat-Seng Chua. 2021. Should Graph Convolution Trust Neighbors? A Simple Causal Inference Method. In *SIGIR*. ACM, 1208–1218.
- [17] Fuli Feng, Jizhi Zhang, Xiangnan He, Hanwang Zhang, and Tat-Seng Chua. 2021. Empowering Language Understanding with Counterfactual Reasoning. In *ACL-IJCNLP Findings*. ACL, 2226–2236.
- [18] Seth Flaxman, Sharad Goel, and Justin M Rao. 2016. Filter bubbles, echo chambers, and online news consumption. *Public opinion quarterly* 80, S1 (2016), 298–320.
- [19] Yingqiang Ge, Shuya Zhao, Honglu Zhou, Changhua Pei, Fei Sun, Wenwu Ou, and Yongfeng Zhang. 2020. Understanding Echo Chambers in E-Commerce Recommender Systems. In *SIGIR*. ACM, 2261–2270.
- [20] Matthew Gentzkow and Jesse M Shapiro. 2011. Ideological segregation online and offline. *The Quarterly Journal of Economics* 126, 4 (2011), 1799–1839.
- [21] Siyuan Guo, Lixin Zou, Yiding Liu, Wenwen Ye, Suqi Cheng, Shuaiqiang Wang, Hechang Chen, Dawei Yin, and Yi Chang. 2021. Enhanced Doubly Robust Learning for Debiasing Post-Click Conversion Rate Estimation. In *SIGIR*. ACM, 275–284.
- [22] Xiangnan He and Tat-Seng Chua. 2017. Neural Factorization Machines for Sparse Predictive Analytics. In *SIGIR*. ACM, 355–364.
- [23] Yoshinori Hijikata, Yuki Kai, and Shogo Nishida. 2012. The relation between user intervention and user satisfaction for information recommendation. In *SAC*. ACM, 2002–2007.
- [24] Dietmar Jannach, Sidra Naveed, and Michael Jugovac. 2016. User control in recommender systems: Overview and interaction challenges. In *EC-Web*. Springer, 21–33.
- [25] Ray Jiang, Silvia Chiappa, Tor Lattimore, András György, and Pushmeet Kohli. 2019. Degenerate feedback loops in recommender systems. In *AIES*. AAAI press, 383–390.
- [26] Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. 2017. Unbiased Learning-to-Rank with Biased Feedback. In *WSDM*. ACM, 781–789.
- [27] Bart P Knijnenburg, Niels JM Reijmer, and Martijn C Willemsen. 2011. Each to his own: how different users call for different interaction methods in recommender systems. In *RecSys*. ACM, 141–148.
- [28] Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. 2021. Towards personalized fairness based on causal notion. In *SIGIR*. ACM, 1054–1063.
- [29] Masoud Mansoury, Himan Abdollahpouri, Mykola Pechenizkiy, Bamshad Mobasher, and Robin Burke. 2020. Feedback loop and bias amplification in recommender systems. In *CIKM*. ACM, 2145–2148.
- [30] Marco Morik, Ashudeep Singh, Jessica Hong, and Thorsten Joachims. 2020. Controlling Fairness and Bias in Dynamic Learning-to-Rank. In *SIGIR*. ACM, 429–438.
- [31] Tien T Nguyen, Pik-Mai Hui, F Maxwell Harper, Loren Terveen, and Joseph A Konstan. 2014. Exploring the Filter Bubble: The Effect of Using Recommender Systems on Content Diversity. In *WWW*. ACM, 677–686.
- [32] Liqiang Nie, Wenjie Wang, Richang Hong, Meng Wang, and Qi Tian. 2019. Multimodal dialog system: Generating responses via adaptive decoders. In *MM*. ACM, 1098–1106.
- [33] Eli Pariser. 2011. *The filter bubble: What the Internet is hiding from you*. Penguin UK.
- [34] Judea Pearl. 2009. *Causality*. Cambridge university press.
- [35] Steffen Rendle. 2010. Factorization Machines. In *ICDM*. IEEE, 995–1000.
- [36] Donald B Rubin. 2005. Causal Inference Using Potential Outcomes: Design, Modeling, Decisions. *JASA* 100, 469 (2005), 322–331.
- [37] Yuta Saito, Suguru Yaginuma, Yuta Nishino, Hayato Sakata, and Kazuhide Nakata. 2020. Unbiased Recommender Learning from Missing-Not-At-Random Implicit Feedback. In *WSDM*. ACM, 501–509.
- [38] J Ben Schafer, Joseph A Konstan, and John Riedl. 2002. Meta-recommendation systems: user-controlled integration of diverse recommendations. In *CIKM*. ACM, 43–51.
- [39] Harald Steck. 2018. Calibrated recommendations. In *RecSys*. ACM, 154–162.
- [40] Yueming Sun and Yi Zhang. 2018. Conversational recommender system. In *SIGIR*. ACM, 235–244.
- [41] Cass R Sunstein. 2009. *Going to extremes: How like minds unite and divide*. Oxford University Press.
- [42] Kirsten Swearingen and Rashmi Sinha. 2001. Beyond algorithms: An HCI perspective on recommender systems. In *SIGIR*, Vol. 13. ACM, 1–11.
- [43] Taavi T Tajala, Martijn C Willemsen, and Joseph A Konstan. 2018. Movieexplorer: building an interactive exploration tool from ratings and latent taste spaces. In *SAC*. ACM, 1383–1392.
- [44] Juntao Tan, Shuyuan Xu, Yingqiang Ge, Yunqi Li, Xu Chen, and Yongfeng Zhang. 2021. Counterfactual explainable recommendation. In *CIKM*. ACM, 1784–1793.
- [45] Antonela Tommasel, Juan Manuel Rodriguez, and Daniela Godoy. 2021. I Want to Break Free! Recommending Friends from Outside the Echo Chamber. In *RecSys*. ACM, 23–33.
- [46] Khanh Hiep Tran, Azin Ghazimatin, and Rishiraj Saha Roy. 2021. Counterfactual Explanations for Neural Recommenders. In *SIGIR*. ACM, 1627–1631.
- [47] Chun-Hua Tsai and Peter Brusilovsky. 2018. Beyond the ranked list: User-driven exploration and diversification of social recommendation. In *23rd international conference on intelligent user interfaces*. 239–250.
- [48] Wenjie Wang, Fuli Feng, Xiangnan He, Xiang Wang, and Tat-Seng Chua. 2021. Deconfounded Recommendation for Alleviating Bias Amplification. In *KDD*. ACM, 1717–1725.

- [49] Wenjie Wang, Fuli Feng, Xiangnan He, Hanwang Zhang, and Tat-Seng Chua. 2021. Click can be Cheating: Counterfactual Recommendation for Mitigating Clickbait Issue. In *SIGIR*. ACM.
- [50] Wenjie Wang, Xinyu Lin, Fuli Feng, Xiangnan He, Min Lin, and Tat-Seng Chua. 2022. Causal Representation Learning for Out-of-Distribution Recommendation. In *WWW*. ACM.
- [51] Yixin Wang, Dawen Liang, Laurent Charlin, and David M Blei. 2018. The Deconfounded Recommender: A Causal Inference Approach to Recommendation. In *arXiv:1808.06581*.
- [52] Zhenlei Wang, Jingsen Zhang, Hongteng Xu, Xu Chen, Yongfeng Zhang, Wayne Xin Zhao, and Ji-Rong Wen. 2021. Counterfactual Data-augmented Sequential Recommendation. In *SIGIR*. ACM, 347–356.
- [53] Yao Wu, Jian Cao, Guandong Xu, and Yudong Tan. 2021. TFROM: A Two-Sided Fairness-Aware Recommendation Model for Both Customers and Providers. In *SIGIR*. ACM, 1013–1022.
- [54] Ke Yang and Julia Stoyanovich. 2017. Measuring fairness in ranked outputs. In *SSDBM*. ACM, 1–6.
- [55] Mengyue Yang, Quanyu Dai, Zhenhua Dong, Xu Chen, Xiuqiang He, and Jun Wang. 2021. Top-N Recommendation with Counterfactual User Preference Simulation. In *CIKM*. ACM, 2342–2351.
- [56] Yongfeng Zhang and Xu Chen. 2020. Explainable Recommendation: A Survey and New Perspectives. *Found. Trends Inf. Retr.* 14, 1 (2020), 1–101.
- [57] Yang Zhang, Fuli Feng, Xiangnan He, Tianxin Wei, Chonggang Song, Guohui Ling, and Yongdong Zhang. 2021. Causal Intervention for Leveraging Popularity Bias in Recommendation. In *SIGIR*. ACM, 11–20.
- [58] Yu Zheng, Chen Gao, Liang Chen, Depeng Jin, and Yong Li. 2021. DGCN: Diversified Recommendation with Graph Convolutional Networks. In *WWW*. ACM, 401–412.
- [59] Cai-Nicolas Ziegler, Sean M McNee, Joseph A Konstan, and Georg Lausen. 2005. Improving recommendation lists through topic diversification. In *WWW*. ACM, 22–32.
- [60] Hao Zou, Peng Cui, Bo Li, Zheyang Shen, Jianxin Ma, Hongxia Yang, and Yue He. 2020. Counterfactual Prediction for Bundle Treatment. In *NeurIPS*.