# Causal Inference in Recommender Systems: A Survey of Strategies for Bias Mitigation, Explanation, and Generalization

Yaochen Zhu, Jing Ma, and Jundong Li

**Abstract** In the era of information overload, recommender systems (RSs) have become an indispensable part of online service platforms. Traditional RSs estimate user interests and predict their future behaviors by utilizing correlations in the observational user historical activities, user profiles, and the content of interacted items. However, since the inherent causal reasons that lead to the observed user behaviors are not considered, multiple types of biases could exist in the generated recommendations. In addition, the causal motives that drive user activities are usually entangled in these RSs, where the explainability and generalization abilities of recommendations cannot be guaranteed. To address these drawbacks, recent years have witnessed an upsurge of interest in enhancing traditional RSs with causal inference techniques. In this survey, we provide a systematic overview of causal RSs and help readers gain a comprehensive understanding of this promising area. We start with the basic concepts of traditional RSs and their limitations due to the lack of causal reasoning ability. We then discuss how different causal inference techniques can be introduced to address these challenges, with an emphasis on debiasing, explainability promotion, and generalization improvement. Furthermore, we thoroughly analyze various evaluation strategies for causal RSs, focusing especially on how to reliably estimate their performance with biased data if the causal effects of interests are unavailable. Finally, we provide insights into potential directions for future causal RS research.

Yaochen Zhu

Department of Electrical and Computer Engineering, University of Virginia, e-mail: `uqp4qh@virginia.edu`

Jing Ma

Department of Computer Science, University of Virginia, e-mail: `jm3mr@virginia.edu`

Jundong Li

Department of Electrical and Computer Engineering, Department of Computer Science, and School of Data Science, University of Virginia, e-mail: `jl6qk@virginia.edu`

# 1 Introduction

With information growing exponentially on the web, recommender systems (RSs) are playing an increasingly pivotal role in modern online services, due to their ability to automatically deliver items[1] to users based on their personalized interests. Traditional RSs can be mainly categorized into three classes [9]: Collaborative filtering-based methods [10], content-based methods [11], and hybrid methods [12]. Collaborative filtering-based RSs estimate user interests and predict their future behaviors by exploiting their past activities, such as browsing, clicking, purchases, etc. Content-based methods, on the other hand, predict new recommendations by matching user interests with item content. Hybrid methods combine the advantages of both worlds, where collaborative information and user/item feature information are comprehensively considered to generate more accurate recommendations.

Although recent years have witnessed substantial achievements for all three classes of RSs introduced above, a great limitation of these methods is that they can only estimate user interests and predict future recommendations based on correlations in the observational user historical behaviors and user/item features, which guarantee no causal implications [13, 14]. For example, a collaborative filtering-based RS may discover that several drama shows from a certain genre *tend to* have high ratings from a group of users, and conclude that we should keep recommending drama shows from the same genre to these users. But there is an important question: Are the high ratings caused by the fact that the users indeed like drama shows from this genre, or they were limitedly exposed to drama shows from the same genre (i.e., exposure bias), and if given a chance, they would prefer something new to watch? In addition, a content-based RS may observe that micro-videos with certain features *are associated with* more clicks and conclude that these features may reflect the current trend of user interests. But are the clicks because these micro-videos tend to have sensational titles as clickbait where users could be easily deceived? Moreover, if the titles of these micro-videos are changed to the ones that reflect their true content, would users still click them? The above questions are causal in nature because they either ask about the effects of an intervention (e.g., what the rating would be if a new drama show **is made exposed** to the user) or a counterfactual outcome (e.g., would the user still click a micro-video if its title **had been changed** to faithfully reflect the content), rather than mere associations in the observational data. According to Pearl [15], these questions lie on Rungs 2 and 3 of the Ladder of Causality, i.e., interventional and counterfactual reasoning, and they cannot be answered by traditional RSs that reason only with associations, which lie on Rung 1 of the ladder.

Why are these causal questions important for RSs? The first reason is that failing to address them may easily incur bias in recommendations, which can get unnoticed for a long time. If the collaborative filtering-based RSs mentioned above mistake exposure bias for user interests, they would amplify the bias by continuously recommending users with similar items; eventually, recommendations will lose serendipity,

---

[1] We use the term item in a broad sense to refer to anything recommendable to users, such as news [1], jobs [2], articles [3], music [4], movies [5], micro-videos [6], PoIs [7], hashtags [8], etc.

and users' online experience can be severely degraded. Similarly, for the content-based micro-video RSs, if they cannot distinguish clicks due to user interests from the ones deceived by clickbait, they may over-recommend micro-videos with sensational titles, which is unfair to the uploaders of high-quality micro-videos who put much effort into designing the content. In addition, understanding the cause of user activities can help improve the explainability of recommendations. Consider the causal question of whether a user purchases an item due to its quality or its low price. Pursuing the causal explanations behind user behaviors can help service providers to enhance the RS algorithm based on users' personalized preferences. Finally, causal inference allows us to identify and base recommendations on causal relations that are stable and invariant, while discarding other correlations that are undesirable or susceptible to change. Take restaurant recommendations as an example. Users can choose a restaurant because of its convenience (e.g., going to a nearby fast food shop to quickly grab a bite, but they do not necessarily like it, a non-stable correlation) or due to their personal interests (e.g., traveling far away for a hot-pot restaurant, a stable causal relation). If an RS can properly disentangle users' intent that causally affects their previous restaurant visits, even if the convenience levels of different restaurants may change due to various internal or external reasons such as users' moving to a new place, the system can still adapt well to the new situation. From this aspect, the generalization ability of the causal RSs can be substantially improved.

This survey provides a systematic overview of recent advances in causal RS research. The organization is illustrated in Fig. 1. We start with the fundamental concepts of traditional RSs and their limitation of correlational reasoning in Section 2. Then Section 3 recaps two important causal inference paradigms in machine learning and statistics, and shows their connections with the recommendation task. Section 4 thoroughly discusses how different causal inference techniques can be introduced to address the limitations of traditional RSs, with an emphasis on debiasing, explainability promotion, and generalization improvement. Section 5 summarizes the general evaluation strategies for causal RSs. Finally, Sections 6 and 7 discuss prospective open questions and future directions for causal RSs and conclude this survey.

## 2 Recommender System Basics

To keep this survey compact, we confine our discussions to simple RSs with $I$ users and $J$ items. The main data for the RSs, i.e., users' historical behaviors, are represented by a user-item rating matrix $\mathbf{R} \in \mathbb{R}^{I \times J}$, where a non-zero element $r_{ij}$ denotes user $i$'s rating to item $j$, and a zero element $r_{ik}^0$ indicates the rating is missing [2]. To make the discussions of RSs compatible with causal inference, we take a probabilistic view of $\mathbf{R}$ [17], where $r_{ij}$ is assumed to be the realized value of the

---

[2] We use rating to refer to any user-item interaction that can be represented by a numerical value. This includes both explicit feedback such as likes/dislikes, and implicit feedback such as views and clicks. When $r_{ij}$ represents implicit feedback, the missing elements $r_{ik}^0$ in $\mathbf{R}$ may be used as weak negative feedback in the training phase [16]. This may complicate the causal problems. Therefore, we assume RSs are trained on observed ratings to simplify the discussion unless specified otherwise.
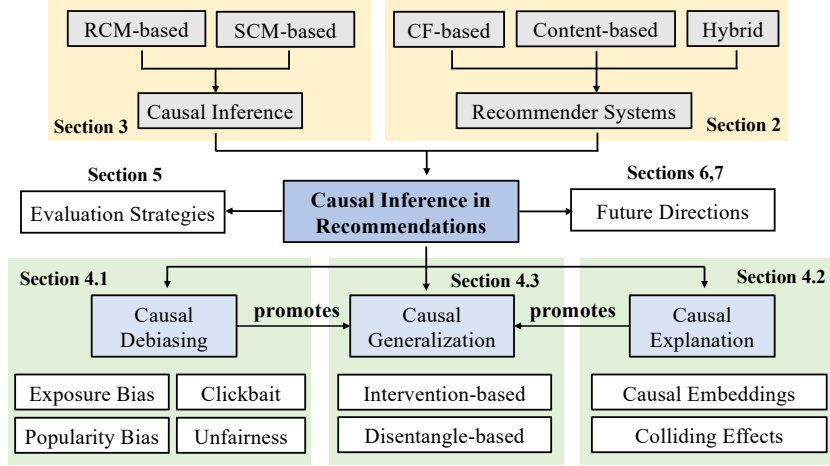
Fig. 1: An overview of the structure of this survey and connections between different sections.

random variable $R$ dependent on user $i$ and item $j$[3]. In addition to $\mathbf{R}$, an RS usually has access to side information like user features $\mathbf{f}_i^u \in \mathbb{R}^{K_F^u}$, such as her age, gender, location, etc., or item features $\mathbf{f}_j^v \in \mathbb{R}^{K_F^v}$, such as its content and textual description. $K_F^u$ and $K_F^v$ are the dimensions of user and item features, respectively. **The main purpose of an RS** is to predict users' ratings for previously uninteracted items (i.e., the missing values $r_{ik}^0$ in $\mathbf{R}$) based on the observed ratings $r_{ij}$ in $\mathbf{R}$ and the available user and item side information such as $\mathbf{f}_i^u$ and $\mathbf{f}_j^v$, such that new relevant items can be properly recommended based on users' personalized interests.

## 2.1 Collaborative Filtering

Collaborative filtering-based RSs recommend new items by leveraging user ratings in the past. They generally consider the ratings $r_{ij}$ as being generated from a user latent variable $\mathbf{u}_i \in \mathbb{R}^K$ that represents user interests and an item latent variable $\mathbf{v}_j \in \mathbb{R}^K$ that encodes the item attributes (i.e., item latent semantic information), where $K$ is the dimension of the latent space. Here we list three widely-used collaborative filtering-based RSs, which will be frequently used as examples in this survey:

- **Matrix Factorization (MF)** [18]. MF models $r_{ij}$ with the inner product between $\mathbf{u}_i$ and $\mathbf{v}_j$, where $r_{ij} \sim \mathcal{N}(\mathbf{u}_i^T \cdot \mathbf{v}_j, \sigma_{ij}^2)$ and $\sigma_{ij}^2$ is the predetermined variance[4].

---

[3] However, we do not distinguish random variables and their specific realizations if there is no risk of confusion. For simplicity, we assume $R$ to be Gaussian unless specified otherwise.

[4] For works that do not explicitly treat $r_{ij}$ as a random variable, we assume it follows a Gaussian distribution with zero variance. The generative process then becomes as $r_{ij} = \mathbf{u}_i^T \cdot \mathbf{v}_j$.

- **Deep Matrix Factorization (DMF) [19]**. DMF extends MF by applying deep neural networks (DNNs) [20], i.e., $f_{nn}^u, f_{nn}^v : \mathbb{R}^K \to \mathbb{R}^{K'}$, to $\mathbf{u}_i$ and $\mathbf{v}_j$, where the ratings are assumed to be generated as $r_{ij} \sim \mathcal{N}(f_{nn}^u(\mathbf{u}_i)^T \cdot f_{nn}^v(\mathbf{v}_j), \sigma_{ij}^2)$.

- **Auto-encoder (AE)-based RSs [21, 22]** model user $i$'s ratings to all $J$ items as $\mathbf{r}_i \sim \mathcal{N}(f_{nn}^u(\mathbf{u}_i), \sigma_i^2 \cdot \mathbf{I}_K)$, where $f_{nn}^u : \mathbb{R}^K \to \mathbb{R}^J$ is a DNN and item latent variables $\mathbf{v}_j$ for all items are implicit in last layer weights of the decoder [23].

In the training phase, the models learn the latent variables $\mathbf{u}_i$, $\mathbf{v}_j$ and the associated function $f_{nn}$ by fitting on the **observed ratings** $r_{ij}$ (e.g., via maximum likelihood estimation, which essentially estimates the conditional distribution $p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)$ from the observational data [24]). Afterward, we can use them to predict new ratings for previously uninteracted items $k$, e.g., $\hat{r}_{ik}^{\text{MF}} = \mathbf{u}_i^T \cdot \mathbf{v}_k$ for MF, $\hat{r}_{ik}^{\text{DMF}} = f_{nn}^u(\mathbf{u}_i)^T \cdot f_{nn}^v(\mathbf{v}_k)$ for DMF, and $\hat{r}_{ik}^{\text{AE}} = f_{nn}^u(\mathbf{u}_i)_k$ for AE-based RSs, where the top ones that best match users' interests can be selected as recommendations.

> **Traditional collaborative filtering-based RSs reasons with correlations.** Ideally, we would expect $\mathbf{u}_i$, $\mathbf{v}_j$ and $f_{nn}$ to capture the causal influence of user interests and item attributes on ratings, i.e., what the rating would be if item $j$ is made exposed to user $i$ [24]. However, since the collected rating data are **observational** rather than experimental, what actually learned by $\mathbf{u}_i$, $\mathbf{v}_j$, and $f$ are the co-occurrence patterns in users' past behaviors, which guarantee no causal implications. Consequently, spurious correlations and biases can be captured by the model, which will be amplified in future recommendations [25]. Furthermore, the learned user latent variable $\mathbf{u}_i$ generally entangles different factors that causally determine user interests. From this perspective, the explainability and generalization of these methods cannot be guaranteed.

## 2.2 Content-Based Recommender Systems

Personalized content-based RSs (CBRSs) estimate user interests based on the features of the items they have interacted with. These models typically encode user interests into user latent variables $\mathbf{u}_i \in \mathbb{R}^K$ and assume that the ratings are generated by matching user interests with item content, i.e., $r_{ij} \sim \mathcal{N}(f(\mathbf{u}_i, \mathbf{f}_j^v), \sigma_{ij})$, where $f$ is a matching function. The training of personalized CBRSs follow similar steps as collaborative filtering, where $\mathbf{u}_i$ and $f$ are learned by fitting on the **observed ratings** (which essentially estimates the conditional distribution $p(r_{ij}|\mathbf{u}_i, \mathbf{f}_j^v)$ from the observational data), and new ratings can be predicted by $\hat{r}_{ik} = f(\mathbf{u}_i, \mathbf{f}_k^v)$. The key step of building a CBRS is to create item features $\mathbf{f}_j^v$ that can best reflect user interests, which crucially depends on the item being recommended. For example, for micro-videos, the visual, audio, and textual modalities are comprehensively considered such that users' interest in different aspects of a micro-video can be well captured [26].

**Traditional content-based RSs cannot model the causal influence of user interests $\mathbf{u}_i$ and item content $\mathbf{f}_j^v$ on user rating $r_{ij}$.** The reason is that, factors other than users' interests in the item content, such as users' being deceived by clickbaits (e.g., sensational titles of micro-videos) [27], etc., can create an undesirable association between item content $\mathbf{f}_j^v$ and user ratings $r_{ij}$ in the observed dataset, where the bias can be captured by the user latent variables $\mathbf{u}_i$ and the matching function $f$, and perpetuates into future recommendations.

## 2.3 Hybrid Recommendation

Hybrid RSs combine user/item side information with collaborative filtering to enhance the recommendations. A commonly-used hybrid strategy is to augment user and item latent variables $\mathbf{u}_i$ and $\mathbf{v}_j$ with user/item side information $\mathbf{f}_i^v$ and $\mathbf{f}_j^v$ in existing collaborative filtering methods by replacing $\mathbf{u}_i$ and $\mathbf{v}_j$ with $\mathbf{u}_i^+ = [\mathbf{u}_i || \mathbf{f}_i^u]$ and $\mathbf{v}_j^+ = [\mathbf{v}_j || \mathbf{f}_j^v]$ in MF, DMF, and AE-based RSs, where $[\cdot || \cdot]$ represents vector concatenation [28, 29]. The dimensions of $\mathbf{u}_i$ and $\mathbf{v}_j$ that encode the collaborative information are adjusted accordingly to make $\mathbf{u}_i^+$ and $\mathbf{v}_j^+$ compatible in the model. Another important class of hybrid RS is the factorization machine (FM) [30] and its extensions like [31, 32], which can be viewed as learning a bi-linear function $f_{fm}$ where the ratings are generated by $r_{ij} \sim \mathcal{N}(f_{fm}(\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_i^u, \mathbf{f}_j^v), \sigma_{ij}^2)$.

**Simple hybrid strategies cannot break the correlational reasoning limitation of collaborative filtering and content-based RSs**, because the objective of the hybridization is still to improve the models' fitting on the observational user historical behaviors (i.e., estimating conditional distribution $p(r_{ij} | \mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_i^u, \mathbf{f}_j^v)$ from the data), where the causal reasons that lead to the observed user behaviors are not considered. However, the idea of introducing extra user/item side information is important for building causal RSs. The reason is that, combined with the domain knowledge of human experts, the side information can help form more comprehensive causal relations among the variables of interests, such as user interests, item attributes, historical ratings, and other important covariates that may lead to spurious correlations and biases, which is usually a crucial step for causal reasoning in recommendations.

## 3 Causal Recommender Systems: Preliminaries

In the previous section, we discussed the recommendation strategies of the traditional RSs and their limitations due to correlational reasoning on observational user

behaviors. In this section, we introduce two causal inference frameworks, i.e., Rubin's potential outcome framework (also known as the Rubin causal model, RCM) [33] and Pearl's structural causal model (SCM) [34], in the context of RSs, aiming to provide a theoretically rigorous basis for reasoning with correlation and causation in recommendations. We show that both RCM and SCM are powerful frameworks to build RSs with causal reasoning ability (i.e., causal RSs), but they are best suited for different tasks and questions. The discussions in this section serve as the foundation for more in-depth discussions of the state-of-the-art causal RS models in Section 4.

## 3.1 Rubin's Potential Outcome Framework

### 3.1.1 Motivation of Applications in RSs

To understand the correlational reasoning nature of traditional RSs, we note that naively fitting models on the observed ratings can only answer the question "what the rating would be **if we observe an item was exposed to the user**". Since item exposure is not randomized in the collected dataset [5], the predicate "the item was exposed to the user" *per se* contains extra information about the user-item pair (e.g., the item could be more popular than other non-exposed items), which cannot be generalized to the rating predictions of **arbitrary** user-item pairs. Therefore, what RS asks is essentially an interventional question (and therefore a causal inference question), i.e., what the rating would be **if an item is made exposed to the user**. To address this question, RCM-based RSs draw inspiration from clinical trials, where exposing a user to an item is compared to subjecting a patient to a treatment, and the user ratings are analogous to the outcomes of the patients after the treatment [39, 40]. Accordingly, RCM-based RSs aim to estimate the causal effects of the treatments (exposing a user to an item) on the outcomes (user ratings), despite the possible correlations between the treatment assignment and the outcome observations [39].

### 3.1.2 Definitions and Objectives

We first introduce necessary symbols and definitions to connect RCM with RSs. We consider the unit as the user-item pair $(i, j)$ that can receive the treatment "exposing user $i$ to item $j$", and the population as all user-item pairs $\mathcal{PO} = \{(i, j), 1 \leq i, j \leq I, J\}$ [41]. We start by using a binary scalar $a_{ij}$ to denote the exposure status of item $j$ for user $i$, i.e., the assigned treatment. We further define the **rating potential outcome** $r_{ij}(a_{ij} = 1)$ as user $i$'s rating to item $j$ if the item is made exposed to the user and $r_{ij}(a_{ij} = 0)$ as the rating if the item is not exposed [42]. For user $i$, if she rated item $j$, we observe $r_{ij}(a_{ij} = 1) = r_{ij}$. Otherwise, we observe the baseline potential outcome $r_{ij}(a_{ij} = 0) = 0$, which is usually ignored in debias-oriented

---

[5] which can be attributed to multiple reasons such as users' self-search [35], the recommendations of previous models [36], the position where the items are displayed [37], item popularity [38], etc. Generally, RCM-based causal RSs are agnostic to the specific reason that causes the exposure bias.

(a) Observed Ratings        (b) Rating Potential Outcomes        (c) Predicted Ratings
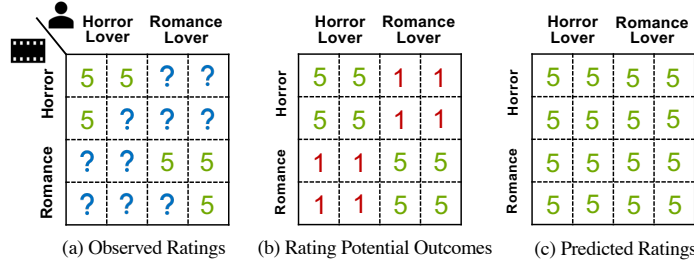
Fig. 2: A classical example of exposure bias in RSs [43]. The example is composed of two horror lovers who always rate horror movies with five while hating romance movies, and two romance lovers would who do exactly the opposite. (a) shows the observed ratings $r_{ij}$. (b) shows the rating potential outcomes $r_{ij}(a_{ij} = 1)$. (c) shows the rating predictions of an RS that maximizes the likelihood of the observed ratings in (a), but the RS is bad because it predicts all ratings to five.

causal RS research [39, 43][6]. Similar to clinical trials, we can define the treatment group $\mathcal{T} = \{(i, j) : a_{ij} = 1\}$ as the set of user-item pairs where user $i$ is exposed to item $j$, and define the non-treatment group $\mathcal{NT} = \{(i, k) : a_{ik} = 0\}$ accordingly. **The purpose of RSs, under the RCM framework**, can be framed as utilizing the observed ratings from units in the treatment group $\mathcal{T}$ to unbiasedly estimate the rating potential outcomes for units from the population $\mathcal{PO}$, despite the possible correlations between item exposures $a_{ij}$ and user ratings $r_{ij}$ in the collected data.

### 3.1.3 Causal Analysis of Traditional RSs

Traditional RSs naively train a rating prediction model that best fits the ratings in the treatment group $\mathcal{T}$ (e.g., via maximum likelihood introduced in Section 2) to estimate the unobserved rating potential outcomes $r_{ij}(a_{ij} = 1)$ for user-item pairs in $\mathcal{NT}$ [46], which neglect the fact that exposure bias can lead to a systematic difference in the distribution of $r_{ij}(a_{ij} = 1)$ between $\mathcal{T}$ and $\mathcal{NT}$. For example, users tend to rate items they like in reality, which could lead to the following spurious correlation between item exposure $a_{ij}$ and rating potential outcome $r_{ij}(a_{ij} = 1)$:

$$p(r_{ij}(a_{ij} = 1) \text{ is high}|a_{ij} = 1) > p(r_{ij}(a_{ij} = 1) \text{ is high}|a_{ij} = 0), \quad (1)$$

i.e., users who have rated an item $j$ may have systematically higher ratings than users who haven't rated it yet. In this case, traditional RSs may have a tendency to overestimate the ratings for units in $\mathcal{NT}$ (see Fig. 2 for an intuitive example). Theoretically, RCM attributes the exposure bias in the collected dataset to the violation of the **unconfoundedness assumption** [33] defined as follows:

$$r_{ij}(a_{ij} = 1) \perp a_{ij}. \quad (2)$$

---

[6] In the uplift evaluation of RSs that aims to estimate how recommendations change user behaviors [44], $r_{ij}(a_{ij} = 0)$ may be used to represent user $i$'s rating to item $j$ through self-searching [45].

The rationale is that, if Eq. (2) holds, the exposure of user $i$ to item $j$ (i.e., $a_{ij}$) is independent of the rating potential outcome $r_{ij}(a_{ij} = 1)$, which implies that $r_{ij}(a_{ij} = 1)$ in $\mathcal{T}$ and $\mathcal{NT}$ follows the same distribution. Therefore, the exposure of the items is randomized, and exposure bias such as Eq. (1) will not exist [42].

### 3.1.4 Potential Outcome Estimation with the RCM Framework

One classic solution from the RCM-based framework to address the exposure bias is that we find user and item covariates $C$, such that in each data stratum specified by $C = \mathbf{c}$, users' exposure to items are randomized [33]. The property of the covariates $C$ can be formulated as the conditional unconfoundedness assumption as follows:

$$r_{ij}(a_{ij} = 1) \perp a_{ij} \mid \mathbf{c}. \tag{3}$$

$C$ is sometimes non-rigorously referred to as **confounder** in the literature, but we will see its formal definition in the next subsection. If Eq. (3) holds, the item exposures are independent of the rating potential outcomes in each data stratum specified by $C = \mathbf{c}$, and the exposure bias can be attributed solely to the discrepancy in the distribution of the covariates $C = \mathbf{c}$ between the treatment group $\mathcal{T}$ and the population $\mathcal{PO}$, i.e., $p(\mathbf{c}|a_{ij} = 1)$ and $p(\mathbf{c})$[7] Therefore, we can reweight the observed ratings in $\mathcal{T}$ based on the covariates $C$ to address the bias, such that they can be viewed as pseudo randomized samples. This leads to inverse propensity weighting (IPW), which eliminates the exposure bias from the data's perspective [39]. In addition, we can also adjust the influence of $C$ in the RS model, where the exposure bias is addressed from the model side [42]. Both methods will be discussed in Section 4.1.1.

### ❗ Attention: Extra Assumptions Required by Most RCM-based RSs

In addition to unconfoundedness, most RCM-based RS need two extra assumptions to identify the causal effects of item exposures on ratings: (1) The **stable unit treatment assumption (SUTVA)**, which states that items exposed to one user cannot affect ratings of another user. (2) The **positivity assumption**, which states that every user has a positive chance of being exposed to every item [33]. For RCM-based causal RSs introduced in this survey, these two assumptions are tacitly accepted.

---

[7] ❗ We can gain an intuition of this claim from Fig. 2. Suppose covariates $C$ represent the two-dimensional features (user type, movie type). Given $C = \mathbf{c}$, $r_{ij}(a_{ij} = 1) \perp a_{ij} \mid \mathbf{c}$ described in Eq. (3) is satisfied because in each data stratum specified by $C = \mathbf{c}$ (i.e., the four $2 \times 2$ blocks in Fig. 2-(b)), $r_{ij}(a_{ij} = 1)$ is constant. Fig. 2-(a) shows that for the treatment group $\mathcal{T}$, $p(\mathbf{c}|a_{ij} = 1) = 1/2$ for $\mathbf{c} \in C_1 = \{(\text{horror fan, horror movie}), (\text{romance fan, romance movie})\}$ and $p(\mathbf{c}|a_{ij} = 1) = 0$ for $\mathbf{c} \in C_2 = \{(\text{horror fan, romance movie}), (\text{romance fan, horror movie})\}$. In contrast, for the population $\mathcal{PO}$, $p(\mathbf{c}) = 1/4$ for $\mathbf{c} \in C_1 \cup C_2$. Therefore, in the treatment group $\mathcal{T}$, user-item pairs with covariates in $C_1$ are over-represented, while those with covariates in $C_2$ are under-represented. However, we also note that this case is too extreme to be addressed by RCM, as $p(\mathbf{c}|a_{ij} = 1) = 0$ for $C \in C_2$ **violates** the positivity assumption mentioned in the above attention box.

## 3.2 Pearl's Structural Causal Model

### 3.2.1 Motivation of Applications in RS

Different from RCM that uses rating *potential outcomes* to reason with causality and attributes the biases in observed user behaviors to non-randomized item exposures, Pearl's structural causal model (SCM) delves deep into the causal mechanism that generates the *observed outcomes* (and biases) and represents it with a causal graph $G = (\mathcal{N}, \mathcal{E})$. The nodes $\mathcal{N}$ specify the variables of interests, which in the context of RS could be user interests $U$, item attributes $V$, observed ratings $R$, and other important covariates $C$, such as item popularity, user features, etc[8]. The directed edges $\mathcal{E}$ between nodes represent their causal relations determined by researchers' domain knowledge. Each node $X \in \mathcal{N}$ is associated with a structural equation $p_G(X|Pa(X))$[9], which describes how the parent nodes $Pa(X)$ causally influence $X$ (i.e., the response of $X$ when setting nodes in $Pa(X)$ to specific values)

Although RCM and SCM are generally believed to be fundamentally equivalent [34], both have their unique advantages. Compared to RCM, the key advantage of SCM is that causal graph offers an intuitive and straightforward way to encode and communicate domain knowledge and substantive assumptions of researchers, which is beneficial even for the RCM-based RSs [42]. Furthermore, SCM is more flexible as it can represent and reason with the causal effects between any subset of nodes in the causal graph (e.g., between two causes $U, V$ and one outcome $R$), as well as the causal effects along specific paths (e.g., $U \rightarrow R$ and $U_c \rightarrow R$). Therefore, SCMs are broadly applicable to multiple problems in RSs (not limited to exposure bias), such as clickbait bias, unfairness, entanglement, domain adaptation, etc [14].

---

### ! Attention: Two Caveats of SCM-based Causal RSs.

There are two caveats of SCM-based causal RSs. (1) Causal graphs for RSs often involve user, item **latent variables** $U, V$ that encode user interests and item attributes. Most works infer them alongside the estimation of structural equations and treat them as if they were observed when analyzing the causal relations. Alternatively, this can be viewed as representing users and items with their IDs (i.e., $i$ and $j$) in the causal graph and subsuming the embedding process into the structural equations [47]. (2) Generally, the causal graph should describe the causal mechanism that generates the **observed data**, because it allows us to distinguish invariant, causal relations from undesirable correlations. For example, we may argue that item popularity $C$ should be determined by item attributes $V$, i.e., $V \rightarrow C$. But to describe the generation of the observed ratings, causal relation $C \rightarrow V$ is usually assumed instead as item popularity causally influences the exposure probability of each item [48].

---

[8] In causal graphs, the subscripts $i$, $j$ for each node are omitted for simplicity.

[9] We also omit the mutually independent exogenous variables for each node and summarize their randomness into the structural equations with probability distributions [14]. Subscript $G$ is used to distinguish structural equations from other conditional relationships that can be inferred from $G$.
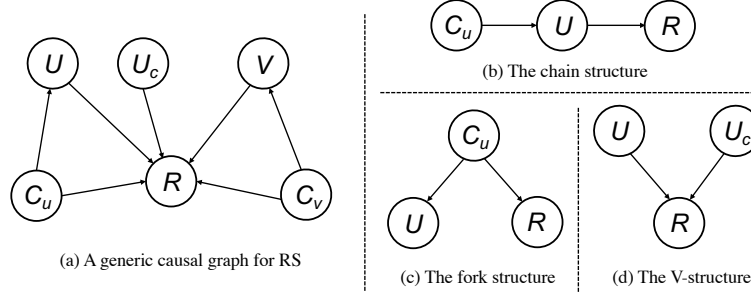
(a) A generic causal graph for RS

(b) The chain structure

(c) The fork structure

(d) The V-structure

Fig. 3: (a): A generic causal graph for RS that depicts the causal influence of user interests $U$, user conformity to the popularity trend $U_c$, and item attributes $V$ on the observed ratings $R$. Specifically, the causal paths $U \to R$ and $V \to R$ are confounded by $C_u$ and $C_v$, which represent user features and item popularity, respectively. (b)(c)(d): Three atomic structures identified from (a).

### 3.2.2 Atomic Structures of Causal Graphs

The structure of causal graphs represents researchers' domain knowledge regarding the causal generation process of the observational data, which is the key to distinguishing stable, causal relations from other undesirable correlations between variables of interest. Here, we use a generic causal graph applicable to RSs in Fig. 3-(a) as a running example to illustrate three atomic graph structures:

- **Chain,** e.g., $C_u \to U \to R$. In a chain, the successor node is assumed to be causally influenced by the ancestor nodes. In the example, $U$ is a direct cause of $R$, whereas $C_u$ indirectly influences $R$ via $U$ as a mediator.

- **Fork,** e.g., $U \leftarrow C_u \to R$. In the fork, $C_u$ is called a **confounder** as it causally influences two children $U$ and $R$. From a probabilistic perspective, $U$ and $R$ are **not** independent unless conditioned on the confounder $C_u$ [49]. This leads to the tricky part of a fork structure, i.e., **confounding effect** [34], where an unobserved $C_u$ can lead to spurious correlations between $U$ and $R$.

- **V-structure**, e.g., $U \to R \leftarrow U_c$. In the V-structure, $R$ is called a **collider** because it is under the causal influence of two parents, i.e., $U$ and $U_c$. An interesting property of the V-structure is the colliding effects [34], where observing $R$ creates a dependence on $U$ and $U_c$, even if they are marginally independent.

Confounders can lead to non-causal dependencies among variables in the observational dataset. This could introduce bias in traditional RSs, where the confounding effects are mistaken as causal relations. Confounding bias is a generic problem in RSs [24], which will be further analyzed in the following subsections. In addition, abstracted V-structure usually leads to the entanglement of causes, which could jeopardize the explainability of RSs. For example, a user's purchase of an item may be due to her interest, i.e., $U$, or her conformity to the popularity trend, i.e., $U_c$. Since most RSs summarize both into a user latent variable $U$, the V-structure $U \to R \leftarrow U_c$ is abstracted away, where the two causes of the purchase cannot be distinguished.

(a) SCM assumed by non-causal RS     (b) Confounded true SCM     (c) SCM under intervention
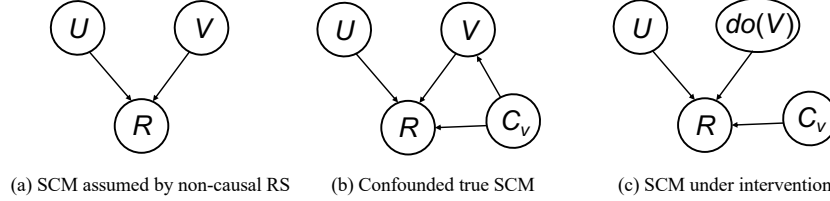
Fig. 4: (a): SCM assumed by non-causal collaborative filtering-based RS. (b): The confounded SCM that depicts the true data generation process. (c): SCM under intervention $do(V)$.


### 3.2.3 Causal Analysis of Traditional RSs

In this section, we investigate the susceptibility of traditional collaborative filtering-based RSs to the confounding bias. As discussed in Section 2.1, a commonality of these models is that they estimate conditional distribution $p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)$ from observed ratings and use it to predict new ratings. For $p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)$ to represent the causal influence of user interests $\mathbf{u}_i$ and item attributes $\mathbf{v}_j$ on ratings $r_{ij}$ (which, in the context of collaborative filtering, means the rating of any arbitrary item $j$ that is made exposed to user $i$ [24]), the causal graph $G_1$ of Fig. 4-(a) is tacitly assumed, i.e., no unobserved confounders for causal paths $U \rightarrow R$ and $V \rightarrow R$[10].

However, in reality, both $U \rightarrow R$ [25, 51] and $V \rightarrow R$ [52, 53] can be confounded, where the confounding effects can be implicitly captured by $p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)$ that bias future recommendations. To reveal the bias, we consider the scenario where the causal path $V \rightarrow R$ is confounded by $C_v$ (e.g., item popularity). We assume the causal path $C_v \rightarrow V$ denotes the causal influence of $C_v$ on the exposure probability of item $V$ [48]. In this case, the observed ratings are generated according to the causal graph $G_2$ in Fig. 4-(b). Utilizing the law of total probability, the conditional distribution $p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)$ estimated from the confounded data can be calculated as:

$$p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j) = \sum_{\mathbf{c}} p(\mathbf{c}|\mathbf{v}_j) \cdot p_{G_2}(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, \mathbf{c}) = \mathbb{E}_{p(C_v|\mathbf{v}_j)}[p_{G_2}(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, C_v)].$$
(4)

The issue of Eq. (4) is that, the $p(\mathbf{c}|\mathbf{v}_j)$ term is not causal (as we only have an edge $C_v \rightarrow V$ in the causal graph but not $V \rightarrow C_v$). In fact, $p(\mathbf{c}|\mathbf{v}_j)$ represents abductive reasoning because it infers the cause $\mathbf{c}$ (e.g., item popularity) from the effect $\mathbf{v}_j$ (i.e., item $j$ is exposed to user $i$) and uses the inferred $\mathbf{c}$ to support the prediction of the rating $r_{ij}$. However, such reasoning cannot be generalized to the rating prediction of an arbitrary item $\mathbf{v}_j$, i.e., an item that **is made exposed** to the user. In other words, uncontrolled confounder $C_v$ leaves open a **backdoor path** (i.e., non-causal path) between $V$ and $R$, such that non-causal dependence of $R$ on $V$ exists in the data, which can be captured by traditional RSs and bias future recommendations. [11]

---

[10] This corresponds to the case where item exposures are randomized (see the discussions in Section 3.1.3), as the user-item pair $(U, V)$ is not determined by other factors associated with $R$ [50].

[11] The similarity between this section and Section 3.1.1 shows us the connection between RCM-based and SCM-based causal RSs, where the claim that when item exposure is not randomized,

### 3.2.4 Causal Reasoning with SCM

To calculate the causal effect of $\mathbf{u}_i$ and $\mathbf{v}_j$ on $r_{ij}$, we should conduct **intervention** on $U$ and $V$. This means that we set $U$, $V$ to $\mathbf{u}_i$, $\mathbf{v}_j$ regardless of the values of their parent nodes in the causal graph, including the confounder $C_v$ (because these nodes determine the exposure of item $j$ to user $i$ in the observed data). SCM denotes the intervention with **do-operator** as $p(r_{ij}|do(\mathbf{u}_i, \mathbf{v}_j))$ to distinguish it from the conditional distribution $p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)$ that reasons with correlations in the observational data. Consider again the causal graph $G_2$ illustrated in Fig. 4-(b). The intervention on node $V$ can be realized by removing all the incoming edges for node $V$ and setting the structural equation $p_{G_2}(V|C_v)$ deterministically as $V = \mathbf{v}_j$, while other structural equations remain intact (Fig. 4-(c)). If the confounder $C_v$ can be determined and measured for each item, the interventional distribution $p(r_{ij}|do(\mathbf{u}_i, \mathbf{v}_j))$ can be directly calculated from the confounded data via **backdoor adjustment** [34] as:

$$p(r_{ij}|do(\mathbf{u}_i, \mathbf{v}_j)) = \sum_{\mathbf{c}} p_{G_2}(\mathbf{c}) \cdot p_{G_2}(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, \mathbf{c}) = \mathbb{E}_{p_{G_2}(C_v)}[p_{G_2}(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, C_v)],$$

(5)

which, compared with Eq. (4), blocks the abductive inference of $\mathbf{c}$ from $\mathbf{v}_j$, such that the causal influence of $\mathbf{u}_i$, $\mathbf{v}_j$ on $r_{ij}$ can be properly identified.

Backdoor adjustment requires all confounders to be determined and measured in advance, but there are other SCM-based causal inference methods that can estimate causal effects with unknown confounders, and we refer readers to [54, 55] for details. Moreover, causal graphs allow us to conduct other types of causal reasoning based on the encoded causal knowledge, such as debiasing for non-confounder-induced biases (e.g., clickbait bias and unfairness), causal disentanglement, and causal generalization [56]. These will be thoroughly discussed in the next section.

## 4 Causal Recommender Systems: The State-of-the-Art

Based on the preliminary knowledge of RSs and causal inference discussed in previous sections, we are ready to introduce the state-of-the-art causal RSs. Specifically, we focus on three important topics, i.e., bias mitigation, explainability promotion, and generalization improvement, as well as their inter-connections, where various limitations of traditional RSs due to correlational reasoning can be well addressed.

### 4.1 Causal Debiasing for Recommendations

The correlational reasoning of traditional RSs can inherit multiple types of biases in the observational user behaviors and amplify them in future recommendations [46].

---

"observing that an item was exposed to the user *per se* contains extra information about the user-item pair" is mathematically transformed into the abductive inference of $\mathbf{c}$ from $\mathbf{v}_j$ by $p(\mathbf{c}|\mathbf{v}_j)$.

The biases may result in various consequences, such as the discrepancy between offline evaluation and online metrics, loss of diversity, reduced recommendation quality, offensive recommendations, etc. Causal inference can distinguish stable causal relations from spurious correlations and biases that could negatively influence the recommendations, such that the robustness of recommendations can be improved.

### 4.1.1 Exposure Bias

Exposure bias in RSs broadly refers to the bias in observed ratings due to non-randomized item exposures. From the RCM's perspective, exposure bias can be defined as the bias *where users are favorably exposed to items depending on their expected ratings for them (i.e., rating potential outcomes)* [43]. Exposure bias occurs due to various reasons, such as users' self-search or the recommendation of the previous RSs [36], which leads to the down-weighting of items less likely to be exposed to users. Since item exposures can be naturally compared with treatments in clinical trials, we discuss the debiasing strategies with the RCM framework.

**Inverse Propensity Weighting (IPW).** IPW-based causal RSs aim to reweight the biased observed ratings $r_{ij}$ for user-item pairs in the treatment group, i.e., $\mathcal{T} = \{(i, j) : a_{ij} = 1\}$, to create pseudo randomized samples [57] for unbiased training of RS models that aim to predict the rating potential outcomes $r_{ij}(a_{ij} = 1)$ for the population $\mathcal{PO} = \{(i, j), 1 \leq i, j \leq I, J\}$. Intuitively, we can set the weight of $r_{ij}$ for units in $\mathcal{T}$ to be the inverse of item $j$'s exposure probability to user $i$, such that under-exposed items can be up-weighted and vice versa. If for each user-item pair, the covariates $\mathbf{c}$ that satisfy the conditional unconfoundedness assumption in Eq. (3) are available, the exposure probability $e_{ij}$ can be unbiasedly estimated from $\mathbf{c}$ via

$$e_{ij} = p(a_{ij} = 1|\mathbf{c}) = \mathbb{E}[a_{ij}|\mathbf{c}], \tag{6}$$

which is formally known as **propensity score** in causal inference literature [58].

> **Background: The Balancing Property of Propensity Scores.**

Propensity scores have the following property called balancing [33, 59], which is the key to proving the unbiasedness of IPW-based RSs:

$$
\begin{aligned}
\mathbb{E}\left[\frac{r_{ij}}{e_{ij}}\bigg|a_{ij} = 1\right] &= \mathbb{E}\left[\frac{r_{ij} \cdot a_{ij}}{e_{ij}}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{r_{ij} \cdot a_{ij}}{e_{ij}}\bigg|\mathbf{c}\right]\right] \\
&= \mathbb{E}\left[\mathbb{E}\left[\frac{r_{ij}(a_{ij} = 1) \cdot a_{ij}}{e_{ij}}\bigg|\mathbf{c}\right]\right] \stackrel{(a)}{=} \mathbb{E}\left[\frac{\mathbb{E}[r_{ij}(a_{ij} = 1) \mid \mathbf{c}] \cdot \mathbb{E}[a_{ij} \mid \mathbf{c}]}{e_{ij}}\right] \\
&= \mathbb{E}\left[\frac{\mathbb{E}[r_{ij}(a_{ij} = 1) \mid \mathbf{c}] \cdot e_{ij}}{e_{ij}}\right] = \mathbb{E}[r_{ij}(a_{ij} = 1)],
\end{aligned} \tag{7}
$$

where the step $(a)$ follows the conditional unconfoundedness assumption in Eq. (3).

We first discuss the implementation of IPW-based RS and its unbiasedness if user and item covariates **c** that satisfy Eq. (3) are available and the propensity scores $e_{ij}$ can be calculated exactly as Eq. (6). We denote the rating predictor of an RS that aims to predict the rating potential outcome $r_{ij}(a_{ij} = 1)$ as $\hat{r}_{ij}$ and assume $r_{ij}(a_{ij} = 1)$ follows the unit-variance Gaussian distribution. Ideally, we would like $\hat{r}_{ij}$ to maximize the log-likelihood on the rating potential outcomes $r_{ij}(a_{ij} = 1)$ for all user-item pairs in $\mathcal{PO}$, which is equivalent to the minimization of the mean squared error (MSE) loss between $\hat{r}_{ij}$ and $r_{ij}(a_{ij} = 1)$ as follows:

$$\mathcal{L}^{\text{True}} = \frac{1}{I \times J} \sum_{i,j} (\hat{r}_{ij} - r_{ij}(a_{ij} = 1))^2. \tag{8}$$

However, since $r_{ij}(a_{ij} = 1)$ is unobservable for user-item pairs in the non-treatment group $\mathcal{NT}$, $\mathcal{L}^{\text{True}}$ is impossible to calculate. Therefore, traditional RSs only maximize the log-likelihood of the observed ratings for user-item pairs in the treatment group $\mathcal{T}$, which leads to the empirical MSE loss as follows:

$$\mathcal{L}^{\text{Obs}} = \frac{1}{|(i, j) : a_{ij} = 1|} \sum_{(i,j):a_{ij}=1} (\hat{r}_{ij} - r_{ij})^2, \tag{9}$$

where $|(i, j) : a_{ij} = 1|$ is the number of observed ratings. When exposure bias exists, item exposure $a_{ij}$ depends on the rating potential outcome $r_{ij}(a_{ij} = 1)$. Therefore, $\mathcal{L}^{\text{Obs}}$ is a biased estimator for $\mathcal{L}^{\text{True}}$, because the observed ratings for user-item pairs in the treatment group $\mathcal{T}$ are biased samples from the rating potential outcomes of the population $\mathcal{PO}$ (see Fig. 5-(a) and Fig. 5-(b) for an example). To remedy the bias, IPW-based causal RSs reweight the observed ratings $r_{ij}$ in $\mathcal{T}$ by the inverse of the propensity scores, i.e., $\frac{1}{e_{ij}}$, which leads to the following new training objective:

$$\mathcal{L}^{\text{IPW}} = \frac{1}{I \times J} \sum_{(i,j):a_{ij}=1} \frac{1}{e_{ij}} \cdot (\hat{r}_{ij} - r_{ij})^2. \tag{10}$$

The proof for the unbiasedness of $\mathcal{L}^{\text{IPW}}$ for $\mathcal{L}^{\text{True}}$ can be achieved by utilizing the balancing property of propensity scores in Eq. (7), where we substitute $(\hat{r}_{ij} - r_{ij})^2$ for $r_{ij}$ in the LHS of Eq. (7) and treat the rating predictor $\hat{r}_{ij}$ as constant [39]. We also provide a toy example in Fig. 5 to intuitively show the calculation of $e_{ij}$, the biasedness of $\mathcal{L}^{\text{Obs}}$ and the unbiasedness of $\mathcal{L}^{\text{IPW}}$. The objective for IPW-based RSs defined in Eq. (10) is model-agnostic. Therefore, it is applicable to all traditional RSs we introduced in Section 2. For example, for MF-based RSs, we can plug in $\hat{r}_{ij}^{\text{MF}} = \mathbf{u}_i^T \cdot \mathbf{v}_j$, for DMF-based RSs, we plug in $\hat{r}_{ij}^{\text{DMF}} = f_{nn}^u(\mathbf{u}_i)^T \cdot f_{nn}^v(\mathbf{v}_j)$, etc.

In practice, since the conditional unconfoundedness assumption in Eq. (3) is untestable, it is usually infeasible to calculate the exact value of $e_{ij}$ based on user/item covariates that satisfy Eq. (3). Nevertheless, we can still calculate approximate propensity scores $\tilde{e}_{ij}$ and reweight the observed ratings by $1/\tilde{e}_{ij}$, but the unbiasedness of Eq. (10) after the reweighting cannot be guaranteed. Here we introduce two strategies for the approximate estimation. If user/item features $\mathbf{f}_i^u$ and $\mathbf{f}_j^v$ are available,

|  | Horror Lover | Romance Lover |  |  | Horror Lover | Romance Lover |  |  | Horror Lover | Romance Lover |  |  | Horror Lover | Romance Lover |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

(a) Observed Ratings   (b) Rating Potential Outcomes   (c) Propensity Scores   (d) Predicted Ratings
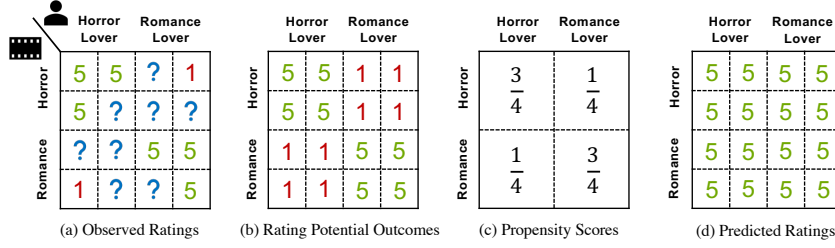
Fig. 5: An example adapted from Fig. (2) where the positivity assumption holds. Suppose again covariates $C$ represent the two-dimensional features (user type, movie type). (a) shows the observed ratings; (b) shows rating potential outcomes; (d) shows the predicted rating potential outcome of an RS model. The propensity scores $e_{ij} = p(a_{ij}|\mathbf{c}) = \mathbb{E}[a_{ij}|\mathbf{c}]$ are shown in (c). Based on (a)(d) and Eq. (9), $\mathcal{L}^{\mathrm{Obs}} = (5-1)^2 \times 2/8 = 4$. Based on (b)(d) and Eq. (8), $L^{\mathrm{True}} = (5-1)^2 \times 8/16 = 8$. Based on (a)(c)(d) and Eq. (10), $\mathcal{L}^{\mathrm{IPW}} = \frac{1}{1/4}(5-1)^2 \times 2/16 = 8$, which is unbiased for $L^{\mathrm{True}}$.

$\tilde{e}_{ij}$ can be estimated with *logistic regression* [39] as follows:

$$\tilde{e}_{ij}^{\mathrm{LR}} = \mathrm{Sigmoid}\left(\left(\sum_k w_k^u f_{ik}^u\right) + \left(\sum_k w_k^v f_{jk}^v\right) + b_i + b_j\right), \tag{11}$$

where $\mathrm{Sigmoid}(x) = (1+\exp(-x))^{-1}$, $w_k^u$ and $w_k^v$ are the regression coefficients, and $b_i$, $b_j$ are the user and item-specific offsets, respectively. If user/item features $\mathbf{f}_i^u$ and $\mathbf{f}_j^v$ are not available, we can crudely approximate $e_{ij}$ based on the exposure data alone. For example, we can estimate $\tilde{e}_{ij}$ with *Poisson factorization* [60] as:

$$\tilde{e}_{ij}^{\mathrm{PF}} \approx 1 - \exp\left(-\boldsymbol{\pi}_i^T \cdot \boldsymbol{\gamma}_j\right), \tag{12}$$

where $\boldsymbol{\pi}_i$ and $\boldsymbol{\gamma}_j$ are trainable user and item embeddings with Gamma prior, and they can be inferred from the exposure data as discussed in [61]. Additional strategies to calculate the propensity scores can be found in [62, 63, 64, 65].

The advantage of IPW is that the unbiasedness of Eq. (10) for rating potential outcome estimation can be guaranteed if the propensity scores $e_{ij}$ are correctly estimated. However, the accuracy of the propensity score estimation models relies heavily on the domain knowledge and expertise of human experts, which is untestable by experiments. In addition, IPW suffers from a large variance and numerical instability issues, especially when the estimated propensity scores $e_{ij}$ are very small. Therefore, variance reduction techniques such as clipping and multi-task learning are usually applied to improve the stability of the training dynamics [66, 67, 68].

**Substitute Confounder Adjustment.** IPW-based RSs address exposure bias from the data's perspective: They reweight the biased observational dataset to create a pseudo randomized dataset that allows unbiased training of RSs. Confounder adjustment-based methods, in contrast, estimate confounders $C$ that cause the exposure bias and adjust their effects in the rating prediction model (A simple adjustment

strategy is to control $C$ as extra covariates[12]). For the adjustment to be unbiased, classical causal inference requires the conditional unconfoundedess assumption in Eq. (3) hold, i.e., no unobserved confounders [33], which is generally infeasible in practice. Fortunately, recent advances in multi-cause causal inference [69] have shown that we can control substitute confounders estimated from item co-exposure data instead, where exposure bias can be mitigated with weaker assumptions.

We use $\mathbf{a}_i = [a_{i1}, \cdots, a_{iJ}]$ to denote the exposure status of all $J$ items to user $i$, which can be viewed as a bundle treatment in clinical trials [70]. Wang et al. [42] showed that if we can estimate user-specific latent variables $\boldsymbol{\pi}_i$, such that conditional on $\boldsymbol{\pi}_i$, the exposures of different items to the user are mutually independent, controlling $\boldsymbol{\pi}_i$ can eliminate the influence of multi-cause confounders $\mathbf{c}_i^m$ (i.e., confounders that simultaneously affect the exposure of multiple items and ratings). A simple proof of the claim is that, if $\mathbf{c}_i^m$ can still influence $\mathbf{a}_i$ and $\mathbf{r}_i$ after conditioning on $\boldsymbol{\pi}_i$, since $\mathbf{c}_i^m$ is an unobserved common cause for the exposure of different items, $a_{ij}$ cannot be conditionally independent (see the discussion of the fork structure in section 3.2.2), which renders a contradiction. The rigorous proof can be found in [69]. Wang et al. further assumed that $p(\mathbf{a}_i|\boldsymbol{\pi}_i) = \Pi_j p(a_{ij}|\boldsymbol{\pi}_i) = \Pi_j \text{Poission}(\boldsymbol{\pi}_i^T \cdot \boldsymbol{\gamma}_j)$ and used the Poisson factorization to infer $\boldsymbol{\pi}_i$ and $\boldsymbol{\gamma}_j$. Afterward, exposure bias can be mitigated by controlling $\boldsymbol{\pi}_i$ as extra covariates in the RS model [33]. For example, controlling $\boldsymbol{\pi}_i$ in MF-based RSs leads to the following adjustment:

$$r_{ij}^{\text{adj}}(a_{ij} = 1) \sim \mathcal{N}\left( \underbrace{\mathbf{u}_i^T \cdot \mathbf{v}_j}_{\text{user interests}} + \underbrace{\sum_k w_k \pi_{ik}}_{\text{adj. for expo. bias}}, \sigma_{ij}^2 \right). \tag{13}$$

The property of propensity scores can be utilized to further simplify Eq. (13): If unconfoundedness in Eq. (3) holds for $C = \boldsymbol{\pi}_i$, it will also hold for $C = \tilde{e}_{ij} = p(a_{ij}|\boldsymbol{\pi}_i)$ [58]. Therefore, we can control the approximate propensity scores estimated by $\boldsymbol{\pi}_i$, i.e., $\tilde{e}_{ij} = \boldsymbol{\pi}_i^T \cdot \boldsymbol{\gamma}_j$, which leads to the simplified adjustment formula:

$$r_{ij}^{\text{adj}}(a_{ij} = 1) \sim \mathcal{N}\left( \mathbf{u}_i^T \cdot \mathbf{v}_j + w_i \cdot \tilde{e}_{ij}, \sigma_{ij}^2 \right), \tag{14}$$

where $w_i$ is a user-specific coefficient that captures the influence of $\tilde{e}_{ij}$ on ratings.

Despite the success in addressing exposure bias with weaker assumptions, one limitation of the above method is that, since Poisson factorization is a shallow model, it may fail to capture the complex influences of multi-cause confounders on item co-exposures. To address this problem, recent works have introduced deep neural networks (DNNs) to infer the user-specific substitute confounders $\boldsymbol{\pi}_i$ from bundle treatment $\mathbf{a}_i$ [71, 72]. These methods generally assume that $\mathbf{a}_i$ are generated from $\boldsymbol{\pi}_i$ via $p(\mathbf{a}_i|\boldsymbol{\pi}_i)$ parameterized by a deep generative network $f_{nn}^{\text{exp}}$ as:

$$p(\mathbf{a}_i|\boldsymbol{\pi}_i) = \Pi_j \text{Bernoulli}(\text{Sigmoid}(f_{nn}^{\text{exp}}(\boldsymbol{\pi}_i)_j)), \tag{15}$$

---

[12] Consider again the toy example in Fig. 5. If we know exactly the user type and item type $\mathbf{c}$ for each user-item pair, the predictions can be unbiased even if the item exposures are non-randomized.
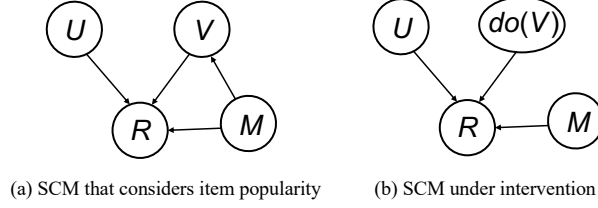
(a) SCM that considers item popularity      (b) SCM under intervention

Fig. 6: (a): SCM that explicitly models item popularity. (b): SCM under intervention $do(V)$.

where the intractable posterior of $\pi_i$ is then approximated with a Gaussian distribution parameterized by DNNs via the variational auto-encoding Bayes algorithm [73], i.e., $q(\pi_i|\mathbf{a}_i) = \mathcal{N}(f_{nn}^{\mu}(\mathbf{a}_i), \mathrm{diag}(f_{nn}^{\sigma^2}(\mathbf{a}_i)))$, where $f_{nn}^{\mu}$ and $f_{nn}^{\sigma^2}$ are two DNNs that calculate the posterior mean and variance (before diagonalization) of $\pi_i$. With deep generative models introduced to estimate the substitute confounders $\pi_i$, non-linear influences of multi-cause confounders on item exposures can be adjusted in the RS models, where exposure bias can be further mitigated in recommendations.

The key advantage of substitute confounder estimation-based causal RSs is that controlling confounders in the potential outcome prediction model generally leads to lower variance than IPW-based methods [42]. However, these models need to estimate substitute confounders $\pi_i$ from the item co-exposures and introduce extra parameters in the RS models to adjust their influences, which may incur extra bias if the confounders and the parameters are not correctly estimated. In addition, exposure bias due to single-cause confounders cannot be addressed by these methods.

### 4.1.2 Popularity Bias

Popularity bias can be viewed as a special kind of exposure bias where *users are overly exposed to popular items* [74, 75]. Therefore, it can be addressed with techniques introduced in the previous section, especially the IPW-based methods [76]. The reason is that, if we define the popularity of an item as its exposure rate:

$$m_j = \frac{\sum_i a_{ij}}{\sum_j \sum_i a_{ij}}, \tag{16}$$

we can view $m_j$ as pseudo propensity scores and use IPW to reweight the observed ratings. Alternatively, we can also analyze and address popularity bias with the structural causal model (SCM), where the causal mechanism that generates the observed ratings under the influence of item popularity is deeply investigated.

The discussion is mainly based on the popularity-bias deconfounding (PD) algorithm proposed in [48]. PD assumes that the relations among user interests $\mathbf{u}_i$, item latent attributes $\mathbf{v}_j$, item popularity $m_j$, and observed ratings $r_{ij}$ can be represented by the causal graph illustrated in Fig. 6, where item popularity can be clearly identified as a confounder that spuriously correlates the item attributes and the user

ratings. PD aims to eliminate such spurious correlations with backdoor adjustment, such that the causal influences of $\mathbf{u}_i$ and $\mathbf{v}_j$ on $r_{ij}$ (which represents users' interests on intrinsic item properties) can be properly identified. Recall that backdoor adjustment with SCM involves two stages: (1) During the training phase, the relevant structural equations in the causal graph are estimated from the collected dataset. (2) Afterward, we adjust the influence of confounders according to Eq. (5) to remove the spurious correlations. Therefore, we need to estimate $p_G(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, m_j)$ with the observed ratings $r_{ij}$ and item popularty $m_j$ and infer the latent variables $\mathbf{u}_i$ and $\mathbf{v}_j$. In PD, $p_G(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, m_j)$ is modeled as a variant of MF as follows:

$$p_G(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, m_j) \propto \underbrace{\mathrm{Elu}(\mathbf{u}_i^T \cdot \mathbf{v}_j)}_{\text{user interests}} \times \underbrace{m_j^\lambda}_{\text{pop. bias}}, \qquad (17)$$

where $\lambda$ is a hyper-parameter that denotes our belief toward the strength of influence of item popularity on ratings, and the function Elu (defined as $\mathrm{Elu}(x) = e(x)$ if $x < 0$ else $x + 1$) makes the RHS of Eq. (17) a proper unnormalized probability density function. After $\mathbf{u}_i$, $\mathbf{v}_j$ are estimated from the datasets with Eq. (17), we conduct an intervention on the item node $V$ in the causal graph (see Eq. (5)), where the spurious correlation due to item popularity can be eliminated with backdoor adjustment:

$$p(r_{ij}|do(\mathbf{u}_i, \mathbf{v}_j)) \propto \mathbb{E}_{p(m_j)}[\mathrm{Elu}(\mathbf{u}_i^T \cdot \mathbf{v}_j) \times m_j^\lambda] = \mathrm{Elu}(\mathbf{u}_i^T \cdot \mathbf{v}_j) \times \mathbb{E}_{p(m_j)}[m_j^\lambda]. \quad (18)$$

Since the second term $\mathbb{E}_{p(m_j)}[m_j^\lambda]$ in Eq. (18) is a constant and Elu is a monotonically increasing function, they have no influence on the ranking of the uninteracted items in the prediction phase. Therefore, we can drop them and use $\hat{r}_{ij} = \mathbf{u}_i^T \cdot \mathbf{v}_j$ as the unbiased rating predictor to generate future recommendations.

Generally, the debiasing mechanism of PD is very intuitive and universal among backdoor adjustment-based causal RSs [25, 24]: When fitting the RS model on the biased training set, we explicitly introduce the item popularity $m_j$ (i.e., the confounder) in Eq. (17) to explain away the spurious correlation between item attributes and the observed user ratings. Therefore, the user/item latent variables $\mathbf{u}_i$ and $\mathbf{v}_j$ used to generate future recommendations, i.e., $\hat{r}_{ij} = \mathbf{u}_i^T \cdot \mathbf{v}_j$, can focus exclusively on estimating users' true interests on intrinsic item properties.

**Is popularity bias always bad?** Recently, more researchers have begun to believe that popularity bias is not necessarily bad for RSs, because some items are popular because they *per se* have better quality than other items or they catch the current trends of user interests, where more recommendations for these items can be well-justified [77, 78]. Therefore, rather than setting the interventional distribution of item popularity to $p(m_j)$, PD introduced above as well as some other methods [48] further propose to make it correspond to item qualities or reflect the future popularity predictions. We will introduce these strategies in Section 4.3 regarding causal generalizations of RSs.
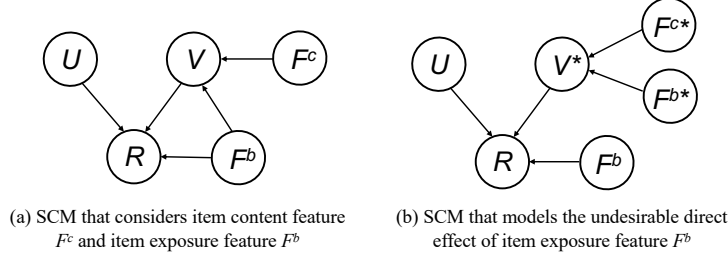
(a) SCM that considers item content feature
$F^c$ and item exposure feature $F^b$

(b) SCM that models the undesirable direct
effect of item exposure feature $F^b$

Fig. 7: (a): The SCM that considers both the causal influences of item content feature $F^c$ and item exposure feature $F^b$ on item latent variable $V$. (b): The counterfactual SCM where $V^*$ is determined by baseline value $F^{b*}$ and $F^{c*}$ to model the undesirable direct effects of $F^b$.

### 4.1.3 Clickbait Bias

Different from previous subsections that mainly focus on causal debiasing strategies for collaborative filtering-based RSs, this section discusses content-based recommendations. Specifically, we discuss the clickbait bias, which is defined as the bias of *overly recommending items with attractive exposure features such as sensational titles but with low content qualities*. The discussion is mainly based on [27]. We assume that item features $\mathbf{f}_j^v$ can be further decomposed into the item content feature $\mathbf{f}_j^c$ that captures item content information and the item exposure feature $\mathbf{f}_j^b$ whose main purpose is to attract users' attention. Taking micro-video as an example, item content feature $\mathbf{f}_j^c$ can be the audiovisual content of the video, whereas item exposure feature $\mathbf{f}_j^b$ can be its title, which is not obliged to describe its content faithfully.

The relations among user interests $\mathbf{u}_i$, item exposure feature $\mathbf{f}_j^b$, item content feature $\mathbf{f}_j^c$, item fused features $\mathbf{v}_j$, and the observed ratings $r_{ij}$ are depicted in the causal graph in Fig. 7-(a). We note that clickbait bias occurs when a user's recorded click on an item because she was cheated by the item exposure feature $\mathbf{f}_j^b$ before viewing the item content $\mathbf{f}_j^c$. Therefore, the bias can be defined as the **direct influence** of $\mathbf{f}_j^b$ on ratings $r_{ij}$ represented by the causal path $F^b \rightarrow R$. To eliminate the clickbait bias, we need to block the direct influence of $F^b$ on rating predictions, such that the item content quality can be comprehensively considered in recommendations.

As with SCM-based causal RSs, we first estimate structural equations of interest in the causal graph, i.e., $p_G(\mathbf{v}_j|\mathbf{f}_j^b, \mathbf{f}_j^c)$ and $p_G(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b)$. Since distributions in [27] are reasoned in a deterministic manner (i.e., Gaussian distributions with infinite precision), we keep the discussion consistent with them. Specifically, we use $\mathbf{v}_j(\mathbf{f}_j^b, \mathbf{f}_j^c) = f^{ff}(\mathbf{f}_j^b, \mathbf{f}_j^c)$ to represent the structural equation $p_G(\mathbf{v}_j|\mathbf{f}_j^b, \mathbf{f}_j^c)$, where $f^{ff}$ is the feature fusion function that aggregates $\mathbf{f}_j^b, \mathbf{f}_j^c$ into $\mathbf{v}_j$, and use $r_{ij}(\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b)$ to represent the structural equation $p_G(r_{ij}|\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b)$, respectively. To explicitly disentangle the influence of item exposure feature $\mathbf{f}_j^b$ and item latent variable $\mathbf{v}_j$ on the observed ratings, $r_{ij}(\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b)$ is assumed to factorize as follows:

$$r_{ij}(\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b) = \underbrace{f_{nn}^{uv}(\mathbf{u}_i, \mathbf{v}_j)}_{\text{user interests}} \cdot \underbrace{\text{Sigmoid}\left(f_{nn}^{uf}(\mathbf{u}_i, \mathbf{f}_j^b)\right)}_{\text{potential clickbait bias}}, \tag{19}$$

where the Sigmoid function provides necessary non-linearity in the fusion process. Essentially, Eq. (19) represents the causal mechanism that generates the observed ratings, which entangles both user interests in item content and clickbait bias.

However, after learning the latent variables $\mathbf{u}_i, \mathbf{v}_j$ and functions $f_{nn}^{uf}, f_{nn}^{uv}$ via Eq. (19), removing clickbait bias from the rating predictions is not as straightforward as the PD algorithm, because we should eliminate only the direct influence of item exposure feature $\mathbf{f}_j^b$ on ratings $r_{ij}$, while preserving its indirect influence mediated by item latent variable $\mathbf{v}_j$, such that all available item features can be comprehensively considered in recommendations. To achieve this purpose, we first calculate the natural direct effect (NDE) [79] of item exposure feature $\mathbf{f}_j^b$ on ratings $r_{ij}$ as follows:

$$\text{NDE}(\mathbf{u}_i, \mathbf{v}_j^*, \mathbf{f}_j^b) = r_{ij}(\mathbf{u}_i, \mathbf{v}_j^*, \mathbf{f}_j^b) - r_{ij}(\mathbf{u}_i, \mathbf{v}_j^*, \mathbf{f}_j^{b*}), \tag{20}$$

where $\mathbf{v}_j^* = f_{nn}^{ff}(\mathbf{f}_j^{b*}, \mathbf{f}_j^{c*})$, and the baseline values $\mathbf{f}_j^{b*}, \mathbf{f}_j^{c*}$ are treated as if the corresponding features are missing from the item [27]. Since the second term $r_{ij}(\mathbf{u}_i, \mathbf{v}_j^*, \mathbf{f}_j^{b*})$ in Eq. (20) denotes the user's rating to a "void" item and can be viewed as a constant, it will not affect the rank of the items. So we only adjust the first term of Eq. (20), which reasons with user $i$'s rating to item $j$ in a counterfactual world where item $j$ has only the exposure feature $\mathbf{f}_j^b$ but no content and fused features $\mathbf{f}_j^{c*}$ and $\mathbf{v}_j^*$, in Eq. (19) (Fig. 7-(b)). The adjustment leads to the following estimator,

$$\hat{r}_{ij} = r_{ij}(\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b) - r_{ij}(\mathbf{u}_i, \mathbf{v}_j^*, \mathbf{f}_j^b) \triangleq \underbrace{r_{ij}(\mathbf{u}_i, \mathbf{v}_j, \mathbf{f}_j^b)}_{\text{user interests + clickbait}} - \underbrace{r_{ij}(\mathbf{u}_i, \mathbf{v}_j^*, \mathbf{f}_j^b)}_{\text{clickbait bias}}. \tag{21}$$

Eq. (21) removes the direct influence of $\mathbf{f}_j^b$ on rating predictions, such that item content quality can be comprehensively considered in future recommendations.

### 4.1.4 Unfairness

Recently, with the growing concern of algorithmic fairness, RSs are expected to show no discrimination against users from certain demographic groups [80, 81, 82]. However, traditional RSs may capture the undesirable associations between users' sensitive information and their historical activities, which leads to potentially offensive recommendations to the users. Causal inference can help identify and address such unfair associations, where fairness can be promoted in future recommendations. This section focuses on the user-oriented fairness discussed in [83], which is defined as *the bias where RS discriminately treats users with certain sensitive attributes*.

When considering the user-oriented fairness for RSs, a subset of user features $\mathbf{f}_i$, which we denote as $\mathbf{s}_i$, is assumed to contain the sensitive information of users, such

(a) Causal Generation Process of   (b) Causal Decision Process of
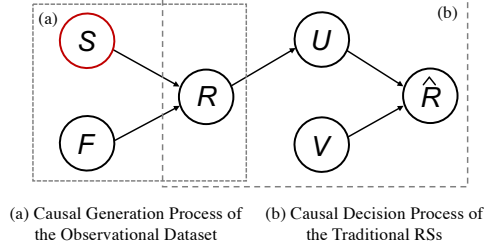    the Observational Dataset        the Traditional RSs

Fig. 8: The SCM that reasons with the causal decision mechanism of traditional RSs. Observed user ratings $R$ can be causally driven by user features $F$, including sensitive features $S$, which can then unfairly influence the inference of user latent variables $U$ and new rating predictions $\hat{R}$.

as gender, race, and age. Features $\mathbf{s}_i$ are sensitive because recommendations that improperly rely on these features may be offensive to users, which degrade both their online experiences and their trust in the system. The causal graph that depicts the causal decision mechanism of most traditional RSs is illustrated in Fig. 8 [83]. From Fig. 8 we can find that the user historical behaviors, i.e., the observed ratings $r_{ij}$, are causally driven by user features $\mathbf{f}_i$, including user sensitive features $\mathbf{s}_i$. Therefore, the user latent variables $\mathbf{u}_i$ inferred from $r_{ij}$ could capture sensitive user information in $\mathbf{s}_i$, which unfairly influences the rating predictions $\hat{r}_{ij}$ in the future.

To address this problem, Li et al. [83] proposed to disentangle the user sensitive features $\mathbf{s}_i$ from the user latent variable $\mathbf{u}_i$, such that the unfair influence of $\mathbf{s}_i$ on $\mathbf{u}_i$ represented by the causal chain $S \rightarrow R \rightarrow U$ can be maximally suppressed in the future recommendations. A common strategy to achieve the disentanglement is adversarial training [84], where we train a discriminator $f_{nn}^{\mathrm{cls}}(\mathbf{u}_i) \rightarrow \mathbf{s}_i$ that predicts the sensitive features $\mathbf{s}_i$ from user latent variables $\mathbf{u}_i$ alongside the RS. While fitting the RS on the observe ratings $r_{ij}$, we constrain the inferred $\mathbf{u}_i$ to fool the discriminator $f_{nn}^{\mathrm{cls}}$ by making wrong predictions about $\mathbf{s}_i$, which discourages $\mathbf{u}_i$ from capturing sensitive information in $r_{ij}$ due to its unfair correlations with $\mathbf{s}_i$. Here we take the MF-based RS as an example to show the details. We use $\mathcal{L}^{\mathrm{Rec}}$ to denote the original training objective of the MF-based RS that maximizes the log-likelihood on observed ratings $r_{ij}$ and use $\mathcal{L}^{\mathrm{cls}}$ to denote the loss function of the discriminator $f_{nn}^{\mathrm{cls}}$. The adjusted training objective $\mathcal{L}^{\mathrm{Fair}}$ with fairness constraint becomes the following:

$$\mathcal{L}^{\mathrm{Fair}} = \underbrace{\mathcal{L}^{\mathrm{Rec}}(\mathbf{u}_i^T \cdot \mathbf{v}_j, r_{ij})}_{\text{user interests}} - \lambda \cdot \underbrace{\mathcal{L}^{\mathrm{cls}}(f_{nn}^{\mathrm{cls}}(\mathbf{u}_i), \mathbf{s}_i)}_{\text{fairness constraint}}, \qquad (22)$$

where $\lambda$ is a hyper-parameter that balances the recommendation performance and the fairness objective. Generally, a higher $\lambda$ leads to better fairness, but it also restricts the capacity of the user latent variables $\mathbf{u}_i$, which could negatively impact the recommendation performance. Although here we use the MF-based RS as an example, it is straightforward to generalize Eq. (22) to DMF or AE-based RS by replacing the $\mathbf{u}_i^T \cdot \mathbf{v}_j$ term with the corresponding rating estimators.
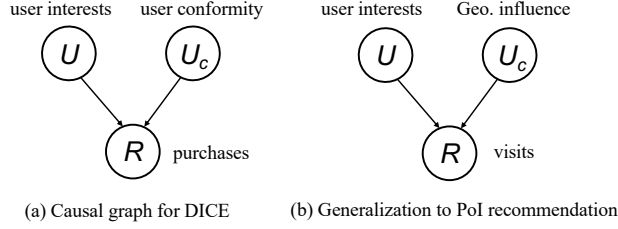
(a) Causal graph for DICE      (b) Generalization to PoI recommendation

Fig. 9: Causal Graphs for DICE (a) and its generalization to PoI recommendations (b).

## 4.2 Causal Explanation in Recommendations

In previous sections, we have introduced causality to address various types of bias and spurious correlation issues for traditional RSs. In this section, we use causality to explain the user decision process. Specifically, we discuss an interesting question aiming to disentangle users' intent that causally explains their past behaviors, i.e., *did a user purchase an item because she conformed to the current trend or because she really liked it*? The tricky part of this question is that: in reality, we only observe the effects, i.e., the purchases, which can be explained by both causes.

### 4.2.1 Disentangling Interest and Conformity with Causal Embedding

The discussion is based on DICE proposed in [56]. To simplify the discussion, we consider $r_{ij}$ as implicit feedback and define the set of user, positive item ($j : r_{ij} = 1$), negative item ($k : r_{ik} = 0$) triplets as $\mathcal{R}_{pn} = \{(i, j, k) | r_{ij} = 1 \wedge r_{ik} = 0\}$. The popularity of each item $j$, i.e., $m_j$, which reflects the current trend, can be calculated with Eq. (16). Observing that the causal relation between user interests $U$, user conformity $U_c$ and observed ratings $R$ can be represented as a V-structure in Fig. 9-(a), DICE exploits the *colliding effect* to achieve the disentanglement, i.e., outcomes that cannot be explained by one cause are more likely caused by another (see discussions in Section 3.2.2). Therefore, although users' interests cannot be directly estimated from their ratings $r_{ij}$ due to entanglement, their conformity to the trend can be estimated by the popularity level of item $j$, and positive feedback not likely caused by conformity has a higher chance of reflecting users' true interests.

In implementation, DICE assumes that the observed ratings $r_{ij}$ can be decomposed into the sum of a conformity part $r_{ij}^c = f^c(\mathbf{u}_i^c, \mathbf{v}_j^c)$ and a user interests part $r_{ij}^i = f^i(\mathbf{u}_i^i, \mathbf{v}_j^i)$, where $\mathbf{u}_i^{c,i}, \mathbf{v}_j^{c,i}$ are learnable user, item embeddings that reflect user $i$'s interests in (i.e., superscript $i$) and conformity to (i.e., superscript $c$) item $j$, respectively. According to the colliding effect of causal graphs, we can split the triplets in $\mathcal{R}_{pn}$ into two parts: In the first part $\mathcal{R}_{pn}^{(1)}$, positive item $a$ in the triplet has a higher popularity level than the negative item $b$, i.e., $m_a > m_b$. In this case, we can draw two general conclusions from this triplet: (1) Overall, the user prefers item $a$ over $b$; (2) She is more likely to conform to item $a$ than item $b$ due to $a$'s higher

popularity. These conclusions lead to the two inequalities as follows:

$$\forall (i, a, b) \in \mathcal{R}_{pn}^{(1)}, \text{ we have } \begin{cases} r_{ia}^c > r_{ib}^c \text{ (conformity)} \\ r_{ia}^i + r_{ia}^c > r_{ib}^i + r_{ib}^c \text{ (overall preference)}, \end{cases} \tag{23}$$

where the dependency of $r_{i\{a,b\}}^{c,i}$ on latent variables $\mathbf{u}_i^{c,i}$, $\mathbf{v}_{\{a,b\}}^{c,i}$ are omitted for simplicity. The second part, i.e., $\mathcal{R}_{pn}^{(2)}$, is the **key** to achieving disentanglement, because for every triplet $(i, c, d)$ in $\mathcal{R}_{pn}^{(2)}$, the negative item $d$ is more popular than the positive item $c$. In this case, user $i$ *could have simply conformed to the trend* and chosen item $d$ to consume, but instead, she actively chose the less popular item $c$. Therefore, we can draw one more specific conclusion that leads to the disentanglement between user interests and conformity: The choice of item $c$ over $d$ is more likely due to user interests. Therefore, we can form three inequalities as:

$$\forall (i, c, d) \in \mathcal{R}_{pn}^{(2)}, \text{ we have } \begin{cases} r_{ic}^i > r_{id}^i \text{ (interests)}, r_{ic}^c < r_{id}^c \text{ (conformity)}, \\ r_{ic}^i + r_{ic}^c > r_{id}^i + r_{id}^c \text{ (overall preference)}. \end{cases} \tag{24}$$

The inequalities in Eqs. (23) and (24) can be solved by ranking-based loss in RSs, such as Bayesian personalized ranking (BPR) [85], where the disentangled embeddings $\mathbf{u}_i^{c,i}$, $\mathbf{v}_j^{c,i}$ and the match functions $f^{c,i}(\cdot, \cdot)$ can be learned from $\mathcal{R}_{pn}^{(1)}$ and $\mathcal{R}_{pn}^{(2)}$. Finally, we form a rating predictor $\hat{r}_{ij} = f^i(\mathbf{u}_i^i, \mathbf{v}_j^i) + f^c(\mathbf{u}_i^c, \mathbf{v}_j^c)$ for future recommendations.

### 4.2.2 Generalizations of DICE

DICE disentangles the user intent and promotes the explainability of RSs from the data's perspective: It partitions the triplets $(i, j, k)$ in $\mathcal{R}_{pn}$ into two disjoint subsets $\mathcal{R}_{pn}^{(1)}$ and $\mathcal{R}_{pn}^{(2)}$ based on the relative popularity of the positive and negative items, and shows that the triplets in $\mathcal{R}_{pn}^{(2)}$ are informative to distinguish the user interests from their conformity to the popularity trend. The basic idea of DICE is generalizable to promote explainability for other types of recommendation tasks, if we can find **alternative causal explanations** to challenge the assumption that the observed positive feedback in these tasks can be attributed solely to user interests.

For example, in point-of-interests (PoI) recommendations, the target items are specific point locations that users may find useful or interesting to visit, such as restaurants, grocery stores, and malls [7]. In this task, the location of a PoI is an important alternative explanation for users' visits to the PoI other than user interests, because nearby POIs are more convenient to visit than the remote ones [86]. Therefore, to disentangle user interests from potential geographical factors that could causally influence users' choices, we can take a similar strategy as DICE and partition the user historical visit records according to the distance of positive and negative PoIs to users. Then, the disentangled user interest embeddings can be estimated based on the partitioned dataset with the same ranking-based approach.

### 4.2.3 Other Works on Explainable RSs

Explanable recommendation is a broad topic [87], where disentangling user's intent based on data partitioning is a small part. There are also plenty of works that focus on improving the explainability of RSs from the model's side, where specific disentanglement modules, such as prototype learning [88], context modeling [89], and aspect modeling [90], are designed and integrated with traditional RS models to further enhance their transparency and explainability. We refer interested readers to the corresponding papers as well as [91, 92] for further investigation.

## 4.3 Causal Generalization of Recommendations

After estimating the causal relations from potentially biased and entangled observational datasets, the generalization ability of RSs can be substantially enhanced, because even if the context (or environment) in which we make recommendations changes (e.g., item popularity, user conformity, etc.), we can still basing the recommendations on causal relations that are *stable and invariant*, while discarding or correcting other undesirable correlations that are transient and susceptible to change [56, 93]. In this section, we use the PD algorithm for popularity bias and the DICE algorithm for causal explainability as two examples to show how the generalization of RSs can be improved with causal intervention and disentanglement.

### 4.3.1 Generalization Based on Intervention

First, we take the PD algorithm as an example to show how causal intervention can improve the generalization of RSs within a dynamic environment. In RS, it is generally assumed that user interests can remain unchanged for a certain period of time, i.e., the causal structure $U \to R \leftarrow V$ in Fig. (6) represents the stable user interests on intrinsic item properties. However, the popularity of different items, i.e., the context in which we make recommendations, can shift rapidly during the same period [78]. Recall that PD disentangles the causal influences of user interests and item popularity on ratings via the product of two terms, i.e., $\text{Elu}(\mathbf{u}_i^T \cdot \mathbf{v}_j)$ and $m_j^\lambda$, as Eq. (17). Suppose $m_j$ represents the current popularity level of item $j$. If we predict that the popularity of item $j$ will change to $m_j'$ in the future [6], we can conduct an intervention that sets $M$ to the predicted value $m_j'$ in the structural equation $p_G(R|U, V, M)$ and predict future ratings $r_{ij}'$ via the following formula:

$$p_G(r_{ij}'|\mathbf{u}_i, \mathbf{v}_j, do(m_j')) \propto \underbrace{\text{Elu}(\mathbf{u}_i^T \cdot \mathbf{v}_j)}_{\text{stable user interests}} \times \underbrace{(m_j')^\lambda}_{\text{future popularity}} \quad , \qquad (25)$$

where the user, item latent variables $\mathbf{u}_i$ and $\mathbf{v}_j$ learned from the current time step remain unaltered. With the influence of future changes in item popularity on ratings considered in the predictions, service providers can make strategic decisions to allocate resources for items with different popularity potentials. In contrast, traditional RSs could mistakenly capture the influence of the current popularity level of items on ratings as user interests. Therefore, they will not generalize well when the item popularity $m_j$ changes to a different level $m'_j$ due to time evolution.

### 4.3.2 Generalization Based on Disentanglement

In addition, causal disentanglement can promote the generalization of RSs by identifying and basing recommendations on causes that are more robust to potential changes in the environments [94, 95]. For example, if users' conformity and interest are disentangled based on their historical behaviors, if a user's conformity reduces to a low level due to certain reasons, since user interests are assumed to be stable within a certain period of time, we can still use the learned user/item interest variables $\mathbf{u}_i^i$, $\mathbf{v}_j^i$ to make recommendations based on their interests, where the previously estimated unreliable user conformity information can be discarded or down-weighted. In contrast, for traditional RSs, different factors that causally influence their historical behaviors are entangled as a single user latent variable $\mathbf{u}_i$. Therefore, even if some less stable causes of user behaviors are known to change (e.g., in the PoI RS introduced above, a user could move to a new place where the convenience levels of different PoIs change for the user), these models will still utilize the outdated causes to make recommendations, which could fail to generalize to the new environment.

## 5 Evaluation Strategies for Causal RSs

In the previous sections, we have discussed various causal inference techniques that are promising to address multiple types of biases, entanglement, and generalization problems in traditional RSs. However, without a well-designed model evaluation strategy, it is difficult to tell whether the proposed causal RS model is indeed effective, nor can we guarantee that the model will perform reliably after deploying in a real-world environment. The evaluation of causal models is particularly difficult, because the groundtruths, i.e., the causal effects of interest, are usually infeasible. Despite the challenges, there are several strategies that can reliably evaluate causal RSs with biased real-world data, and we will thoroughly discuss them in this section. In addition, we also compile the available real-world datasets that conduct randomized experiments to eliminate exposure bias to facilitate future causal RS research.

## 5.1 Evaluation Strategies for Traditional RSs

The assessment of traditional RSs generally follows three steps: First, the observed ratings $r_{ij}$ in the rating matrix $\mathbf{R}$ are split into the non-overlapping training set $\mathbf{R}_{tr}$ and test set $\mathbf{R}_{te}$, usually by randomly holding out a certain percentage of the observed ratings from each user. Then, the proposed RS is trained on ratings in $\mathbf{R}_{tr}$ to learn the latent variables and the associated functional models (see Section 2). Finally, the trained RS predicts the missing ratings in $\mathbf{R}_{tr}$ for each user, where the results are compared with the held-out ratings in $\mathbf{R}_{te}$ to evaluate the model performance. The quality of rating predictions can be measured by accuracy-based metrics such as mean squared error (MSE) and mean absolute error (MAE), and ranking-based metrics such as recall, precision, normalized discounted cumulative gain (NDCG), etc. More information on these evaluation metrics can be found in [96].

## 5.2 Challenges for the Evaluation of Causal RSs

The above evaluation strategy, however, is not directly applicable to causal RSs, because ratings in $\mathbf{R}_{te}$ may have the same spurious correlation and bias as ratings in $\mathbf{R}_{tr}$, which makes the evaluation on $\mathbf{R}_{te}$ a biased measure of the true model performance. Therefore, to unbiasedly evaluate the effectiveness of causal RSs, it is ideal that we have a biased/entangled training set $\mathbf{R}_{tr}^{b}$ to learn the model, and an unbiased/disentangled test set $\mathbf{R}_{te}^{ub}$ for model evaluation, such that the effectiveness of the causal debiasing/disentangling algorithm can be directly verified from experiments. However, such unbiased/disentangled test set $\mathbf{R}_{te}^{ub}$ can be difficult to acquire and expansive to establish. Therefore, we first introduce common data simulation strategies for causal RS evaluation. We then discuss how real-world datasets can be directly utilized to further promote the credibility of causal RS research.

## 5.3 Evaluation Based on Simulated Datasets

A good dataset simulation strategy to evaluate causal RSs should have the following properties: (1) The generation mechanisms of the bias and entanglement to be studied are clearly identified, credibly designed, and can be adjusted in a flexible manner; (2) The available real-world information is utilized as much as possible.

### 5.3.1 Simulation Based on Generative Models

One promising dataset simulation strategy that satisfies the above criteria is to use deep generative models. Here we take exposure bias as an example to demonstrate how it can be simulated from real-world datasets [71]. The simulation is composed of

two phases. In the training phase, two variational auto-encoders (VAEs) [22, 73] are trained on the exposure and rating data in a real-world dataset (e.g., the MovieLens dataset [5]), which results in two decoder networks $f_{nn}^a$ and $f_{nn}^r$ that generate item exposures $\mathbf{a}_i \in \{0, 1\}^J$ and user ratings $\mathbf{r}_i \in \mathbb{R}^J$ from $K$-dimensional Gaussian user latent variables $\mathbf{u}_i^a \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ and $\mathbf{u}_i^r \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$, respectively. The decoders capture the generative distributions of item exposures and user ratings based on the data of real users, where the available real-world information is effectively utilized. In the generation phase, for each hypothetical user $i'$, we draw a confounder $\mathbf{c}_{i'} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ that simultaneously affects $\mathbf{u}_{i'}^a$ and $\mathbf{u}_{i'}^r$. Then, to simulate the exposure bias, we set $\mathbf{u}_{i'}^a = \mathbf{c}_{i'}$ and $\mathbf{u}_{i'}^r = \lambda \cdot \mathbf{c}_{i'} + (1 - \lambda)\boldsymbol{\epsilon}_{i'}$ and use $f_{nn}^a$, $f_{nn}^r$ to generate the simulated item exposures $\mathbf{a}_{i'}$ and ratings $\mathbf{r}_{i'}$, where $\boldsymbol{\epsilon}_{i'} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ and hyperparameter $\lambda$ controls the strength of the confounding bias. Finally, we mask $\mathbf{r}_{i'}$ with $\mathbf{a}_{i'}$ to form the biased training set $\mathbf{R}_{tr}^b$, and keep the generated ratings $\mathbf{r}_{i'}$ unmasked in the test set $\mathbf{R}_{te}^{ub}$ for an unbiased evaluation of model performance.

The advantage of dataset simulation strategies based on generative models is that the true causal mechanisms of interest, such as the rating potential outcomes, are available in the evaluation stage, which is generally impossible for real-world datasets. Therefore, the effectiveness of causal RSs can be easily verified based on the simulated groundtruths. In addition, the simulations are flexible as the strength of biases and entanglements can be set into different levels (e.g., $\lambda$ in the example), where the sensitivity and robustness of causal RSs can be thoroughly investigated.

### 5.3.2 Test Set Intervention

Another reliable dataset simulation strategy is test set intervention, where an intervened test set is created from the original test set, such that it has a different bias/entanglement distribution from the training set [56, 60, 97]. For example, to study the popularity bias, we can first select observed ratings from $\mathbf{R}$ such that 90% of the interacted items are popular and 10% are unpopular to form the training set $\mathbf{R}_{tr}$ [98]. We then select from the remaining ratings, i.e., the original test set $\mathbf{R}_{te}$, a subset with a different ratio of popular and unpopular items (e.g., 10% popular and 90% unpopular) to form the intervened test set $\mathbf{R}_{te}^{int}$. If the causal RSs trained on $\mathbf{R}_{tr}$ can still perform well on the intervened test set $\mathbf{R}_{te}^{int}$, the model's invariance to the popularity bias can be supported. A similar test set intervention strategies can be used to evaluate the disentanglement of user interests and conformity for DICE [56].

The advantage of the test set intervention-based causal RS evaluation strategy is that extra assumptions that cannot be justified by real-world information are minimally introduced, because the intervention is usually achieved by selecting samples from the original test set to change the data distribution, which does not introduce extra assumptions of the generative mechanisms or hypothetical users, items, and ratings. From this perspective, the evaluation results based on test set intervention may be more credible compared with the generative model-based strategies.

## 5.4 Evaluation Based on Real-world Datasets

### 5.4.1 Randomized Experiments

For the study of exposure bias, it is feasible to establish-bias free real-world datasets, where ratings for either every item or randomly selected items are collected from a subset of users. This can be extremely expansive and user-unfriendly, but recent years have witnessed a growing interest in causal RS research from the industry, where more such randomized datasets are established and released to facilitate causal RS research. The available real-world datasets are compiled as follows:

- **Coat datasets**[13] [39] (2016). The Coat dataset is a small-scale dataset crowd-sourced from the Amazon Mechanical Turkers platform with 300 users and 290 items. Specifically, each Turker is first asked to self-select 24 coats to rate, where the ratings form the biased training set $\mathbf{R}_{tr}^b$. Then each Turker is asked to rate 16 random coats, and these ratings form the unbiased test set $\mathbf{R}_{te}^{ub}$.

- **Yahoo! R3 dataset**[14] [99, 100] (2009). The Yahoo! R3 dataset is collected from the Yahoo! Music platform. The biased training set $\mathbf{R}_{tr}^b$ is composed of 300,000 self-supplied ratings from 15,400 users to 1,000 items. In addition, a subset of 5,400 users is presented with ten randomly selected items to rate, and the ratings are used to create the unbiased test set $\mathbf{R}_{te}^{ub}$.

- **KuaiRec dataset**[15] [101] (2022). The KuaiRec dataset is established based on a popular micro-video sharing platform, KuaiShou, in China (known as Kwai internationally). The dataset records self-supplied ratings from 7,176 users to 10,728 items as the biased training set $\mathbf{R}_{tr}^b$. The unbiased test set $\mathbf{R}_{te}^{ub}$ is composed of a subset of 1,411 users and 3,327 items, where the ratings between these users and items are almost fully observed (with 99.6% density).

The statistics of the datasets are summarized in Table 1 for reference. There are also randomized datasets for some related topics such as click-through rate prediction [102], i.e., Criteo Ads datasets[16] [103], and bandit-based RS [104], i.e., Open Bandit dataset[17] [105], where the sources are also provided in case the readers are interested.

From Table 1 we can find that, the Coat dataset is small in scale. While for the Yahoo! R3 dataset, the training set is comparatively large (15,400 users and 1,000 items), the randomized experiment conducted to establish the unbiased test set is small-scale in comparison (16 and 10 randomly exposed items per user, respectively). Therefore, although these ratings are unbiased due to randomization, they may not capture well-rounded user interests and therefore induce a high evaluation variance.

---

[13] https://www.cs.cornell.edu/~schnabts/mnar/

[14] https://webscope.sandbox.yahoo.com/catalog.php?datatype=r&did=3

[15] https://github.com/chongminggao/KuaiRec

[16] http://cail.criteo.com/criteo-uplift-prediction-dataset/

[17] https://research.zozo.com/data.html

| Dataset | # Users | # Items | Item Type | Training Sets | Test Sets |
|---------|---------|---------|-----------|---------------|-----------|
| Coat | 300 | 290 | Coat | 24 i/u (self-supplied) | 16 i/u (random) |
| Yahoo! R3 | 15,400 | 1,000 | Music | 300,000 r (self-supplied) | 10 i/u (random) for 5,400 u |
| KuaiRec | 7,176 | 10,728 | Video | 16.3% r (self-supplied) | 99.6% r for 1,411 u and 3,327 i |

Table 1: Characteristics of the currently available real-world causal recommendation datasets, where the test sets are devoid of exposure bias either due to randomized item exposures or fully observed ratings. In the table, terms like 24 i/u mean that every user rates 24 items, the term 300,000 r denotes the number of observed ratings, and terms like 16.3% r represent the density of interactions.

For the recently released KuaiRec datasets, large-scale experiments are conducted on users to establish the bias-free test set, where the 1,411 users' ratings for 3,327 items are almost fully collected. Therefore, it may be a promising new benchmark that allows the evaluation of more complex causal RS models with a lower variance.

### 5.4.2 Qualitative Evaluation and Case Study

For other types of biases in RSs that cannot be attributed to non-randomized item exposures (e.g., clickbait bias and unfairness), the establishment of bias-free test sets is more challenging. For example, when studying the clickbait bias, it is difficult to determine whether a user clicked an item due to interests or clickbait. Similarly, when examining the user-oriented fairness of RSs, we cannot know if the generated items are offensive to the users. Under such circumstances, we can still conduct case studies for qualitative model evaluations, where we manually select some representative samples from the original test set and observe whether the trained causal RS model would respond as expected to these samples [106].

Consider the evaluation of the robustness of a causal RS to clickbait bias. We can select some representative items with low-quality content but highly-attractive exposure features and other items with high-quality content but normal exposure features from the original test set. Then, we obtain rating predictions for items from these two groups and draw comparisons. If the studied causal RS indeed ranks items in the second group higher than those in the first group, we can likely conclude that the model is robust to clickbait bias because the quality of the item content, not its exposure features, is prioritized in recommendations. In addition, to evaluate the user-oriented fairness of a causal RS, we can analyze the generated recommendation for users from certain demographic groups. If the recommended items tend to capture the social stereotypes that are negatively associated with user sensitive features, we can conclude that the model is still discriminatory against users.

# 6 Future Directions

Despite the recent achievements in marrying causal inference with traditional RSs to address their various limitations of correlational reasoning on observational user data, causal RS research is still in its emerging stage. Several promising directions could be pursued to further advance this field. In this section, we identify four interesting and important open problems worthy of exploration in the future.

First, the assumptions required by existing causal RSs could be too strong, which may not hold in reality. For example, most RCM-based causal RSs rely on SUTVA to exclude the interference of item exposures for different users. However, if users are connected by a social network, they may interact closely with each other or be heavily affected by the influencers in the network [107]. Consequently, SUTVA can be violated because the recommendations made to one user may causally affect the ratings of others (i.e., the spill-over effects [108, 109]). In addition, the positivity assumptions may also be violated if some users never click certain types of items (i.e., non-compliance and defiers [33]). Therefore, it is crucial to further weaken the assumptions of causal RSs to make them more practical for real-world applications.

In addition, there currently lacks a universal causal model for RSs that can be applied for different causal reasoning purposes. Most SCM-based causal RSs are designed to address one specific type of bias or entanglement problem, where other issues are tacitly assumed to be absent and omitted from the causal graph. Moreover, even for causal RSs that address the same problem, several varieties of causal graphs that include different sets of variables and relationships can be assumed, which leads to inconsistency between different works. Therefore, it would be promising and beneficial to have a generic and widely-accepted causal model that is able to comprehensively address multiple types of causal problems in recommendations.

Furthermore, certain types of biases in RSs are double-blade swords, where the positive side is seldom investigated. Consider the item exposure bias discussed in Section 4.1.1. We should note that some items are more likely to be exposed because they have higher quality than other items. Therefore, the higher exposure rate of these items can be well justified and may be utilized to further enhance the recommendation performance. In addition, recent research also found that confounders that spuriously correlate item exposures and user ratings may also help explain the co-occurrence patterns of different items [71]. Therefore, how to properly identify and utilize the positive side of biases while maximally suppressing their negative effects is of great importance and deserves more in-depth investigations in the future.

Finally, although recent years have witnessed the establishment and release of more real-world datasets for causal RS research from the industry, many causal RS models still rely heavily on simulated datasets for evaluation. The simulation can lead to the over-simplification of the problem and is often designed to correspond exactly with the debiasing/disentanglement mechanism of the proposed model. Therefore, the effectiveness of these methods in more complicated real-world scenarios is still uncertain due to the lack of model deployment and online tests. As such, to more convincingly demonstrate the practical utility of causal RSs, more collaborations with the industry are highly expected.

## 7 Summary

In this survey, we provide a comprehensive overview of recent advances in causal inference for RSs. We start by pointing out issues of traditional RSs that rely on correlations in observed user behaviors and user/item features. We then introduce two mainstream causal inference frameworks, i.e., Rubin's RCM and Pearl's SCM, which provide deeper insights into these issues and the foundation for moving traditional RSs to the upper rungs of the Ladder of Causality. Specifically, we thoroughly discuss several state-of-the-art causal RS models that lead to enhanced robustness to various biases and improved explainability. In addition, since causal RSs can base recommendations on causal relationships that are stable and invariant, we also demonstrate that their generalization abilities can be significantly improved. Finally, we introduce evaluation strategies for causal RSs, with an emphasis on how to reliably estimate the model performance based on biased real-world data. We further compile real-world datasets where expensive randomized experiments are conducted on users, which reflects growing attention to causal RSs from the industry.

Overall, causal RS is still a relatively new and under-explored research topic. More efforts are urgently demanded to systematize the existing works and conduct deeper investigations for further improvements. Accordingly, we point out four interesting and practically important open problems in causal RSs. We hope that this survey can help readers gain a comprehensive understanding of the main idea of applying causality in RSs and encourage further progress in this promising area.

## References

1. Jiahui Liu, Peter Dolan, and Elin Rønby Pedersen. Personalized news recommendation based on click behavior. In *Proceedings of the 15th International Conference on Intelligent User Interfaces*, pages 31–40, 2010.
2. Ioannis Paparrizos, B Barla Cambazoglu, and Aristides Gionis. Machine learned job recommendation. In *Proceedings of the 5th ACM Conference on Recommender Systems*, pages 325–328, 2011.
3. Chong Wang and David M Blei. Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 448–456, 2011.
4. Jing Yi, Yaochen Zhu, Jiayi Xie, and Zhenzhong Chen. Cross-modal variational auto-encoder for content-based micro-video background music recommendation. *IEEE Transactions on Multimedia*, 2021.
5. F Maxwell Harper and Joseph A Konstan. The MovieLens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems*, 5(4):1–19, 2015.

6. Jiayi Xie, Yaochen Zhu, Zhibin Zhang, Jian Peng, Jing Yi, Yaosi Hu, Hongyi Liu, and Zhenzhong Chen. A multimodal variational encoder-decoder framework for micro-video popularity prediction. In *Proceedings of The Web Conference 2020*, pages 2542–2548, 2020.

7. Mao Ye, Peifeng Yin, Wang-Chien Lee, and Dik-Lun Lee. Exploiting geographical influence for collaborative point-of-interest recommendation. In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 325–334, 2011.

8. Yuyun Gong and Qi Zhang. Hashtag recommendation using attention-based convolutional neural network. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, pages 2782–2788, 2016.

9. Francesco Ricci, Lior Rokach, and Bracha Shapira. Introduction to recommender systems handbook. In *Recommender Systems Handbook*, pages 1–35. Springer, 2011.

10. Yehuda Koren, Steffen Rendle, and Robert Bell. Advances in collaborative filtering. In *Recommender Systems Handbook*, pages 91–142. Springer, 2022.

11. Pasquale Lops, Marco de Gemmis, and Giovanni Semeraro. Content-based recommender systems: State of the art and trends. In *Recommender Systems Handbook*, pages 73–105. Springer, 2011.

12. Erion Çano and Maurizio Morisio. Hybrid recommender systems: A systematic literature review. *Intelligent Data Analysis*, 21(6):1487–1524, 2017.

13. Liuyi Yao, Zhixuan Chu, Sheng Li, Yaliang Li, Jing Gao, and Aidong Zhang. A survey on causal inference. *ACM Transactions on Knowledge Discovery from Data*, 15(5):1–46, 2021.

14. Chen Gao, Yu Zheng, Wenjie Wang, Fuli Feng, Xiangnan He, and Yong Li. Causal inference in recommender systems: A survey and future directions. *arXiv preprint arXiv:2208.12397*, 2022.

15. Judea Pearl and Dana Mackenzie. *The book of why: The new science of cause and effect*. Basic books, 2018.

16. Yifan Hu, Yehuda Koren, and Chris Volinsky. Collaborative filtering for implicit feedback datasets. In *The 8th IEEE International Conference on Data Mining*, pages 263–272, 2008.

17. Andriy Mnih and Russ R Salakhutdinov. Probabilistic matrix factorization. In *Advances in Neural Information Processing Systems*, 2007.

18. Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.

19. Hong-Jian Xue, Xinyu Dai, Jianbing Zhang, Shujian Huang, and Jiajun Chen. Deep matrix factorization models for recommender systems. In *International Joint Conference on Artificial Intelligence*, volume 17, pages 3203–3209, 2017.

20. Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. Deep learning based recommender system: A survey and new perspectives. *ACM Computing Surveys*, 52(1):1–38, 2019.

21. Yao Wu, Christopher DuBois, Alice X Zheng, and Martin Ester. Collaborative denoising auto-encoders for top-N recommender systems. In *Proceedings of the 9th ACM International Conference on Web Search and Data Mining*, pages 153–162, 2016.

22. Dawen Liang, Rahul G Krishnan, Matthew D Hoffman, and Tony Jebara. Variational autoencoders for collaborative filtering. In *Proceedings of the World Wide Web Conference*, pages 689–698, 2018.

23. Yaochen Zhu and Zhenzhong Chen. Mutually-regularized dual collaborative variational autoencoder for recommendation systems. In *Proceedings of The ACM Web Conference 2022*, pages 2379–2387, 2022.

24. Shuyuan Xu, Yingqiang Ge, Yunqi Li, Zuohui Fu, Xu Chen, and Yongfeng Zhang. Causal collaborative filtering. *arXiv preprint arXiv:2102.01868*, 2021.

25. Wenjie Wang, Fuli Feng, Xiangnan He, Xiang Wang, and Tat-Seng Chua. Deconfounded recommendation for alleviating bias amplification. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1717–1725, 2021.

26. Yinwei Wei, Xiang Wang, Liqiang Nie, Xiangnan He, Richang Hong, and Tat-Seng Chua. MMGCN: Multi-modal graph convolution network for personalized recommendation of micro-video. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 1437–1445, 2019.

27. Wenjie Wang, Fuli Feng, Xiangnan He, Hanwang Zhang, and Tat-Seng Chua. Clicks can be cheating: Counterfactual recommendation for mitigating clickbait issue. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1288–1297, 2021.

28. Yehuda Koren. Factorization meets the neighborhood: A multifaceted collaborative filtering model. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 426–434, 2008.

29. Yaochen Zhu and Zhenzhong Chen. Variational bandwidth auto-encoder for hybrid recommender systems. *IEEE Transactions on Knowledge and Data Engineering*, 2022.

30. Steffen Rendle. Factorization machines. In *IEEE International Conference on Data Mining*, pages 995–1000. IEEE, 2010.

31. Xiangnan He and Tat-Seng Chua. Neural factorization machines for sparse predictive analytics. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 355–364, 2017.

32. Yunqi Li, Hanxiong Chen, Juntao Tan, and Yongfeng Zhang. Causal factorization machine for robust recommendation. In *Proceedings of the 22nd ACM/IEEE Joint Conference on Digital Libraries*, pages 1–9, 2022.

33. Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.

34. Judea Pearl. *Causality*. Cambridge University Press, 2009.

35. Xiaojie Wang, Rui Zhang, Yu Sun, and Jianzhong Qi. Combating selection biases in recommender systems with a few unbiased ratings. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 427–435, 2021.

36. Dugang Liu, Pengxiang Cheng, Zhenhua Dong, Xiuqiang He, Weike Pan, and Zhong Ming. A general knowledge distillation framework for counterfactual recommendation via uniform data. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 831–840, 2020.

37. Xuanhui Wang, Nadav Golbandi, Michael Bendersky, Donald Metzler, and Marc Najork. Position bias estimation for unbiased learning to rank in personal search. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining*, pages 610–618, 2018.

38. Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. Controlling popularity bias in learning-to-rank recommendation. In *Proceedings of the 11th ACM Conference on Recommender Systems*, pages 42–46, 2017.

39. Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims. Recommendations as treatments: Debiasing learning and evaluation. In *International Conference on Machine Learning*, pages 1670–1679, 2016.

40. Peng Wu, Haoxuan Li, Yuhao Deng, Wenjie Hu, Quanyu Dai, Zhenhua Dong, Jie Sun, Rui Zhang, and Xiao-Hua Zhou. On the opportunity of causal learning in recommendation systems: Foundation, estimation, prediction and challenges. In *Proceedings of the International Joint Conference on Artificial Intelligence, Vienna, Austria*, pages 23–29, 2022.

41. Stephen Bonner and Flavian Vasile. Causal embeddings for recommendation. In *Proceedings of the 12th ACM Conference on Recommender Systems*, pages 104–112, 2018.

42. Yixin Wang, Dawen Liang, Laurent Charlin, and David M Blei. Causal inference for recommender systems. In *The 14th ACM Conference on Recommender Systems*, pages 426–431, 2020.

43. Harald Steck. Training and testing of recommender systems on data missing not at random. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 713–722, 2010.

44. Amit Sharma, Jake M Hofman, and Duncan J Watts. Estimating the causal impact of recommendation systems from observational data. In *Proceedings of the 16th ACM Conference on Economics and Computation*, pages 453–470, 2015.

45. Masahiro Sato, Janmajay Singh, Sho Takemori, Takashi Sonoda, Qian Zhang, and Tomoko Ohkuma. Uplift-based evaluation and optimization of recommenders. In *Proceedings of the 13th ACM Conference on Recommender Systems*, pages 296–304, 2019.

46. Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. Bias and debias in recommender system: A survey and future directions. *arXiv preprint arXiv:2010.03240*, 2020.

47. Shuyuan Xu, Juntao Tan, Zuohui Fu, Jianchao Ji, Shelby Heinecke, and Yongfeng Zhang. Dynamic causal collaborative filtering. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management*, pages 2301–2310, 2022.

48. Yang Zhang, Fuli Feng, Xiangnan He, Tianxin Wei, Chonggang Song, Guohui Ling, and Yongdong Zhang. Causal intervention for leveraging popularity bias in recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 11–20, 2021.

49. Daphne Koller and Nir Friedman. *Probabilistic graphical models: Principles and techniques*. MIT press, 2009.

50. Thomas S Richardson and James M Robins. Single world intervention graphs (SWIGs): A unification of the counterfactual and graphical approaches to causality. *Center for the Statistics and the Social Sciences, University of Washington Series*, 128(30):2013, 2013.

51. Tianxin Wei, Fuli Feng, Jiawei Chen, Ziwei Wu, Jinfeng Yi, and Xiangnan He. Model-agnostic counterfactual reasoning for eliminating popularity bias in recommender system. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1791–1800, 2021.

52. Aman Agarwal, Kenta Takatsu, Ivan Zaitsev, and Thorsten Joachims. A general framework for counterfactual learning-to-rank. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 5–14, 2019.

53. Jiawei Chen, Hande Dong, Yang Qiu, Xiangnan He, Xin Xin, Liang Chen, Guli Lin, and Keping Yang. Autodebias: Learning to debias for recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 21–30, 2021.

54. Shuyuan Xu, Juntao Tan, Shelby Heinecke, Jia Li, and Yongfeng Zhang. Deconfounded causal collaborative filtering. *arXiv preprint arXiv:2110.07122*, 2021.

55. Xinyuan Zhu, Yang Zhang, Fuli Feng, Xun Yang, Dingxian Wang, and Xiangnan He. Mitigating hidden confounding effects for causal recommendation. *arXiv preprint arXiv:2205.07499*, 2022.

56. Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Yong Li, and Depeng Jin. Disentangling user interest and conformity for recommendation with causal embedding. In *Proceedings of The Web Conference 2021*, pages 2980–2991, 2021.

57. Masahiro Sato, Sho Takemori, Janmajay Singh, and Tomoko Ohkuma. Unbiased learning for the causal effect of recommendation. In *The 14th ACM Conference on Recommender Systems*, pages 378–387, 2020.

58. Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.

59. Yang Zhang, Wenjie Wang, Peng Wu, Fuli Feng, and Xiangnan He. Causal recommendation: Progresses and future directions. Tutorial for The Web Conference 2022. `https://causalrec.github.io/file/www2022-tutorial-CausalRec.pdf`. 2022-04-26.

60. Dawen Liang, Laurent Charlin, and David M Blei. Causal inference for recommendation. In *Causation: Foundation to Application, Workshop at UAI. AUAI*, 2016.

61. Prem Gopalan, Jake M Hofman, and David M Blei. Scalable recommendation with hierarchical poisson factorization. In *Proceedings of the 31th Conference on Uncertainty in Artificial Intelligence*, pages 326–335, 2015.

62. Qingyao Ai, Keping Bi, Cheng Luo, Jiafeng Guo, and W Bruce Croft. Unbiased learning to rank with unbiased propensity estimation. In *The 41st International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 385–394, 2018.

63. Wenhao Zhang, Wentian Bao, Xiao-Yang Liu, Keping Yang, Quan Lin, Hong Wen, and Ramin Ramezani. Large-scale causal approaches to debiasing post-click conversion rate estimation with multi-task learning. In *Proceedings of The Web Conference 2020*, pages 2775–2781, 2020.

64. Zhenlei Wang, Shiqi Shen, Zhipeng Wang, Bo Chen, Xu Chen, and Ji-Rong Wen. Unbiased sequential recommendation with latent confounders. In *Proceedings of the ACM Web Conference 2022*, pages 2195–2204, 2022.

65. Chang Zhou, Jianxin Ma, Jianwei Zhang, Jingren Zhou, and Hongxia Yang. Contrastive learning for debiased candidate generation in large-scale recommender systems. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3985–3995, 2021.

66. Léon Bottou, Jonas Peters, Joaquin Quiñonero-Candela, Denis X Charles, D Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning systems: The example of computational advertising. *Journal of Machine Learning Research*, 14(11), 2013.

67. Yuta Saito, Suguru Yaginuma, Yuta Nishino, Hayato Sakata, and Kazuhide Nakata. Unbiased recommender learning from missing-not-at-random implicit feedback. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pages 501–509, 2020.

68. Ziwei Zhu, Yun He, Yin Zhang, and James Caverlee. Unbiased implicit recommendation and propensity estimation via combinational joint learning. In *The 14th ACM Conference on Recommender Systems*, pages 551–556, 2020.

69. Yixin Wang and David M Blei. The blessings of multiple causes. *Journal of the American Statistical Association*, 114(528):1574–1596, 2019.

70. Hao Zou, Peng Cui, Bo Li, Zheyan Shen, Jianxin Ma, Hongxia Yang, and Yue He. Counterfactual prediction for bundle treatment. In *Advances in Neural Information Processing Systems*, pages 19705–19715, 2020.

71. Yaochen Zhu, Jing Yi, Jiayi Xie, and Zhenzhong Chen. Deep causal reasoning for recommendations. *arXiv preprint arXiv:2201.02088*, 2022.

72. Jing Ma, Ruocheng Guo, Aidong Zhang, and Jundong Li. Multi-cause effect estimation with disentangled confounder representation. In *International Joint Conference on Artificial Intelligence*, pages 2790–2796, 2021.

73. Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.

74. Harald Steck. Item popularity and recommendation accuracy. In *Proceedings of the 5th ACM Conference on Recommender Systems*, pages 125–132, 2011.

75. Himan Abdollahpouri, Masoud Mansoury, Robin Burke, and Bamshad Mobasher. The unfairness of popularity bias in recommendation. In *RecSys Workshop on Recommendation in Multistakeholder Environments*, 2019.

76. Ziwei Zhu, Yun He, Xing Zhao, and James Caverlee. Popularity bias in dynamic recommendation. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2439–2449, 2021.

77. Zihao Zhao, Jiawei Chen, Sheng Zhou, Xiangnan He, Xuezhi Cao, Fuzheng Zhang, and Wei Wu. Popularity bias is not always evil: Disentangling benign and harmful bias for recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 2022.

78. Zhihong Chen, Jiawei Wu, Chenliang Li, Jingxu Chen, Rong Xiao, and Binqiang Zhao. Co-training disentangled domain adaptation network for leveraging popularity bias in recommenders. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 60–69, 2022.

79. Judea Pearl. Direct and indirect effects. In *Proceedings of the 17th Conference on Uncertainty in Artificial Intelligence, 2001*, pages 411–420, 2001.

80. Yunqi Li, Hanxiong Chen, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. User-oriented fairness in recommendation. In *Proceedings of The Web Conference 2021*, pages 624–632, 2021.

81. Yushun Dong, Jing Ma, Chen Chen, and Jundong Li. Fairness in graph mining: A survey. *arXiv preprint arXiv:2204.09888*, 2022.

82. Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, Juntao Tan, Shuchang Liu, and Yongfeng Zhang. Fairness in recommendation: A survey. *arXiv preprint arXiv:2205.13619*, 2022.

83. Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. Towards personalized fairness based on causal notion. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1054–1063, 2021.

84. Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

85. Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR: Bayesian personalized ranking from implicit feedback. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence*, pages 452–461, 2009.

86. Hao Wang, Huawei Shen, Wentao Ouyang, and Xueqi Cheng. Exploiting POI-specific geographical influence for point-of-interest recommendation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 3877–3883, 2018.

87. Yongfeng Zhang, Xu Chen, et al. Explainable recommendation: A survey and new perspectives. *Foundations and Trends in Information Retrieval*, 14(1):1–101, 2020.

88. Jianxin Ma, Chang Zhou, Peng Cui, Hongxia Yang, and Wenwu Zhu. Learning disentangled representations for recommendation. In *Advances in Neural Information Processing Systems*, 2019.

89. Xiangmeng Wang, Qian Li, Dianer Yu, Peng Cui, Zhichao Wang, and Guandong Xu. Causal disentanglement for semantics-aware intent learning in recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 2022.

90. Juntao Tan, Shuyuan Xu, Yingqiang Ge, Yunqi Li, Xu Chen, and Yongfeng Zhang. Counterfactual explainable recommendation. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management*, pages 1784–1793, 2021.

91. Shuyuan Xu, Yunqi Li, Shuchang Liu, Zuohui Fu, Yingqiang Ge, Xu Chen, and Yongfeng Zhang. Learning causal explanations for recommendation. In *The 1st International Workshop on Causality in Search and Recommendation*, 2021.

92. Paras Sheth, Ruocheng Guo, Kaize Ding, Lu Cheng, K Selçuk Candan, and Huan Liu. Causal disentanglement with network information for debiased recommendations. In *International Conference on Similarity Search and Applications*, pages 265–273, 2022.

93. Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.

94. Raphael Suter, Djordje Miladinovic, Bernhard Schölkopf, and Stefan Bauer. Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness. In *International Conference on Machine Learning*, pages 6056–6065, 2019.

95. Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye Hao, and Jun Wang. CausalVAE: Disentangled representation learning via neural structural causal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9593–9602, 2021.

96. Guy Shani and Asela Gunawardana. Evaluating recommendation systems. In *Recommender Systems Handbook*, pages 257–297. Springer, 2011.

97. Yaochen Zhu, Xubin Ren, Jing Yi, and Zhenzhong Chen. Deep deconfounded content-based tag recommendation for UGC with causal intervention. *arXiv preprint arXiv:2205.14380*, 2022.

98. Jing Yi and Zhenzhong Chen. Debiased cross-modal matching for content-based micro-video background music recommendation. *arXiv preprint arXiv:2208.03633*, 2022.

99. Benjamin M Marlin, Richard S Zemel, Sam Roweis, and Malcolm Slaney. Collaborative filtering and the missing at random assumption. In *Proceedings of the 23rd Conference on Uncertainty in Artificial Intelligence*, pages 267–275, 2007.

100. Benjamin M Marlin and Richard S Zemel. Collaborative prediction and ranking with non-random missing data. In *Proceedings of the 3rd ACM Conference on Recommender Systems*, pages 5–12, 2009.

101. Chongming Gao, Shijun Li, Wenqiang Lei, Jiawei Chen, Biao Li, Peng Jiang, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. KuaiRec: A fully-observed dataset and insights for

evaluating recommender systems. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management*, 2022.

102. Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1059–1068, 2018.

103. Eustache Diemert, Artem Betlei, Christophe Renaudin, and Massih-Reza Amini. A large scale benchmark for uplift modeling. In *Proceedings of AdKDD and TargetAd Workshop*, 2018.

104. Djallel Bouneffouf, Amel Bouzeghoub, and Alda Lopes Gançarski. A contextual-bandit algorithm for mobile context-aware recommender system. In *International Conference on Neural Information Processing*, pages 324–331. Springer, 2012.

105. Yuta Saito, Shunsuke Aihara, Megumi Matsutani, and Yusuke Narita. Large-scale open dataset, pipeline, and benchmark for bandit algorithms. *arXiv preprint arXiv:2008.07146*, 2020.

106. Wenjie Wang, Xinyu Lin, Fuli Feng, Xiangnan He, Min Lin, and Tat-Seng Chua. Causal representation learning for out-of-distribution recommendation. In *Proceedings of the ACM Web Conference 2022*, pages 3562–3571, 2022.

107. Jing Ma and Jundong Li. Learning causality with graphs. *AI Magazine*, 2022.

108. Qian Li, Xiangmeng Wang, Zhichao Wang, and Guandong Xu. Be causal: De-biasing social network confounding in recommendation. *ACM Transactions on Knowledge Discovery from Data*, 2022.

109. Jing Ma, Mengting Wan, Longqi Yang, Jundong Li, Brent Hecht, and Jaime Teevan. Learning causal effects on hypergraphs. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1202–1212, 2022.