

A Comprehensive Survey on Trustworthy Recommender Systems

WENQI FAN, The Hong Kong Polytechnic University, Hong Kong
XIANGYU ZHAO*, City University of Hong Kong, Hong Kong
XIAO CHEN, The Hong Kong Polytechnic University, Hong Kong
JINGRAN SU, The Hong Kong Polytechnic University, Hong Kong
JINGTONG GAO, City University of Hong Kong, Hong Kong
LIN WANG, The Hong Kong Polytechnic University, Hong Kong
QIDONG LIU, City University of Hong Kong, Hong Kong
YIQI WANG, Michigan State University, USA
HAN XU, Michigan State University, USA
LEI CHEN, The Hong Kong University of Science and Technology, Hong Kong
QING LI, The Hong Kong Polytechnic University, Hong Kong

As one of the most successful AI-powered applications, recommender systems aim to help people make appropriate decisions in an effective and efficient way, by providing personalized suggestions in many aspects of our lives, especially for various human-oriented online services such as e-commerce platforms and social media sites. In the past few decades, the rapid developments of recommender systems have significantly benefited human by creating economic value, saving time and effort, and promoting social good. However, recent studies have found that data-driven recommender systems can pose serious threats to users and society, such as spreading fake news to manipulate public opinion in social media sites, amplifying unfairness toward under-represented groups or individuals in job matching services, or inferring privacy information from recommendation results. Therefore, systems' trustworthiness has been attracting increasing attention from various aspects for mitigating negative impacts caused by recommender systems, so as to enhance the public's trust towards recommender systems techniques. In this survey, we provide a comprehensive overview of Trustworthy **Recommender** systems (**TRec**) with a specific focus on six of the most important aspects; namely, Safety & Robustness, Nondiscrimination & Fairness, Explainability, Privacy, Environmental Well-being, and Accountability & Auditability. For each aspect, we summarize the recent related technologies and discuss potential research directions to help achieve trustworthy recommender systems in the future.

*Corresponding author.

Authors' addresses: Wenqi Fan, The Hong Kong Polytechnic University, Hong Kong, wenqifan03@gmail.com; Xiangyu Zhao, City University of Hong Kong, Hong Kong, xianzhao@cityu.edu.hk; Xiao Chen, The Hong Kong Polytechnic University, Hong Kong, shawn.chen@connect.polyu.hk; Jingran Su, The Hong Kong Polytechnic University, Hong Kong, jing-ran.su@connect.polyu.hk; Jingtong Gao, City University of Hong Kong, Hong Kong, jt.g@my.cityu.edu.hk; Lin Wang, The Hong Kong Polytechnic University, Hong Kong, comp-lin.wang@connect.polyu.hk; Qidong Liu, City University of Hong Kong, Hong Kong, qidongliu2-c@my.cityu.edu.hk; Yiqi Wang, Michigan State University, East Lansing, MI, USA, wangy206@msu.edu; Han Xu, Michigan State University, East Lansing, MI, USA, xuhan1@msu.edu; Lei Chen, The Hong Kong University of Science and Technology, Hong Kong, leichen@cse.ust.hk; Qing Li, The Hong Kong Polytechnic University, Hong Kong, csqli@comp.polyu.edu.hk.

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

© 2022 Association for Computing Machinery.

XXXX-XXXX/2022/9-ART \$15.00

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

CCS Concepts: • **Computing methodologies** → **Artificial intelligence**; • **General and reference** → *Surveys and overviews*; • **Security and privacy**;

Additional Key Words and Phrases: Recommender Systems, Trustworthiness, Artificial Intelligence, Robustness, Fairness, Explainability, Privacy, Environmental Well-being, Accountability, Auditability.

ACM Reference Format:

Wenqi Fan, Xiangyu Zhao, Xiao Chen, Jingran Su, Jingtong Gao, Lin Wang, Qidong Liu, Yiqi Wang, Han Xu, Lei Chen, and Qing Li. 2022. A Comprehensive Survey on Trustworthy Recommender Systems. 1, 1 (September 2022), 71 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In the past few decades, the growth of information exchange through the Internet has resulted in extreme information explosion. Thus, recommender systems have been playing an increasingly important role in people’s daily lives via their successful deployments in various user-oriented online services, such as online shopping [130, 216], jobs matching [124, 136], financial product recommendations [21], and medical recommendations [453]. As reported in *MIT Technology Review 2021*¹, TikTok (one of the world’s fastest-growing social networks) recommendation algorithm was awarded as one of the “Top 10 Global Breakthrough Technologies”. More recently, inspired by the great success of Deep Neural Networks (DNNs) in powerful representation learning abilities, DNN-based recommendation techniques have shown impressive performance across a wide range of tasks [103, 106]. For example, a simple yet effective deep learning based recommender system is designed for videos recommendations in YouTube mobile app [75]; A BERT based ranking model shows great power on solving an online job search task in LinkedIn [136]; Some works also leverage Graph Neural Networks (GNNs) to learn helpful representations of users and items in social medias [104, 105, 412].

Despite their great achievements in benefiting human daily lives, recent studies have shown that recommender systems could also bring negative consequences to human society. For example, recommender systems are highly vulnerable to adversarial attacks: someone can generate and spread malicious information, thereby fooling the prediction of recommender systems to wrongly promote or demote items [50, 100]. In addition, recommender algorithms can also implicitly inherit and amplify the biased opinions from the data that collected from society. It will cause the models to have discriminatory biases and unfairness towards under-represented groups, such as people from various genders, races and occupations [235]. Moreover, recommender systems are also susceptible to the risk of leaking users’ private information [402]. For example, someone is able to recover the private information from other users, by only exploiting the model parameters. Furthermore, because of the complicated architecture of DNNs, it is extremely hard to decipher and explain the prediction mechanism of recommender systems [58, 435]. These vulnerabilities of recommender systems can make unreliable recommendation results and produce significant harmful effects in various real-world applications, especially those in safety-critical areas such as finance and healthcare, resulting in severe economic, social, and security consequences. Meanwhile, the concerns on the trustworthiness of recommender systems have significantly hindered the development and deployment of recommendation algorithms. Therefore, how to build trustworthy recommender systems has attracted increasing attention from both academy and industry.

More recently, the European Union (EU) has provided ethics guidelines for promoting Trustworthy Artificial Intelligence (TAI) [322], in which a trustworthy AI system should obey certain ethical principles, such as *Prevention of harm*, *Fairness*, and *Explainability*. These ethical principles must be unfolded in practical requirements for achieving the trustworthiness of AI systems. Meanwhile,

¹<https://www.technologyreview.com/2021/02/24/1014369/10-breakthrough-technologies-2021/>



Fig. 1. Six key dimensions of Trustworthy Recommender Systems (TRec).

building AI systems requires considerable effort from various stakeholders, including system's developers and deployers, end-users, as well as civil society and government. It is worth mentioning that such trustworthy principles in the context of AI systems are also suitable to characterize the trustworthiness of recommender systems, since recommender systems are one of the most successful human-centered AI applications in our daily lives. In this survey, we focus on *SIX* of the most crucial dimensions in achieving trustworthy recommender systems: *Safety & Robustness*, *Non-discrimination & Fairness*, *Explainability*, *Privacy*, *Environmental Well-Being*, and *Auditability & Accountability*, as shown in Figure 1.

Take recommender systems in financial applications as an example, it plays a crucial role in various high-stakes scenarios, such as stock market, insurance products, and loan services. Hence, the recommender systems are expected to make particularly *robust* and *accurate* decisions under any potential security threats. Meanwhile, the demographic attributes of customers such as income, occupation, race, and genders are very *private*, which requires recommender systems to avoid leakage. Thus these information require special and careful protection in recommender systems. Furthermore, it is important that recommendation algorithms ought to mitigate *discriminatory bias* or *unfairness* toward certain groups or individuals for credit card and loan approval. Also, considering the reliability of recommender systems, it is desired to provide *explanations* on how certain decisions are made for various stakeholders, and conduct system auditing periodically from different parties. In addition, training and fine-tuning a large-scale recommendation model typically needs huge energy and natural resources, resulting in problems of global environmental deterioration and resource depletion. Thus, it is important to consider the *sustainability and environmental friendliness* of recommender systems for the benefits of our future generations.

Recent years have witnessed a growing awareness of the trustworthiness of recommender systems in both academia and industry, contributing to the emergence of a considerable body of literature that highlights various dimensions of trustworthy recommender systems [100, 235, 435]. For example, to defend against adversarial attacks [100], methods regarding robust recommendation algorithms have been proposed [103, 433]. Debiasing technologies for building fair recommender systems have been designed for various real-world tasks such as online job matching [124]. Explainable recommendations have been proposed to improve transparency and user satisfaction in the recommendation's decision-making process [364]. Privacy-preserving techniques have been explored to reduce the risk of private data leakage [45]. What's more, there are several surveys

regarding the trustworthiness of recommender systems focusing on specific aspects, such as Safety & Robustness [83, 317], Bias and Fairness [49, 220, 365], and Explainability [59, 435]. In addition, recent surveys [77, 236] give a thorough review of trustworthy AI and Graph Neural Networks (GNNs). As one of the most successful application of human-centered AI systems, it is imperative to systematically summarize the existing achievements and challenges of trustworthy recommender systems. Therefore, in this survey, we provide a comprehensive overview of **Trustworthy Recommender Systems (TRec)** to help researchers and practitioners gain a basic understanding of trustworthy recommender systems, and then have a deeper understanding of the latest progress and facilitate the discussion of the future directions on this demanding topic. More specifically, this survey introduces six key dimensions in realizing trustworthy recommender systems. For each dimension, we introduce its concepts and definitions, as well as provide a taxonomy to review representative and state-of-the-art algorithms. It is worth noting that these SIX dimensions are not independent of each other for building trustworthy recommender systems. At last, we also provide discussions about potential interactions among different dimensions and other potential aspects to achieve the trustworthiness of recommender systems in future directions. The remainder of this survey is organized as follows.

In section 2, we describe the dimension of **Safety & Robustness** from adversarial attacks and defenses aspects, in which a recommender system is required to be robust against adversarial perturbations, so as to make reliable recommendation results. Recent works show that deep recommender systems can inherit vulnerability from DNNs by generating small input perturbations. This vulnerability has raised tremendous concerns about adopting recommender systems in safety-critical domains such as finance and healthcare. Therefore, it is urgent and essential to study the safety and robustness for building safe and reliable recommender systems.

As most recommendation models are designed by our humans and trained from user behavior data, recommender systems can easily inherit human discrimination and unfairness toward certain groups or individuals, resulting in trust loss from various stakeholders. Recently, non-discrimination & fairness in recommender systems receives considerable attention from both academia and industry. In section 3, we detail the dimension of **Non-discrimination & Fairness**, which requires a recommender systems to make fair decisions.

In section 4, we introduce the dimension of **Explainability**, which expects that the working mechanism behind the predictions in recommender systems can be understandable to various stakeholders (e.g., system's developers and end-users). The explainability in recommender systems is treated as an effective way to motivate users to interact with online service, increase users' trust during interactions, and assist algorithms' developers to develop and debug systems.

Since most modern recommender systems are driven by data, recent works found that users' private data such as browsing history and credit card numbers is likely stored and exposed, which increases the risk of data leakage. In section 5, we detail the dimension of **Privacy**, which requires a recommender system to prevent any private information leakage.

Modern recommender systems heavily rely on deep learning techniques to achieve promising performance, in which the demands for large recommendation models will constantly increase, leading to long training time, large storage space, and tremendous energy consumption. A recent study [4] shows that training a model on the Taobao dataset needs 621 minutes with 4 GPUs, whose average GPU power consumption is 56.39W per hour. In section 6, we present the dimension of **Environmental Well-being**, which expects that a recommender system can be sustainable and environmentally-friendly.

In section 7, we discuss the dimension of **Auditability & Accountability**, which expects that the responsibility distribution can be clearly determined for many different parties in the function of recommender systems.

An ideal trustworthy recommender system should satisfy six aforementioned dimensions simultaneously, but most researches only focus on one of them and ignore their potential interactions. In section 8, we introduce the complicated interactions among different dimensions for achieving trustworthy recommender systems. At last, we discuss some future directions to be explored for achieving trustworthy recommender systems in section 9.

Concurrent to our survey, Ge et al. [119] review trustworthy recommender systems from five perspectives, namely explainability, fairness, privacy, robustness, and controllability. Wang et al. [359] describe trustworthy recommender systems in four stages, including data preparation, data representation, recommendation generation, and performance evaluation. In contrast, our work provides a comprehensive survey of trustworthy recommender systems from a computational perspective, discusses the interactions among different dimensions, and provides potential research directions to explore in the future.

2 SAFETY & ROBUSTNESS

Recommender systems play an increasingly important role in high-stake scenarios such as bank loan systems and healthcare recommendations. In recent years, researchers have found that recommendation systems are highly vulnerable to malicious attacks [205], in which modifying a tiny amount of user-item interactions can manipulate recommender systems to produce incorrect results with malicious intentions [50, 100]. These systems cannot be fully trusted and even be denial-of-service attacked if their vulnerabilities are exposed and exploited intentionally. As a result, such vulnerability raises huge concerns when applied to high-stakes tasks, and hinders recommender systems' deployment. For instance, if recommender systems are applied to financial prediction, there may exist some adversaries who attempt to generate fake transactions to deliberately affect the systems' predictions. In healthcare recommender systems, it is possible for an attacker to generate fake cases as a way to mislead the system's diagnosis, posing a threat to patient safety. To prevent attackers from producing harmful effects, recommender systems are required to be robust to artificial perturbations. Besides, it is worth mentioning that understanding the recommendations' weaknesses can provide great opportunities to design new countermeasures against adversarial attacks. Therefore, it is necessary to study adversarial attacks for manipulating recommender systems and to design corresponding defense strategies, so as to improve the reliability and safety of recommender systems.

In this section, we will first introduce the concept and taxonomy of adversarial attacks and defense in recommendation systems. Then we describe how to attack recommendation systems and corresponding defense strategies in detail. All involved methods are summarized in Table 1. Next, we present some practical applications in our daily lives where robustness is critical. Finally, we demonstrate some potential future directions for robust recommendation systems.

2.1 Concepts and Taxonomy

In this subsection, we introduce concepts and taxonomy related to the safety & robustness of recommendation systems from the perspective of adversarial attack and defense.

2.1.1 Attackers' Goal. In recommender systems, according to adversaries' goals, we can divide them into the two following categories.

- **Target Attacks:** The goal is to promote/demote a set of target items in recommendation systems, such that target items can be recommended to as many/few users as possible. This goal is to manipulate the exposure rate of target items to achieve attackers' desires.

- **Untarget Attacks:** In this setting, there is no specific items to be promoted or demoted in untarget attack. Its goal is to degrade a recommendation system's overall performance, so as to reduce users' online experience and satisfaction.

2.1.2 *Attack Stage.* Generally, adversarial attacks in recommendation systems can be divided into two types according to the attack stage: evasion attack and poisoning attacks, which can affect recommendation systems in the inference and training phases, respectively.

- **Evasion Attack (Inference/Test Stage):** Evasion attack happens during the model service (test) stage. For instance, given a fixed well-trained model, attackers can modify a target user's profile, such as historical interaction logs, so its recommendation outcome is changed.
- **Poisoning Attack (Training Stage):** Poisoning attack, also called shilling attack, occurs during the data collection phase before model training. The attacker intends to inject fake users into the training data of recommendation systems, so that trained model's prediction behavior can be controlled with malicious desires.

2.1.3 *Attackers' Knowledge.* To conduct adversarial attacks, the knowledge that attackers are allowed to access target recommender systems can heavily affect the attacking strategies and performance to achieve the adversarial goal. Typically, auxiliary knowledge on a target recommender system includes the target model's architecture and parameters, and datasets, etc. In general, adversaries can conduct three different types of attacking strategies according to their accessibility to the target recommender systems' knowledge, including white-box, grey-box, and black-box attacks.

- **White-box Attacks:** In this setting, attackers can get all information about the target recommendation system, including training data, recommendation architecture and parameters. One widely used strategy is to utilize the gradient to assist in generating adversarial perturbations. Since it is difficult for attackers to obtain such complete knowledge in the real world, this type of attacks cannot pose severe threats. However, researchers can utilize this type of attack to analyze the robustness of the target system in the worst case.
- **Grey-box Attacks:** In this setting, attackers can only get partial information about the target recommender system to conduct attack. Compared with white-box attacks, grey-box attacks are more practical and dangerous, since such limited information is easy for attackers to obtain. For instance, users' reviews and items' ratings information on Amazon are easily collected, which motivates adversaries to take advantage of such available training dataset to train a surrogate model and perform white-box attack subsequently.
- **Black-box Attacks:** In this case, it is challenging for attackers to access the target models of recommender systems and their training data. This setting is more realistic and practical for existing adversarial attacks and have attracted increasing attentions recently. Typically, black-box attacks tend to perform query the target recommender systems for updating the attacking strategies.

2.1.4 *Adversarial Perturbation Type.* With malicious goals, adversaries can add adversarial perturbations in different ways by considering various scenarios. In general, such data perturbations are implemented via adding fake user profiles into user-item interactions, modifying users attributes information (e.g., age, gender, occupation, etc.), and modifying item side information such as the attributes and description of movies.

2.1.5 *Countermeasure Strategies Against Adversarial Attacks (Defense Methods).* To prevent the harm from adversarial attacks, its countermeasure strategies in recommendation systems can be divided into two categories: *Perturbation Detection* and *Adversarial Training*.

Table 1. Taxonomy of related methods.

	Taxonomy	Related Research
Attack	Heuristic Methods	[37, 38, 40, 198, 261, 372]
	Gradient-based Methods	[72, 73, 108, 109, 205, 230, 343, 375]
	Reinforcement Learning-based Methods	[50, 100, 324]
Defense	Detection Methods	[39, 118, 187, 202, 254, 256, 309, 425, 433]
	Adversarial Robust Training Methods	[47, 145, 341, 355, 416]

- **Perturbation Detection.** This kind of defense strategies is to identify perturbations data and remove them for resisting adversarial attacks in recommender systems.
- **Adversarial Training.** This is a widely used strategy to resist adversarial attacks by enhancing the robustness of recommender systems.

2.2 Representative Attack Methods

In this subsection, we mainly introduce poisoning attack (i.e., shilling attack), which is the most widely studied mainstream attacks in recommendation systems. First, we give a unified formulation of poisoning attacks. Then, we present representative methods from various aspect, including heuristic methods, gradient-based methods, and reinforcement learning-based methods.

2.2.1 A Unified Formulation of Poisoning Attack. The attackers' goal is to inject well-designed fake user profiles into recommender systems to manipulate the recommendation's output, as illustrated in Figure 2. In general, this attacking process can be formulated as a bi-level optimization problem. Mathematically, given a set of users $U = \{u_1, u_2, \dots, u_{|U|}\}$, a set of items $I = \{i_1, i_2, \dots, i_{|I|}\}$, and user-item interactions matrix $R \in \mathbb{R}^{|U| \times |I|}$, attackers aims to design a set of fake user $\hat{U} = \{\hat{u}_1, \hat{u}_2, \dots, \hat{u}_{|\hat{U}|}\}$ with the fake user-item interaction data $\hat{R} \in \mathbb{R}^{|\hat{U}| \times |I|}$ to achieve their adversarial goals, which can be formulated as \mathcal{L}_{adv} :

$$\min_{\hat{U}} \mathcal{L}_{adv}(\theta^*), \quad \text{s.t.} \quad \theta^* = \arg \min_{\theta} (\mathcal{L}_{rec}(R, O_{\theta}) + \mathcal{L}_{rec}(\hat{R}, O_{\theta})), \quad (1)$$

where θ is the recommendation system's parameters, O_{θ} is the system's predictions with parameters θ , and \mathcal{L}_{rec} can be a general recommendation objective. The objective function \mathcal{L}_{adv} can be determined by the specific goal, such as promoting/demoting some items or destroy the target system's utility. Generally, the fake users \hat{U} cannot be arbitrarily designed, since defenders can easily detect such fake users who have a large discrepancy with normal users. In addition, the perturbations will be constrained by pre-defined budget for attacking, like the number of fake users and their interactions.

2.2.2 Heuristic Methods. Some early studies are based on hand-engineered fake user profiles to poison a recommendation system. A straightforward solution is that attackers can design fake users that assign high scores to target items and a low score to random others to promote target items [198]. Burke et al. [37] consider to interact with some popular items rather than random items, so that the fake users can have more impact on normal users. Some works [38, 261] introduce low-knowledge promotion attack methods such as the random attack, average attack, bandwagon attack, and segment attack. Similarly, Williams and Mobasher [372] propose several demotion attack methods: love/hate attack and reverse bandwagon attack.

However, these attacking methods have several limitations. First, fake user profiles generated by heuristics tend to have distinct characteristics, sometimes even forming self-established clusters [40],

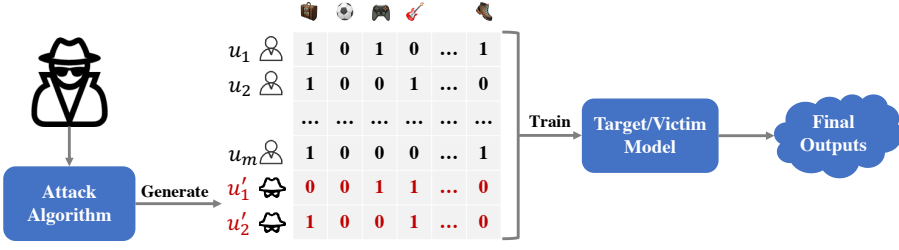


Fig. 2. The poisoning attack. Attackers inject well-designed faker users into training data to manipulate target model's behaviors.

making them easy to be detected [255]. In addition, heuristic methods strongly rely on prior knowledge to generate fake profiles that influence normal users. Finally, adversarial goals cannot be optimized by most existing heuristic methods, leading to low attack success rates.

2.2.3 Gradient-based Methods. Unlike heuristic methods, the process of poisoning attacks can be formulated as the optimization problem as shown in Eq. 1, where the generation process of fake user profiles can be formulated as a bi-level optimization problem, consisting of a constraint called inner objective and the main goal called outer objective. Here we take Matrix Factorization (MF) as the target recommendation model and then introduce it with gradient-based methods. Given the user-item interaction data $R \in \mathbb{R}^{|U| \times |I|}$, the MF model aims to learn the user embedding matrix $P \in \mathbb{R}^{|U| \times d}$ and item embedding matrix $Q \in \mathbb{R}^{|I| \times d}$, such that the inner product PQ^T can approximate the interaction R under observed values (i.e., non-zero values in R). And then, the well-learned PQ^T can be used to predict unobserved user-item pairs. Mathematically, when considering fake users, the objective of MF learning can be formulated as:

$$\begin{aligned} & \min_{\theta} \mathcal{L}_{rec}(R, O_{\theta}) + \mathcal{L}_{rec}(\widehat{R}, O_{\theta}), \\ & \rightarrow \min_{P, Q, \widehat{P}} \sum_{u,i} (R_{u,i} - P_u Q_i^T) + \sum_{\widehat{u},i} (\widehat{R}_{\widehat{u},i} - \widehat{P}_{\widehat{u}} Q_i^T) + \lambda (\|P\|^2 + \|Q\|^2 + \|\widehat{P}\|^2), \end{aligned} \quad (2)$$

where \widehat{P} denotes the embedding matrix of fake users, and λ is a hyper-parameter to regularize model for addressing overfitting.

The adversarial objective function \mathcal{L}_{adv} varies for different malicious purposes. Here, we focus on promoting target item k to all normal users, which can be formulated as follows:

$$\begin{aligned} & \min_{\widehat{R}} \mathcal{L}_{adv}(\theta^*) = - \sum_{u \in \mathcal{U}} \log \left(\frac{\exp(r_{uk})}{\sum_{i \in \mathcal{I}} \exp(r_{ui})} \right), \\ & \text{s.t. } \theta^* = \arg \min_{\theta} \mathcal{L}_{rec}(R, O_{\theta}) + \mathcal{L}_{rec}(\widehat{R}, O_{\theta}). \end{aligned} \quad (3)$$

In this case, the model's parameters θ are embedding matrices P and Q . The goal is to minimum the loss, hoping that all normal users' prediction on target item k is greater than other items. Finding optimal fake users \widehat{U} in Eq. 1 is equivalent to optimizing its rating matrix \widehat{R} in Eq. 3, which can be achieved by the projection gradient descent method as follows:

$$\widehat{R}^{t+1} = \text{Proj}_{\mathbb{R}}(\widehat{R}^t - \alpha \cdot \nabla_{\widehat{R}} \mathcal{L}_{adv}(\theta^*)), \quad (4)$$

$$\nabla_{\widehat{R}} \mathcal{L}_{adv}(\theta^*) = \nabla_{\widehat{R}} \theta^* \nabla_{\theta^*} \mathcal{L}_{adv}, \quad (5)$$

where $\text{Proj}_{\mathbb{R}}$ denotes the projection operator under the feasible region \mathbb{R} , and α is the step size. Note that the second gradient $\nabla_{\theta^*} \mathcal{L}_{adv}$ can be easily obtained, while it is challenging to obtain the first gradient term (i.e., $\nabla_{\widehat{R}} \theta^*$) because of involving the minimization term θ^* .

Li et al. [205] pioneer a gradient-based attack method for factorization-based recommendation systems. The key technique of their method is to approximately compute the Eq. 5 based on first-order KKT conditions. Christakopoulou et al. [73] use the zero-order optimization method in evolutionary algorithms to find the gradient's direction. Specifically, they make several minor changes on fake user profiles \hat{R} to evaluate the adversarial loss \mathcal{L}_{adv} , and then update user profile by the change that makes the loss smaller. In [108], Fang et al. propose the differentiable hit ratio loss to generate fake users for top-N recommendation systems and leverage first-order stationary condition to approximately compute the Eq. 5. To improve imperceptibility, they select filler items relying on the value of the final solution of Eq. 4 and give ratings sampled from the distribution of normal users' interactions to these filler items. Moreover, they only select a subset of critical users with the influence function for efficient computation. Different from the previous approximation methods, Tang et al. [343] compute the exact solution based on the high-order gradient, while the method requires more computing resources.

In addition, some studies take advantage of Generative Adversarial Networks (GAN) to approximate real users behaviors for attacking, so that the generative fake users are undetectable. For instance, Christakopoulou and Banerjee [72, 73] first train a GAN on real user profiles so that the generator of the GAN can generate faker users having the same distribution as normal ones; then the generator's outputs are treated as the initialization for gradient-based attacks. Different from methods [72, 73], Lin et al. [230] propose an end-to-end GAN-based attacking method AUSH by directly training a GAN with a loss that can include attacks as well. Further, Wu et al. [375] propose a TripleAttack method, where an extra influence module provides the guideline of the generator outputting influential fake users.

2.2.4 Reinforcement Learning-based Methods. The gradient-based poisoning attacks have achieved good performance under the white-box setting, which cannot be directly applied into attacking black-box recommender systems due to extremely limited knowledge towards target systems accessed by adversaries. Recently, some studies have leveraged Deep Reinforcement Learning (DRL) to learn attacking policy strategies via query rewards under the black-box setting. More specifically, the DRL based attacking process can be formulated as a Markov Decision Process (MDP) to learn a policy $\pi(s_t) \rightarrow a_t$ for outputting the action a_t under state s_t .

To perform attacks under black-box setting, PoisonRec [324] proposes a model-free reinforcement learning based framework for generating fake user profiles. More specifically, a Biased Complete Binary Tree (BCBT) is constructed to model the item sampling process, which can help significantly reduce the time complexity in a hierarchical action space. Moreover, in order to improve the quality of fake user profiles, Chen et al. propose a knowledge-enhanced black-box attacks for recommendations (KGAttack) [50], which takes advantage of items' attribute features (treated as Knowledge Graph) to enhance the process of sampling items. More specifically, a graph neural networks and a recurrent neural network are introduced to model knowledge graph for enhancing state representation learning. Meanwhile, in order to effectively select items from the large-scale discrete action space (i.e., the massive item sets), hierarchical policy networks are proposed to decompose the selection process into two actions, including anchor item selection and next item picking.

Furthermore, instead of generating fake user profiles from scratch, Fan et al. propose a novel copy mechanism to obtain real user profiles for black-box recommender systems (CopyAttack) [100]. In detail, cross-domain user profiles in source domain, which can share similar online behaviors with target recommender systems, are copied into target domain for promoting a set of items. However, selecting real user profiles in source domain based on reinforcement learning is challenging due to the large-scale user profiles (i.e., discrete action), resulting in inefficiency and ineffectiveness. To

address such challenges, CopyAttack proposes hierarchical-structure policy gradient in balanced hierarchical clustering tree over cross-domain user profiles to search a path from the root to a certain leaf of the tree, where each non-leaf node represents as a policy gradient network and a leaf node represents a user profile. In addition, masking mechanism is introduced to exclude user profiles which does include target items, so as to further reduce the action space. At last, a crafting policy gradient network is introduced to refine the raw user profiles, so as to decrease the attacking budgets and reduce some noise. It is worth noting that KGAttack and CopyAttack use some spy users as proxy to obtain reward for optimizing the proposed DRL based attacking framework, while PoisonRec uses the number of Page View on the target item as the reward.

2.3 Representative Defense Methods

Various attacking methods expose the high vulnerability of modern recommendation systems, which motivates researchers to design countermeasure strategies against adversarial attacks. In this subsection, we introduce some representative defense methods that improve the robustness of recommendation systems. Generally, there are two pathways to defend against an adversarial attack in recommendation systems: (1) *Detection* methods to localize anomalies (e.g., fake user profiles); (2) *Adversarial Robust Training* to make recommender systems more robust against adversarial attacks.

2.3.1 Detection. In early year research, some works utilize machine learning-based classifiers, e.g., SVM and KNN, to detect anomalies and outliers in recommender systems. These methods train a classifier using specific attributes of user profiles, which works well against heuristic attacks. For example, Burke et al. [39] study generic attributes of user profiles and exploit them to conduct defense. In particular, they propose three variants strategies to measure discrepancies between user’s ratings and item’s average ratings. Zhang et al. [425] propose a hybrid detection method that combines SVM and Hilbert–Huang transform, where Hilbert–Huang transform is used to capture spectrum-based features of series rate values of each user and then use the features train an SVM classifier to distinguish fake users. Besides, some studies explore unsupervised learning approaches to cluster outlier data, relying on statistical attributes of the whole dataset. For instance, Mehta [254] finds that the soft-cluster method based on Probabilistic Latent Semantics Analysis is effective to determine fake users. Bhaumik et al. [27] use k-means method to cluster instances and identify fake users from small clusters.

More recently, researchers have adopted deep learning models to develop more effective defense strategies. Gao et al. [118] propose a LSTM-based model to encode a series of user behaviors to indicate whether user profiles are suspicious. Zhang et al. [433] propose a unified framework for both recommendation and attack detection based on GNNs, which can adaptively detect fake users in the process of learning users and items representations. Specifically, the detection component is proposed to dynamically adjust the users’ weights for representation learning according to their probability of being fake. Shahrabi et al. [309] propose a semi-supervised algorithm to detect fake user profiles using SeqGAN [414] that can deal with discrete sequential data compared with vanilla GAN, which can learn the distribution of normal users’ behaviors over a partial dataset that is definitely normal, so as to identify anomalies in recommender systems.

2.3.2 Adversarial Robust Training. Adversarial robust training endows a model with the ability to tolerate adversarial perturbations instead of detecting anomalies and outliers. In general, adversarial training contains two alternating processes: (1) generating perturbations that can confuse a recommendation model; (2) training the recommendation model along with generated perturbations. Mathematically, this process can be formulated as a min-max game as follows:

$$\min_{\theta} \max_{\eta} \mathcal{L}(X + \eta, \theta), \quad (6)$$

where θ is the recommendation model's parameters, η indicates perturbations, and \mathcal{X} denotes the original normal dataset.

In [145], He et al. propose an adversarial training method - Adversarial Personalized Ranking (APR) to enhance the robustness of BPR based Matrix Factorization method, which aims to perturb the embeddings of users and items by leveraging adversarial training strategy, instead of perturb raw data input. Mathematically, the optimization objective can be formulated as:

$$\min_{\theta} \mathcal{L}_{APR}(\theta) = \min_{\theta} \max_{\eta} \mathcal{L}_{BPR}(\theta) + \lambda \mathcal{L}_{BPR}(\theta + \eta), \quad (7)$$

where $\theta = \{P, Q\}$ denotes the parameters (i.e., users and items embeddings) in MF based recommendation methods. η is perturbations added to model parameters θ . They demonstrated that APR is not only an effective defensive strategy but also boosts generalization performance. Further, by extending APR to multimedia recommender systems, Tang et al. [341] propose Adversarial Multimedia Recommendation (AMR) framework by optimizing visual-aware BPR (VBPR) objective. More specifically, adversarial perturbations are incorporated to visually-aware item space extracted by CNN encoder, so as to enhance the robustness of multimedia recommendation. In addition, by considering the robustness of tensor-based recommendations, Chen and Li [47] incorporate adversarial training to enhance the robustness of pairwise interaction tensor factorization [294] for context-aware recommendations.

2.4 Application

In this section, we introduce robustness issues in two real-world applications to demonstrate the necessity of building adversarial robust recommendation systems.

- **E-health recommendation.** Recent work shows that more than 42% of clinical misdiagnoses are caused by inadequate doctors who are unfamiliar with certain drugs [20]. In order to reduce the misdiagnosis rate and the burden on doctors in such safety-critical scenario, intelligent systems are developed to assist doctors in making clinical diagnoses and perform drug package recommendations [453]. People can trust drug recommender systems only if the systems can obtain high accuracy and resist potential attacks without any vulnerability. Thus, it is important to investigate the vulnerability of recommender systems, so as to the enhance their robustness for trustworthy recommender systems.
- **E-commercial recommendation.** Online e-commercial platforms, e.g., Amazon, Taobao, etc., dominate people's daily shopping needs. The prevalence of such platforms relies on a reliable recommender system to continuously recommend products of interest to users by exploiting users' previous purchases. Due to the property of openness in such online services, adversaries can easily generate fake users and reviews to maliciously mislead people behaviors when they shop online [421], which can damage to shoppers and businesses. Thus, enhancing the robustness against such perturbations is becoming more and more important for building trustworthy recommender systems.

2.5 Surveys and Tools

In this subsection, we sort out the existing surveys on safety & robustness in recommender systems and a useful toolkit evaluating the robustness of recommender systems to facilitate researchers in this field.

2.5.1 Surveys. The robustness of recommendation systems has been widely studied for a long time. Zhang et al. [424] gives a comprehensive taxonomy about shilling attack strategies, evaluation metrics, and defense methods in recommender systems. In [134], Gunes et al. provide a summary of attacks against various collaborative filtering recommendations and detection methods. They also

introduce cost/benefit analysis, a new attribute for classification shilling attacks, and discussions on future directions. Similarly, Si and Li [317] summarize shilling attacks and defenses from their style and discuss future directions to improve the robustness of recommendation systems. In [346], Truong et al. study the effect of adversarial training on recommendation systems and analyze its properties and designs. Recently, Deldjoo et al. provide surveys [82, 83] about adversarial machine learning in recommender systems (AML-RecSys). They review AML methods in the traditional machine learning field and further survey AML in recommendation systems from two perspectives: adversarial strategy and the GAN-based model.

2.5.2 Tools. While various toolkits have been developed for researchers to build recommendation systems and evaluate the systems' performance conveniently, there are not too many toolkits for the robustness of recommendation systems. The only one we can find out is RGRRecSys [273], which allows researchers to easily evaluate recommender system robustness with respect to attacks and other dimensions.

2.6 Future Directions

Robustness in recommendation has always been an important research topic; however, many open problems and challenges are still not well explored. In this section, we point out the potentially valuable research directions. The main attack and defense research are aimed at the collaborative filtering-based model and consider manipulate user-item interactions (i.e., generating faker user profiles or perturbing real user profiles). In practice, modern recommendation systems can incorporate many sources data in various scenarios, such as social connections [99, 102] and knowledge graph [50], which motivates to develop recommender systems based on various techniques, such as reinforcement learning [446] and graph neural networks [85, 101]. Thus, an important issue is to investigate the vulnerability of different target recommender systems, so as to improve their trustworthiness from robustness aspect. For adversarial robust training in defense methods, instead of raw data space, most existing methods works on parameters space by adding adversarial perturbations, which may limit the robustness improvement. Another direction is to generate adversarial perturbations on user-item interactions to perform adversarial robust training.

3 NON-DISCRIMINATION & FAIRNESS

To be widely deployed in high-stakes scenarios such as finance and healthcare [348], a trustworthy recommender system should avoid exhibiting discriminatory behaviors in human-machine interactions and guarantee to make fair decisions for users from certain groups. Unlike the general machine learning tasks such as classification [44, 252], fairness in recommendation algorithms has several unique characteristics: first, multi-sided fairness needs to be considered [36] since recommender systems serve users and item providers as a two-sided platform; second, discriminatory bias might exist everywhere in the dynamic feedback loop between human and recommender systems [249] and even get amplified without appropriate interventions. The manifestation of bias and prejudice severely hampers the promotion of trustworthiness, and affects the long-term benefits of recommender systems. For example, in a job recommendation platform like LinkedIn, if the platform exhibits gender discriminatory bias [124], e.g., women are being recommended with fewer job opportunities or lower-payment jobs compared to men, it will cause detrimental effects from both the ethical and legal aspects. In movie recommender systems with popularity bias, popular movies would always be over-recommended, which will not only intensify the homogenization of users but also reduce the exposure opportunities of other equally qualified but less popular movies [3]. To alleviate these issues, it is necessary to analyze the potential bias and mitigate the unfairness in recommender systems [95, 167].

In this section, we will first introduce fundamental definitions and concepts regarding fairness in recommender systems, where we provide a detailed taxonomy of the causes of unfairness, the definitions of fairness criterion and the evaluation metrics of fairness. Then, we review and categorize existing bias mitigation methods, which can enhance the fairness performance from different perspectives in recommender systems. Lastly, we discuss the applications and future directions in this field. We hope that researchers can benefit from the broad overview of bias and fairness issues in recommender systems and reach a consensus on pushing for further advances in this field.

3.1 Concepts and Taxonomy

In this subsection, we first introduce the origins of unfairness in recommender systems, then present a taxonomy of fairness definitions and several standard fairness evaluation metrics. It is worth noting that there exists a large number of fairness definitions and evaluation metrics due to the great magnitude of related works. Therefore, we summarize the categories from several typical perspectives [223, 365].

3.1.1 Bias. The discriminatory bias in recommender systems often leads to unfairness issues [13, 94, 236], which means that the system unfairly treats certain individuals or protected groups by providing poorer recommendation quality. Though the sources of bias in recommender systems can be various [271], we can divide the recommender systems' feedback loop into three parts from a bird's-eye view [49]: user \rightarrow data (*data collection*), data \rightarrow recommendation model (*model training*), and recommendation model \rightarrow user (*model serving*), and categorize the potential bias as follows:

- **Data Bias** is the distribution difference between the collected training data and the ideal test data. It pre-exists in the data generation process [44] and may come from many aspects. Following [49], data bias can be further categorized into the following four groups:
 - *Selection Bias* refers to that users' selective rating behavior [251] and the observed ratings do not fully reveal the true ratings. As a consequence, the collected data is missing not at random (MNAR).
 - *Exposure Bias* means that unobserved interactions in implicit feedback do not necessarily disclose users' disliked items since users are merely exposed to a small portion of items.
 - *Conformity Bias* indicates that users behave similarly to other group members, even if what they do goes against their judgment [49].
 - *Position Bias* refers to the observation that items in the higher positions of a recommendation list are more likely to receive interaction no matter how highly relevant they are to users [49].

In addition, collected feedback data can be biased by other factors, such as *marketing bias* [351], indicating that consumers' interactions may be affected by the human model's profile in a product image (a reflection of a product's marketing strategy) and result in the under-representation of particular niche markets.

- **Model and Result Bias** refers to the bias in the algorithm design and model results [16], in which recommendation algorithms tend to exhibit bias and generate unfair recommendation results (e.g., popularity bias), when optimizing without any fairness constraints.
 - *Popularity Bias* happens when popular items are over-recommended compared to what their popularity warrant [49].
- **Feedback Loop Bias** refers to the reinforced bias introduced by the RS feedback loop mechanism [71]. Unfair recommendations would influence users' behaviors in the online serving process, which makes the observed feedback encode biases. Moreover, biased users'

behavior data would enlarge the model's discrimination when collected for model training. To be specific, popular items attract a large traffic volume in recommender systems [2] and a biased model will provide such popular items with better recommendation quality (i.e., precision) than those unpopular items. Consequently, online serving will lead to a greater traffic volume gap between popular and unpopular items.

There are other causes of unfairness, such as conflicts between different fairness requirements [70, 192]. In this case, fulfillment of one fairness criterion would violate some other fairness requirements. In the subsequent part, we will introduce details of different fairness definitions.

3.1.2 Fairness. Previous works [223, 266] have presented various fairness metrics to quantify the effects of discriminatory bias in recommender systems. Though there is still no consensus on a general definition of fairness, it can usually be classified into procedural fairness and outcome fairness [365].

Procedural Fairness represents procedural justice in decision-making processes, which is critical since it affects people's trust and cooperation of recommender systems. Most works [133, 204] mainly focus on whether the usage of input features in decision-making processes is fair. For example, in [133], new scalar fairness measures are introduced to explicitly account for individuals' moral sense of whether it is fair to use the input features in decision processes.

Outcome Fairness holds that fairness-aware models ought to exhibit fair outcome performance [107], which is also called as distributive fairness [133]. Since there are large amounts of definitions falling under this category, we group related concepts as follows:

- **Group by subject.** Since recommender systems are multi-stakeholders connecting users with items, fairness requirements from different sides need to be considered. *User fairness* refers to whether different users obtain fair recommendation, such as equal recommendation accuracy [95] or equal recommendation explainability [114]. *Item fairness*, also called provider fairness, refers to whether different item groups are fairly treated, such as equal prediction errors for ratings [289], equal recommendation probabilities of different item groups given the items are truly liked by users [463], and equal ranking position that is proportional to relevance score [28]. *Joint-sided fairness* considers both users and items, such as the fairness of recommendation results' quality on the user side and the fairness of providers' exposure on the item side [279, 384].
- **Group by granularity.** When looking into the granularity of resource allocation processes, outcome fairness can be classified as: *Group Fairness* and *Individual fairness*. Group Fairness refers to the performance parity among different sensitive groups, which generally specified by user/item sensitive attributes (i.e., gender or race) [393]; *Individual fairness* [112] requires that similar individuals should be treated equally [178]. In general, individual fairness can be achieved when two similar individuals always have similar predictions in the output space.
- **Others.** There are quantities of works defining fairness from other perspectives, and here we choose the following five representative definitions. *Causal fairness* aims to eliminate the causal relation between sensitive attributes and model predictions. By incorporating additional structural knowledge regarding how variables propagate on a causal graph, causal fairness is achieved when recommendation results remain the same in the factual and counterfactual world for each possible user [221]. *Personalized fairness* takes personalized demands from users into consideration, where users are free to select sensitive attributes they care about [386], and recommendation models provide flexible support in filtering out any chosen sensitive attributes from original feature representations. *Explainable fairness* aims to explain why the recommendation model is unfair and provide insights for further improvement. For instance, in [121], they develop an explainable counterfactual framework to explain which

input features significantly influence the fairness-utility trade-off in recommendations. Then by alleviating the negative influence of the detected features when doing fair learning, the recommendations can achieve a better fairness-utility trade-off. *Rawsal max-min fairness* focuses on maximizing outcome performance of the worst individual or group [462]. *Dynamic fairness* requires guaranteeing fairness under dynamic factors' influence, such as users' evolving preferences or item's popularity degree change as a result of user interactions throughout the recommendation process [423].

3.1.3 Fairness Evaluation Metrics. Next, we present corresponding evaluation metrics for the fairness definitions mentioned above.

- **Absolute Difference (AD)** measures utility differences between the disadvantaged group G_0 and the advantaged group G_1 , which can be formulated as:

$$AD = |u(G_0) - u(G_1)|, \quad (8)$$

where $u(G)$ denotes as the group utility function, which is used to calculate the average rating prediction scores or the average ranking performance (i.e., NDCG or F1-score) of the group G . A low group utility difference value indicates fair recommendation performance.

- **Variance** measures the performance dispersion at the group-level or individual-level [289]. It can be calculated by adding up and averaging the performance difference between any two different groups/individuals, e.g., $v_i, v_j \in \mathcal{V}$. Here, \mathcal{V} represents the whole set of individuals or groups.

$$\text{Variance} = \frac{1}{|\mathcal{V}|^2} \sum_{v_i \neq v_j} (u(v_i) - u(v_j))^2. \quad (9)$$

- **Min-Max Difference (MMD)** measures the difference between the maximum and the minimum score value of all allocated utilities, which can be adopted to reflect the disparity of multiple item groups' exposure opportunities [139].

$$MMD = \max \{u(v), \forall v \in \mathcal{V}\} - \min \{u(v), \forall v \in \mathcal{V}\}. \quad (10)$$

- **Entropy** usually reflects the uncertainty and disorder of a system, which can also be adopted to evaluate the inequality of item exposure opportunities in recommendations [279].

$$\text{Entropy} = - \sum_{v \in \mathcal{V}} p(v) \cdot \log p(v). \quad (11)$$

- **KL-Divergence** measures the difference between two probability distributions. This metric can be used to calculate the difference between the item groups' exposure distribution p and their historical exposure q in recommendations [122]. A lower KL-divergence value indicates fairer recommendations. Note that the JS-divergence can be viewed as the symmetrical version of KL-divergence.

$$D_{KL}(p, q) = - \sum_{v \in \mathcal{V}} \frac{p(v)}{q(v)}. \quad (12)$$

3.2 Methods

In this subsection, we introduce some representative fair recommendation methods. Based on the specific stage that these methods can be applied in the whole recommendation pipeline, we categorize existing methods into the following three types: **Pre-processing**, **In-processing**, and **Post-processing**.

Table 2. Taxonomy of related methods

Taxonomy	Method type	Related research
Pre-processing	Data Re-sampling	[95]
	Adding Antidote Data	[289]
In-processing	Regularization & Constrained Optimization	[26, 351, 393, 409, 461]
	Adversarial Learning	[33, 207, 215, 221, 285, 379, 380]
	Reinforcement Learning	[120, 122, 244]
	Causal Graph	[121, 162, 387, 452]
	Others	[31, 110, 167, 224]
Post-processing	Slot-wise Re-ranking	[124, 185, 189, 243, 262, 300, 305] [306, 323, 328, 405, 419]
	User-wise Re-ranking	[28, 253, 304, 318]
	Global-wise Re-ranking	[87, 114, 219, 250, 279, 335, 384, 462]

3.2.1 Pre-processing Methods. Pre-processing methods usually directly modify the training data, aiming to remove data bias before training recommendation models.

The advantage of these methods is their flexibility, since they are decoupled with recommender systems. However, there are multiple steps between the data and the final output, which also indicates that performance gains in the pre-processing step may not be maintained by the following steps (i.e., re-ranking). There are two typical pre-processing methods as follows:

- **Data Re-sampling.** Data re-sampling aims to balance data sets so that the data size of each sensitive group or individual is close. Ekstrand et al. [95] conduct an empirical study on the effectiveness of several collaborative filtering algorithms across multiple datasets, which are stratified by the users' sensitive attributes. Experiment results show that different demographic groups obtain different utilities due to the imbalanced data distribution. Based on this observation, they propose to balance the ratio of various user groups via a re-sampling strategy and then re-train recommendation algorithms, while this approach achieves minor fairness improvement.
- **Adding Antidote Data.** Augmenting input with additional data is another alternative for pre-processing, where the augmented data is designed to promote recommender systems' fairness and thereby can be viewed as antidote data. Rastegarpanah et al. [289] design data augmentation strategies to address the unfairness issues, where additional antidote data is optimized via gradient descent methods for satisfying the fairness objective function. This method mitigates unfairness more effectively than the re-sampling method but consumes more time.

3.2.2 In-processing Methods. In-processing methods aim to mitigate bias in the model training process. Based on their optimization perspectives, we categorize relevant methods as follows:

- **Regularization and Constrained Optimization.** The in-processing fairness methods [26, 351, 393, 409, 461] are primarily based on regularization and constrained optimization, where various fairness criterion are formulated as constraints or regularizers for guiding model optimization. Formally, the loss function consists of a traditional recommendation loss \mathcal{L}_{rec} and a fairness-related regularization loss $\mathcal{L}_{fair-reg}$ as follows:

$$\min_{\theta} \mathcal{L}_{rec}(\theta) + \mathcal{L}_{fair-reg}(\theta). \quad (13)$$

In general, this direct regularization approach is to integrate fairness metrics into the overall loss function. In [393], several fairness metrics are specifically designed for the group recommendation scenario, which can then be transformed into model regularizers and constitute

a multi-objective optimization problem from the perspective of Pareto Efficiency. Wan et al. [351] design two fairness metrics for evaluating marketing bias, where *rating prediction fairness* measures the global parity of prediction errors across different user-product market segments, and *product ranking fairness* measures the KL-divergence of the frequency distribution of market segments between real and predicted interactions. By regularizing the correlation between prediction errors and market segment distribution, the recommendation model can explicitly calibrate prediction errors' equity across different market segments. In [409], four new metrics are proposed to address different forms of unfairness in collaborative filtering methods. For example, the absolute unfairness metric is denoted as follows:

$$U_{abs} = \frac{1}{n} \sum_{i=1}^n \left| |E_{adv}[y]_i - E_{adv}[r]_i| - |E_{-adv}[y]_i - E_{-adv}[r]_i| \right|, \quad (14)$$

where $E_{adv}[r]_i$ and $E_{-adv}[r]_i$ represent the average ratings for the i -th item from the advantaged and disadvantaged user groups, $E_{adv}[y]_i$ and $E_{-adv}[y]_i$ represent the average predicted score for the i -th item from advantaged users and disadvantaged user groups. By integrating these fairness terms into the overall learning objective function, fairness metrics can be jointly optimized with recommendation quality metrics. However, some fairness metrics such as equal exposure opportunity or statistical parity are non-differentiable. Thus, several methods are developed to impose indirect regularizers in recommendation models. For example, a fairness-aware tensor-based recommendation method (FATR) is proposed to encourage isolating sensitive features from the original latent factor matrix by adding an orthogonal regularization term [461]. Beutel et al. [26] propose a novel pairwise regularizer for pairwise ranking fairness, which decouples the residual between clicked and unclicked items with clicked item's group membership. The overall loss function can be formulated as:

$$\min_{\theta} \left(\sum_{(\mathbf{q}, j, y, z) \in \mathcal{D}} \mathcal{L}_{rec}(f_{\theta}(\mathbf{q}, \mathbf{v}_j), (y, z)) \right) + |\text{Corr}_{\mathcal{P}}(A, B)|, \quad (15)$$

where \mathbf{q} is the query that contains user feature and context feature, \mathbf{v}_j is the feature vector of item j , $f_{\theta}(\mathbf{q}, \mathbf{v}_j)$ is the prediction for query \mathbf{q} on item j , (y, z) represents the user click feedback and post-click engagement, \mathcal{L}_{rec} is the recommendation model training loss (i.e., squared error loss), \mathcal{P} is a gathered dataset that includes random pairs of relevant items shown to the user and is recorded when the user clicks on one of the items, and A, B represent random variables over pairs from \mathcal{P} .

- **Adversary Learning.** The adversary learning methods are also potential choices for mitigating unfairness in recommender systems. A fair recommender system ought to make equal prediction outcomes toward different sensitive groups. In other words, model predictions should be independent from sensitive attributes. Adversary learning is hereby proposed to learn truly fair representations or predictions, from which the sensitive attribute cannot be inferred. Formally, adversary learning can be formulated as a min-max optimization between the main recommendation model $\mathcal{L}_{rec}(\theta)$ and the adversarial discriminator $\mathcal{L}_{adv}(\Psi)$ for predicting sensitive features as follows:

$$\min_{\theta} \max_{\Psi} \mathcal{L}_{rec}(\theta) + \alpha \mathcal{L}_{adv}(\Psi), \quad (16)$$

where α is the adversarial coefficient parameter that controls the trade-off between recommendation quality and fairness performance. A general architecture to mitigate bias in such kinds methods can be illustrated in Figure 3, where sensitive attribute filters are integrated into the original recommendation model for removing sensitive information, and

discriminators are additionally incorporated to predict corresponding sensitive attributes from filtered representation. The learned feature representation can be viewed as fair until the discriminator fails to predict sensitive attributes from filtered representation.

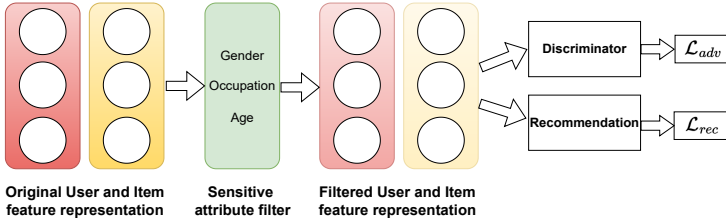


Fig. 3. General framework of fairness-aware adversarial learning in the recommendation field.

In [33], a set of filters are trained to generate graph node embeddings invariant to corresponding sensitive attributes. After training, these filters can be composed flexibly to fulfill different compositional fairness requirements. Wu et al. [380] consider that each user's embedding is composed of ego-centric graph embedding and original embedding learned from the user-item bipartite graph, and propose to utilize discriminators to regularize both embeddings. In this way, node-level fairness and ego-centric fairness can be satisfied simultaneously. Li et al. [221] propose a counterfactual fair recommendation framework to generate sensitive feature-independent user embeddings via adversary learning, which naturally implements individual-level personalized fairness. In [285, 379], they study user and provider fairness in news recommendation with adversarial learning methods. Instead of using sensitive attribute filters, they propose to learn biased user/provider embedding and bias-free user/provider embedding simultaneously and encourage them to be orthogonal to each other.

Adversarial learning can also be used for fairness objectives other than learning fair representations. In [461], a discriminator is designed to predict item-sensitive attributes based on the recommendation model's predicted score, which alternatively helps to learn fair prediction scores. In [215], a discriminator is incorporated to reconstruct user and item textual information from learned feature representations, which forces the representations to preserve their unique properties before being used for rating prediction, thereby reducing the unfairness between mainstream users and non-mainstream users. Li et al. [207] propose a GAN-based model to achieve fair learning without utilizing negative implicit feedback, consisting of a ranker and a controller. Both ranker and controller include a generator and a discriminator. More specifically, the ranker aims to learn user preferences, and its discriminator distinguishes users' real interactions from model-generated interactions. The controller provides fairness signals to the ranker, and its discriminator identifies the generated exposure distribution from the exposure distribution calculated based on the ranker's predictions.

- **Reinforcement Learning.** Since the recommendation feedback loop is both dynamic and sequential, fairness issues in this dynamic process can be modeled as a Markov Decision Process (MDP), which can be addressed by reinforcement learning techniques. In [244], a two-fold reward that measures the system's accuracy and fairness gain is newly designed, facilitating reinforcement learning with the actor-critic architecture. Specifically, time-varying recommendations are performed by an actor network, considering both the system's fairness status and user preferences; the critic network estimates the output of the actor network, which can provide information about whether item groups are over-presented or under-presented. Ge et al. [120] propose a multi-objective reinforcement learning framework, where

the conditioned network can seek the Pareto frontier of fairness and utility, and thereby facilitate decision-makers to control the fairness-utility trade-off. In [122], the dynamic long-term fair recommendation is modeled as a constrained Markov decision process, where the model can dynamically adjust recommendation policy to satisfy fairness requirements when the environment changes, such as the popularity of different item groups due to users' interaction.

- **Causal Graph.** Causal methods have recently been applied to eliminate causal effects between sensitive variables and decisions. In [387], a causal graph is built to identify and remove discrimination in ranked data. Both direct and indirect discrimination would be removed once detected, thereby guaranteeing to reconstruct a fair ranking. In [452], causal inference is applied to solve popularity bias in recommender systems. In order to mitigate popularity bias's effect, a general framework DICE is proposed to disentangle user interests and popularity bias through learning interest embedding and popularity embedding separately. Huang et al. [162] propose to incorporate causal inference into bandits for achieving counterfactual fairness for users in online recommendations. Specifically, they incorporate soft intervention to model the arm selection strategy and adopt the d-separation set identified from the underlying causal graph for developing a fair causal bandit algorithm. Such design can promote fairness by choosing arms that satisfy the fairness constraint. Causality-based methods can also be used to enhance model transparency. In [121], they propose to use counterfactual reasoning to explain which features can cause item exposure unfairness in recommendations.
- **Others.** There are other methods for promoting fairness in recommendations. Borges et al. [31] propose randomness variational autoencoders, which incorporate randomness into the regular operation to alleviate the position bias in multiple-round recommendation. The experimental results indicate that adding noise to VAE in the latent representation sampling process can promote long-term fairness in recommendations with a tolerable trade-off between recommendation quality and fairness. Wu et al. [377] empirically demonstrate that big recommendation models elicit unfairness issues to cold users. Moreover, they propose a self-distillation framework called BigFair, where model predictions on original user data serve as a teacher to regularize predictions on augmented user data generated by randomly dropping historical behaviors. Experimental results indicate that BigFair can encourage big recommendation models to improve cold-start users' performance and achieve better fairness.

3.2.3 Post-processing Methods. The post-processing methods aims to promote fairness based on recommendation models' output, which are primarily in the form of re-ranking. Existing re-ranking methods can be categorized into the subsequent three categories:

- **Slot-wise Re-ranking.** Slot-wise re-ranking methods are typically implemented by adding items sequentially to empty slots of a recommendation list according to specific rules or re-ranking scores. In [306], to promote fairness in the package-to-group recommendation, a greedy algorithm is proposed to add items to the package gradually by considering item category and distance constraints. In order to balance the relevance of items to group members for each prefix of the top-N, Kaya et al. [189] present a new and rank-sensitive definition of fairness (GFAR) for top-N group recommendations, in which a greedy algorithm is designed to find top-N group recommendations according to the GFAR definition. The work of [305] propose an efficient method that enumerates all fair packages, which allows users to select their favorite packages in their own way. Liu et al. [243] propose a personalized re-ranking algorithm to achieve a fair microlending recommendation system, where the objective is formulated as a combination of personalization score $P(v | u)$ and a fairness term, together

with a trade-off controlling hyper-parameter λ as follows:

$$\max_{v \in R(u)} \underbrace{(1 - \lambda)P(v | u)}_{\text{personalization}} + \lambda \underbrace{\sum_c P(\mathcal{V}_c) \mathbb{1}_{\{v \in \mathcal{V}_c\}} \prod_{i \in S(u)} \mathbb{1}_{\{i \notin \mathcal{V}_c\}}}_{\text{fairness}}, \quad (17)$$

where $R(u)$ denotes the initial ranking list, \mathcal{V}_c represents a group of loans with attribute c , $P(\mathcal{V}_c)$ represents the importance of \mathcal{V}_c , $\prod_{i \in S(u)} \mathbb{1}_{\{i \notin \mathcal{V}_c\}}$ indicates the coverage of \mathcal{V}_c for the current generated re-ranked list $S(u)$.

- **User-wise Re-ranking.** User-wise re-ranking aims to find the most appropriate recommendation list for each user guided by the overall optimization objective. In [28], to implement amortized fairness, which refers to attention being fairly distributed across a series of rankings, the optimization is formalized as an integer linear programming problem and solved by an efficient heuristic algorithm Gurobi. Besides, Mehrotra et al. [253] combine fairness with recommendation utility by adopting an interpolation strategy and a probabilistic strategy to generate the candidate recommendation lists.
- **Global-wise Re-ranking.** Global-wise re-ranking aims to re-rank several recommendation lists at the same time. In [219], an integer programming-based approach is proposed for solving user unfairness problems in commercial recommendations, and the overall optimization objective is defined as follows:

$$\begin{aligned} \max_{\mathbf{W}_{ij}} \quad & \sum_{i=1}^n \sum_{j=1}^N \mathbf{W}_{ij} Y_{i,j} \\ \text{s.t.} \quad & \text{UGF}(Z_1, Z_2, \mathbf{W}) < \varepsilon \\ & \sum_{j=1}^N \mathbf{W}_{ij} = K, \mathbf{W}_{ij} \in \{0, 1\}, \end{aligned} \quad (18)$$

where \mathbf{W}_{ij} represents the binary variable that indicates whether item j is recommended to user i , $Y_{i,j}$ represents the preference score of user i to item j , Z_1 and Z_2 refer to the advantaged and disadvantaged user groups, UGF represents user unfairness evaluation metric, K denotes the total length of the recommendation list and ε is a hyper-parameter that controls the strictness of fairness requirements.

3.3 Applications

When recommender systems are deployed in the large-scale and resource-allocated platform, unfairness becomes a severe threat to the platform's trustworthiness. In the following subsection, we illustrate the potential unfairness in two real-world applications from different domains, which featured the necessity of building fair recommender systems.

- **Job Recommendation.** Job recommendation systems are widely deployed in the job recruitment market, such as LinkedIn [124] and Indeed [301]. These systems are expected to guarantee fair opportunities for all qualified candidate users, since the recommendations can be viewed as social resource allocations. If recommendation models fail to meet fairness requirements, over/under-representation of specific user groups or racial/gender stereotypes in the recommendation results will inadvertently occur in practice, raising severe legal and societal issues. Therefore, building fair job recommendation systems is vital to both systems' own benefits and societal benefits.
- **E-commercial Recommendation.** E-commercial recommendation systems, i.e., Amazon, Taobao, etc., are prevalent for their effectiveness in connecting consumers and the relevant products. Users' satisfaction and platforms' interests highly depend on the quality of the generated recommendation results. However, previous studies [219] indicate that most users are disregarded by commercial recommendation engines when we categorize users into

groups according to their different activity levels. Such issues also exist on the provider side [265]. Mitigating unfairness issues on both sides are essential to commercial recommender systems' long-term benefits.

3.4 Surveys and Tools

In this subsection, we sort out the existing surveys and tools on fairness in recommender systems to facilitate researchers in this field.

3.4.1 Surveys. There have been growing concerns regarding fairness in recommender systems in recent years. Chen et al. [49] gives a detailed summary of bias existing in recommender systems, provides a comprehensive taxonomy to organize current recommendation debiasing works, and discuss the strengths and weaknesses of different debiasing methods. In [420], they propose a survey that connects related approaches across various fields, aiming to motivate fairness-enhancing interventions in ranking. Pitoura et al. [283] provide a more technical view of definitions and methods used to guarantee fairness in rankings and recommendations, which has a much broader content coverage and structured content comparison. Li et al. [220] present a comprehensive survey of the foundations for fairness. The content ranges from fairness in general machine learning tasks to more sophisticated ranking and recommender systems. In [365], they mainly focus on fairness in recommender systems.

3.4.2 Tools. Various toolkits have been developed to evaluate or mitigate bias in machine learning models. IBM Fairness 360 [23] is a comprehensive tool that provides more than 70 metrics for quantifying individual or group fairness and nine bias mitigation algorithms. Furthermore, it enables metric explanations to help users understand the fairness evaluation results. Google What-if [369] tool allows users to test model performance in a what-if way, where users can edit the values of data points and see their effects on model performance. In addition, the What-If tool supports comparing multiple models in the same workflow and testing several algorithmic fairness constraints. Fairkit-Learn [176] is an interactive python toolkit that supports developers to reason about and understand model fairness. By comparing different machine learning models concerning quality and fairness metrics, Fairkit-learn can find a model that is both high-quality and fair models. However, evaluation and analysis toolkits for recommender systems are still blank, which motivates further development.

3.5 Future directions

Fairness has attracted intensive attention for achieving trustworthy recommender systems. However, many essential open problems and challenges are still not well addressed. In this subsection, we discuss some critical future directions.

- **Consensus on fairness definitions.** The fairness definition usually varies in different application scenarios, and the biases that lead to unfairness are also typically multi-sourced. Then the fairness demands can arise from multiple sides (i.e., user/provider side) and different perspectives. Some fairness definitions may even conflict with each other [192], such as the calibrated fairness and Rawlsian max-min fairness. Given these situations, it is vital to achieving consensus on fairness definitions. Specifically, there are three key challenges. First, to enhance a recommender system's fairness, how to determine the priority of multiple fairness objectives and make an appropriate balance if there is a conflict? Second, how to determine the most suitable fairness metric for a specific scenario? Third, how to simultaneously incorporate multiple fairness notions into one general framework for achieving the fairness for trustworthy recommender systems? To address these questions, it is highly

desired to evaluate different fairness metrics on benchmark datasets and achieve a consensus on their relationships from a unified view.

- **Fairness-aware algorithm design.** There have been extensive studies conducted on improving fairness of recommender systems. Nevertheless, it is still unclear whether the existing model design implicitly inherits or induces bias. For example, graph neural networks have shown great potential in recommender systems [104, 383], a recent study indicates that given a biased graph topology as input, the information propagation mechanism of graph neural networks may induce bias to the node embeddings. Besides, recent studies have shown that fairness and causality are closely related, and causality-based methods for mitigating bias have become a new trend [269]. Since traditional recommendation models tend to capture spurious associations during collaborative filtering, sensitive features may be encoded into feature representation even if not explicitly used. To solve this problem, causality-based methods for unbiased recommendation are worth exploring to address the model-induced bias.
- **Trade-off between fairness and utility.** Previous studies on fairness in different fields have revealed a trade-off between fairness and utility. While in the industrial recommender system [74], overall performance degradation is unacceptable due to the revenue loss. Therefore, extensive research needs to be conducted to figure out the trade-off mechanism so that the decision-makers can make a better balance.

4 EXPLAINABILITY

Recommendation with explainability, or to say explainable recommendations, refers to the recommendation algorithms focusing on providing interpretation for recommendation results. In fact, it represents the intersection of explainable AI and recommendation algorithms for enhancing the trustworthiness of recommender systems. In recent years, with the introduction of many black-box modules such as Multilayer Perceptron (MLP) [213, 455], Transformer [330, 382] and Reinforcement Learning (RL) [166, 449] in recommender systems, the working mechanisms of advanced recommender systems are obscure without explainability. This problem makes the models hard to be fully trusted and to put into safety-critical applications [25, 268]. Therefore, building a trustworthy recommender system requires explainable recommendation modules.

This section will first introduce the concepts and taxonomy of explainable recommendations. Then, we provide detailed descriptions of some representative methods. Finally, we provide discussions on some applications and open topics in this direction.

4.1 Concepts and Taxonomy

In this subsection, we first introduce the basic concepts about explainable recommendations and then gives taxonomies of research on explainable recommendations.

4.1.1 Concepts. The *explainability* of a system can be described as *the ability to explain or to present in understandable terms to a human* [88]. When such explainability exists in a recommendation model, this framework is called an explainable recommendations model. Explainable recommendations not only provide users with their recommendations but also give reasons why to recommend them. Although the research of explainable recommendations can be dated back to the 2000s [29, 150], the formal analysis [436] and the wide attention from the academia [435] on it have just started recently.

4.1.2 Taxonomy. In this subsection, we introduce the taxonomy of research on explainable recommendations. There are two main parts: the taxonomy for models and the taxonomy for evaluations.

Table 3. Taxonomy of existing explainable recommendations methods and some representative studies of different aspects.

	Model-intrinsic based	Post-Hoc	Characteristics
Structured	[48, 114, 364, 389, 390, 396]	[280, 319]	Logical, Visible
Unstructured	[63, 64, 291]	[211, 315, 338]	Diversified, Fragmented
Focus	Model's reasoning process	Instances' relationship	-

Taxonomy for models. Explainable recommendation models can be classified according to the following two criteria.

- *How to produce explanations: model-intrinsic based or post-hoc.* If a technique seeks to derive explanations from the intrinsic structure of the model, it is called model-intrinsic based explanation. Most of the current studies about explainable recommendations fall into this category, and these techniques are usually important components of relevant recommendation models. [180]. Depending on the model structure, these techniques often pay more attention to the reasoning process of the model. In contrast, another class of techniques called post-hoc methods [353] provide explanations based only on the inputs, outputs and extrinsic conditions of the model [349]. In this case, the recommender system is treated as a completely black-box model [171]. Note that post-hoc methods in recommender systems are more flexible than model-intrinsic methods, since post-hoc methods are model-agnostic and can be applied to any recommendation methods. In addition, such post-hoc explanation methods are widely utilized to explain deep recommender systems (with millions of parameters) which are too complicated to be understood.
- *How the explanations are presented: structured or unstructured.* Structured methods present explanations in the form of logical reasoning based on some particular structures, such as a graph [278], or a knowledge graph [460]. The explanations they provide are highly organized and generally characterized by strong logic, good visualization, and comprehensive reasons. Unstructured methods, on the other hand, do not rely on, or explicitly rely on logical reasoning to give explanations. They tend to generate easy-to-understand sentences, ratings, or features directly from a black-box model [310]. The explanations they provide are often fragmented and structurally uncertain. These two types are not completely opposite. Generally speaking, structured methods focus more on explicit associations between users and items, while unstructured methods focus more on implicit expressions such as emotions and overall evaluations.

Table 3 lists the different categories of explainable recommendations models with some recent representative studies. It also summarizes some common characteristics and focuses of different categories. A detailed description of explainable recommendations model techniques is given in subsection 4.2.

Taxonomy for evaluations. Evaluations of explainable recommendations can be classified according to two main criteria. It is worth noting that most of the current model-intrinsic based explainable models consider the characteristic of explainability and the improvement of the recommendation effect. These two aspects are usually evaluated together. In this case, the improvement of the model itself sometimes can also serve as an evaluation of the explainability.

- *Evaluation perspectives: Effectiveness, Transparency and Scrutability.* Explainable recommendations can be used for different purposes to produce different effects. Therefore, when evaluating explainable recommendations, different perspectives should be taken into full

Table 4. Summary of evaluation perspectives

Evaluation perspective	Evaluation criteria	Related research
Effectiveness	Whether the explanations are useful to users? (e.g. Decision making, Recommendation results)	[8, 58, 337]
Transparency	Whether the explanations can reveal the working principles of the model?	[18, 144, 225]
Scrutability	Whether the explanations contribute to the prediction of the model?	[327, 347, 362]

Table 5. Taxonomy of evaluation forms

Evaluation form	Corresponding perspectives	Related research
Quantitative metrics	Effectiveness; Scrutability	[337, 338]
Case study	Effectiveness; Transparency	[225, 362, 396]
Real-world performance	Effectiveness; Scrutability; Transparency	[58, 347, 392]
Ablation Study	Effectiveness; Transparency	[64, 211, 327]

consideration. Tintarev et al. [344, 345] and Balog et al. [17] propose seven useful perspectives in the previous research². According to our summary of previous papers, we summarize that three of them: Effectiveness, Transparency and Scrutability, are the most representative perspectives that are frequently considered. The characteristics and representative papers of these perspectives have been summarized in Table 4.

- *Evaluation form: Quantitative metrics, Case study, Real-world performance, and Ablation Study.* There are several evaluation forms available for explainable recommendations. According to previous studies, most of them can be divided into the following four forms: Quantitative metrics, Case study, Real-world performance, and Ablation Study. Quantitative metrics focus on the quantification of explainability in mathematics. Case study focuses on examining whether the explanation conforms to human logic. Real-world performance aims at examining the practical effects of the explanation. And ablation study intends to illustrate how algorithmic modules provide explanations and how these modules enhance the recommendation model. A detailed discussion of evaluation techniques is given in subsection 4.3. Research summaries of various forms of evaluations and their most corresponding perspectives are shown in Table 5.

4.2 Methods

In this subsection, we introduce some representative methods of explainable recommendations models according to two taxonomies in Table 3 and summarize some representative characteristics and focuses of these categories. These characteristics and focuses are also shown in Table 3.

4.2.1 Different methods of producing explanations: model-intrinsic based and post-hoc. This category is mainly defined by the underlying principle of the explainable recommendations. Based on this

²i.e. Transparency, Scrutability, Trust, Effectiveness, Persuasiveness, Efficiency and Satisfaction

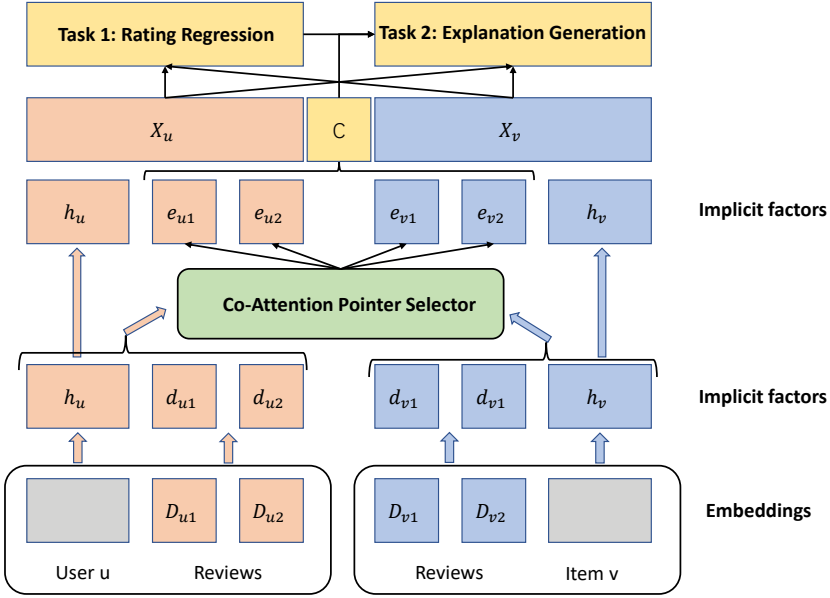


Fig. 4. An example (CAML) of the model-intrinsic based methods. In such explainable recommendations, explanations are often generated along with recommendations by different tasks at the end of the model.

perspective, the explainable recommendations methods can be grouped into model-intrinsic based methods and post-hoc methods.

- *Model-intrinsic based Methods.* Methods in this category are also called model-intrinsic methods. These methods embed explicable ability into the recommendation model and provide explanations for the model that they are attached to by generating explanatory graphs, paths, parameters and other contents in the process of recommendation. In general, such methods are only effective for embedded models and cannot simply be reused in the recommendation process of other models. For example, in Co-Attentive Multi-task Learning (CAML) [64], Chen et al. design an encoder-selector-decoder architecture for multi-task learning to realize a model that considers recommendation and explanation at the same time using transferred cross knowledge. As shown in Figure 4, they first encode the information and reviews of users and items into implicit factors (h for information, d and D for reviews, u for users, v for items), then they feed the review implicit factors d to a selector and generate new factors represented by e . By using different e blocks and h blocks in the following layers, the model can use partially shared knowledge C and partially proprietary knowledge X for multi-task prediction of explanation generation and rating prediction.

Different from multi-task learning, Aspect-aware Latent Factor Model (ALFM) [69] and Multi-Modal based Aspect-aware Latent Factor Model (MMALFM) [68] perform explainable rating prediction by modeling aspect importance to generate aspect ratings in the latent factor model. Specifically, the user-item rating modeled according to the latent factor algorithm can be defined as:

$$\hat{r}_{u,i} = \sum_{a \in \mathcal{A}} \overbrace{\rho_{u,i,a}}^{\text{aspect importance}} \underbrace{r_{u,i,a}}_{\text{aspect rating}}, \quad (19)$$

where $\rho_{u,i,a}$ represents the aspect importance and $r_{u,i,a}$ represents the aspect rating for a user u , an item i and a related aspect a in aspect set \mathcal{A} . ALFM proposes Aspect-aware Topic Model (ATM) to model the relationship between the user, item, aspect and latent topic through probability distribution, and calculates the parameters needed by $\rho_{u,i,a}$ and $r_{u,i,a}$ according to the corpus composed of text words by ATM method. MMALFM further proposes Multi-modal Aspect-aware Topic Model (MATM) to add the input and modeling of visual words for parameter estimation. Finally, the acquired parameters related to aspects such as $\rho_{u,i,a}$ can explain the user's inclination towards objects in different aspects, and the top words about a certain aspect of a user or an item can also explain the user's inclinations or object characteristics.

Another example comes from Neural-Symbolic explainable recommendations (NSER) [389]. In this paper, the author introduces Knowledge Graph (KG) to their model and uses neural symbolic reasoning methods to narrow down the explainable path search area in KG and then generate coarse-to-fine explanations for their recommendation results. Moreover, Temporal Meta-path Guided explainable recommendations (TMER) [48] uses the dynamic KG rather than static KG technique to model the Temporal Meta-path Guided explainable recommendations. Specifically, it uses the multi-head attention module to learn the combinational features from multiple path instances:

$$\text{Attention}(Q_\phi, K_\phi, V_\phi) = \text{Softmax}\left(\frac{Q_\phi K_\phi^T}{\sqrt{d_k}}\right) V_\phi, \quad (20)$$

$$\text{MultiHead}(Q_\phi, K_\phi, V_\phi) = \text{Concat}(\text{head}_1, \dots, \text{head}_m) W,$$

where W is the weight, head_i is the Attention module with dimensionality d_k of Query Q , key K and value V related to path ϕ . After that, the model captures sequential dependencies through time-ordered links for recommendation. Finally, it will be able to provide explanations for recommendations by aggregating multiple item-item instance paths generated by user history.

Although the techniques used by the above methods and other methods in this class may differ greatly, the explanations they provide commonly rely on the recommendation process and serve the recommendation models. Therefore, they mainly focus on the reasoning process of the recommendation models. This characteristic often helps these methods produce more detailed explanations, but also makes them difficult to migrate to different recommendation models.

- *Post-hoc Methods.* Methods in this classification are also called model-agnostic methods. Such methods regard the recommendation model as a black-box or do not even consider the recommendation model. Instead, they only use the known input and prediction information of the recommender system to train a new model similar to the system for explanation generation, or to directly generate explanations for predictions or features viewed by users. An example process of these methods is shown in Figure 5. Such approaches are more general than model-intrinsic based approaches, since in most cases they require only the specific form of input and prediction data. Note that because of their versatility, these methods can also be applied to model-intrinsic based explainable recommendations models.

As an example, Singh and Anand [319] approximate the black-box ranker model used in the web search area by modeling an explainable tree-based model to obtain explanations of instances. Similarly, Shmaryahu et al. [315] provide explanations by approximating complex algorithms using a set of simple explainable recommendations algorithms. To get better explanations, Peake and Wang [280] treat a matrix factorization recommendation model in

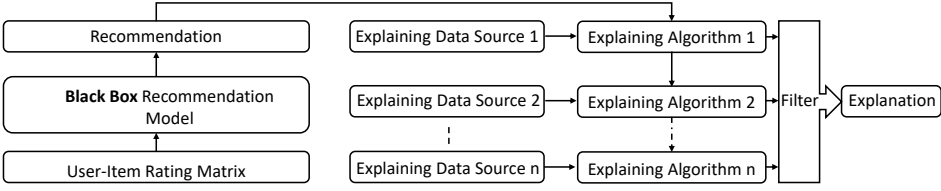


Fig. 5. An example ([315]) process of the post-hoc methods. In this work, the original recommendation model is treated as a black-box model. And their explainable model can generation explanations for recommendations from any recommendation model.

their paper as a black-box. Their model takes the rating prediction matrix of the recommendation model as input and tries to extract association rules for explanation generation while approximating the black-box model.

Another example comes from Ai et al. [7], where they propose to provide personalized recommendation using KG embedding-enhanced Collaborative Filtering (CF) methods, and try to generate explanations using similarity matching between the user and item embeddings. Mathematically, they generate explanation paths with probability P for an arbitrary user e_u and item e_i with relation sets R_α, R_β and the aggregation operation represented by $trans(\cdot)$ through any intermediate e_x , as follows:

$$P(e_x | e_u, R_\alpha, e_i, R_\beta) = P(e_x | trans(e_u, R_\alpha)) P(e_x | trans(e_i, R_\beta)). \quad (21)$$

After that, they choose the best path with top probability for final explanation. Since the explanation is independent of the model reasoning process, this approach is also considered as a post-hoc approach.

In general, most of these methods obtain explanations by approximating simple explainable models to complex models, or by summarizing inputs, outputs and visible features using different techniques. Since they are mostly independent of specific recommendation models, they can generally generate explanations for different recommendation algorithms. However, since they cannot fully obtain the internal reasoning logic of the recommender system, most of the explanations generated by post-hoc methods rely on the association between instances, and the explanation effect is commonly weaker than model-intrinsic based methods.

4.2.2 Different methods of presenting explanations: structured and unstructured. This category is mainly defined by the presentation of explanation. Based on this perspective, explainable recommendations methods can be divided into structured methods and unstructured methods.

- **Structured Methods.** These approaches mainly provide structured explanations with strong logic, and the explanation logic is often visible. To generate structured explanations, these methods usually use structured data as the main part of their input. The most representative technique of this category is explainable recommendations based on knowledge graph [363, 364, 390, 440], and the explanations provided by this kind of explainable recommendations method mainly include explanation path graph [364, 390, 440], association rules [280], decision tree [319], etc.

One of the most representative works in this category is Policy-Guided Path Reasoning (PGPR) [390] whose structure is shown in Figure 6. In this model, a path p_k with k entities in KG is defined as $p_k(e_0, e_k) = \{e_0 \xleftrightarrow{r_1} e_1 \xleftrightarrow{r_2} \dots \xleftrightarrow{r_k} e_k\}$, where e_i represents entity i and r_i represents relation i . This model simulates path reasoning in KG with Markov Decision

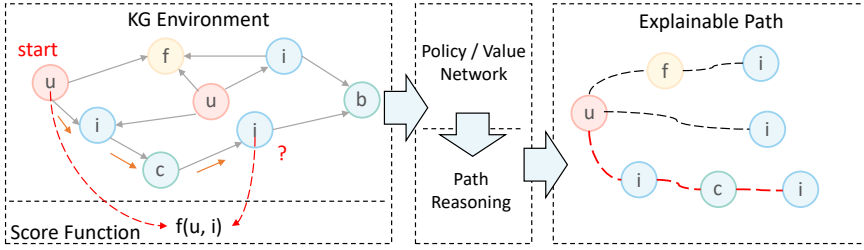


Fig. 6. An example (PGPR) of the reasoning path for structured methods. u for users, i for items, b for brands, c for categories, and f for features in the KG environment.

Process (MDP) with reward function as follows:

$$R_T = \begin{cases} \max \left(0, \frac{f(u, e_T)}{\max_{i \in I} f(u, i)} \right), & \text{if } e_T \in I \\ 0, & \text{otherwise} \end{cases} \quad (22)$$

where $f(\cdot)$ is a scoring function. u is a user, i is an item of item set I , and e_T is the terminal entity. Then, REINFORCE algorithm [373] is used to learn a path finding policy $\pi(\cdot | s, \tilde{A}_u)$, where \tilde{A}_u is a binarized vector of pruned action space for user u and s is the current state. Finally, it can couple recommendation and explainability by providing paths in KG using beam search.

Another representative work is Knowledge-aware Path Recurrent Network (KPRN) [364]. In this study, the path can also be represented as $p_k(e_0, e_k) = \{e_0 \xrightarrow{r_1} e_1 \xrightarrow{r_2} \dots \xrightarrow{r_k} e_k\}$. Differently, KPRN aggregates entity, entity type, and relation type pairs for each possible path and enters them sequentially into a Long Short-Term Memory (LSTM) model to predict a score for each path. Ultimately, these paths provide explanations for the recommendation based on their contribution scores.

Moreover, in ADversarial Actor-Critic (ADAC) [440], to model imperfect paths and relations in KG more efficiently and improve the persuasion of path explainability, an actor-critic model is used to search for persuasive paths for final explanation. In the Reinforcement learning framework for Multi-level recommendation Reasoning (ReMR) [363], to solve the lack of explainability of high-level abstract categories in previous knowledge graph reasoning processes, Wang et al. use their Cascading Actor-Critic module-based multi-level reasoning over KGs to better infer and represent user interests. Different from using KG to generate explanations during the recommendation process, Singh and Anand [319] try to train a structured tree-based model for post-hoc explanations of learning-to-rank-models, which are widely used in most information retrieval and recommender systems.

In general, the explanations provided by structured methods usually tend to have inference paths or rules, which makes them highly logical. This feature also partly limits the use scenarios of structured methods, because in most cases they need a large number of structured data with many features or connections between users and items to provide the information that they need for reasoning.

- *Unstructured Methods.* Unlike structured approaches, unstructured methods do not require explanations to be organized, and thus have few restrictions on how explanations can be presented, as long as they are intuitive to human beings. Because of the loose restrictions on interpretive logic, these approaches can focus on different aspects of a problem and their

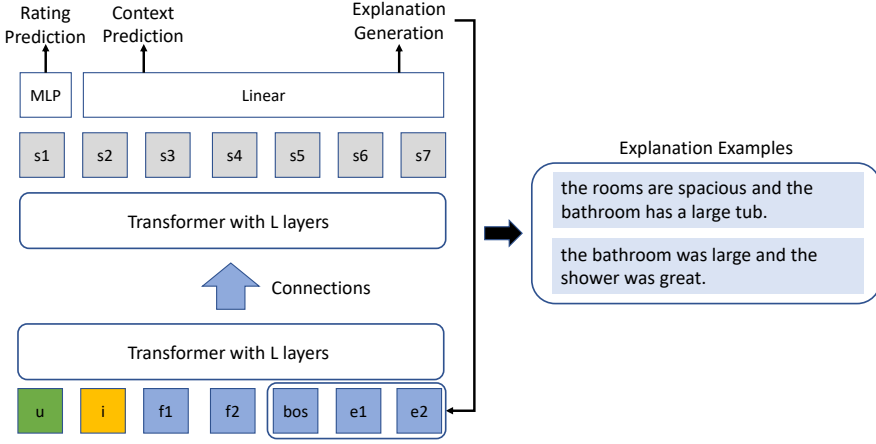


Fig. 7. An example (PETER) of sentence generation for unstructured methods. u for users, i for items, f_1, f_2 for features, bos, e_1 and e_2 for explanation words where bos is the “start” representation.

inputs are diversified. The most representative technique of unstructured methods is the sentence generation model based on RNNs or Transformer [64, 211], and the explanations provided by this kind of method mainly include sentences, scores, important features, etc. An example of sentence generation for unstructured methods (PETER [211]) is illustrated in Figure 7.

One work using unstructured method comes from Explainable Conversational Recommendation (ECR) [63] for explainability in conversational recommendation. This model uses incremental multitask learning and conversation feedback to improve the effectiveness of both recommendation and explanation results from conversational recommendations. The pre-trained loss function constructed in this model is designed as:

$$\mathcal{L}^{\Omega} = \sum_{u \in U} \left(\mathcal{L}_r^{\Omega u} + \lambda_n \mathcal{L}_n^{\Omega u} + \lambda_c \mathcal{L}_c^{\Omega u} \right) + \lambda_{\theta} \|\Theta\|_2^2 \quad (23)$$

where $\mathcal{L}_r^{\Omega u}$ is Factorization Machine (FM) [292] based BPR loss [293], $\mathcal{L}_n^{\Omega u}$ is the negative log-likelihood loss of two Gated Recurrent Unit models (GRUs), $\mathcal{L}_c^{\Omega u}$ is the concept relevance loss [64], and Θ represents the model parameters. λ represents the weight. In this loss function, $\mathcal{L}_n^{\Omega u}$ and $\mathcal{L}_c^{\Omega u}$ are highly related to the generation task of explanation. Finally, the model can make reliable explainable recommendations in the form of conversations. Another example comes from PErsonalized Transformer for explainable recommendations (PETER) [211], which also uses multitask learning for explainable recommendations. Differently, this model mainly uses Transformer as its core module and follows a linear layer to generate personalized explanations for different IDs.

Different from explanations from sentences, Counterfactual explainable recommendations (CountER) [338] attempts to construct counterfactual items for recommended items to provide explanations. Specifically, this model tries to use small changes in item aspects to reverse

the decision. These small changes in item aspects constitute the explanation of the recommended item. The mathematical representation of the method can be written as follows: *Minimize Explanation Complexity, s.t., Explanation Strength is higher enough*. Here Explanation Complexity and Strength should be predefined according to the properties of these concepts.

Generally, the explanations provided by this kind of method are mostly applied in the directions of comments, conversations, sentence generation and other fragmented explanation generations. Their forms are diversified, and the content is often fragmented. Therefore, such models are suitable for a wide variety of data. However, since they usually focus on a few aspects of the data and lack clear reasoning logic, the validity of the explanations they generate depends more on human experts' intuitive evaluation.

4.3 Methods of Evaluations

In this subsection, we introduce some representative methods of explainability evaluations in recommendations, which can be summarized in Table 4 and Table 5.

4.3.1 Quantitative Metrics. To better evaluate the explanation, many quantitative metrics [210, 210, 231, 240, 276] have been designed to numerically approximate the general evaluation of some aspects of the explanations [181, 214, 233]. Most of these evaluations are designed for natural language generation (NLG) tasks [59]. Since sentences can be split into words, it is convenient for the metrics to perform mathematical association matching and numerical accumulation operations. The most common quantitative metrics are ROUGE score [231] and BLEU metric [276] as follows:

- ROUGE score:

$$\text{Rouge-N} = \frac{\sum_{S \in \text{Reference Summaries}} \sum_{gram_n \in S} \text{Count}_{match}(gram_n)}{\sum_{S \in \text{Reference Summaries}} \sum_{gram_n \in S} \text{Count}(gram_n)}, \quad (24)$$

where n represents the length of the n-gram, $gram_n$ is the maximum number of n-grams co-occurring in a candidate summary and $\text{Count}_{match}(gram_n)$ is a set of reference summaries.

- BLEU metric:

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^N \frac{1}{N} \log p_n\right), \quad (25)$$

where BP represents the brevity penalty. p_n represents the average value of the modified n-gram precision. N represents the length of a sentence.

In addition, there are also some other quantitative matrices such as Unique Sentence Ratio (USR) [210], Feature Matching Ratio (FMR) [210], shift scores [240], etc.

4.3.2 Case Study. This approach provides examples to illustrate explanations and even part of the explanation generation process [54, 57, 391], so that people can intuitively judge their validity. Specifically, researchers usually lists one or several examples of explanations, and summarizes the advantages of the proposed explanations by comparing them with each other or with explanations generated by other models. For most structured explanations, researchers can further verify the logic of the explanation by showing the detailed features of the explanation or the path of reasoning. Since explanation is oriented towards people's actual judgment, which includes logic, experience, intuition and other factors, providing concrete examples for people can help improve the credibility of an explanation.

4.3.3 Real-world Performance. This approach focuses on collecting people's actual feedback or carrying out online experiments for evaluation. This approach is usually evaluated by running the model online and gathering feedback, or by recruiting annotators to carefully evaluate the

explanations [57, 392]. This approach focuses on the practical effects of the explanation, i.e., whether the explanation has a positive impact on the recommendation in the application. Moreover, the different evaluation details will further reflect different perspectives of the explanation.

4.3.4 Ablation Study. This approach is often used to analyze the role of different modules in recommender systems [64, 211]. Specifically, the researchers verify the impact of important modules on the overall model by removing some modules or drastically adjusting some parameters. In explainable recommendations, this approach can analyze the specific contributions that explainable modules bring to the model, or specific modules that contribute most to generating the explanation.

4.4 Applications

The explainable recommendations have been widely applied in various scenarios in our daily lives.

- **E-commercial Recommendation.** The main purpose of explainable recommendations in various e-commerce platforms is to recommend different products to each customer for justifiable reasons, so as to help improve the attractiveness of the recommended products to customers and thus to stimulate consumer confidence. For example, Zhang et al. [436] use a big commercial e-commerce platform JD.com to evaluate the effect of their explainable recommendations in real situations. In [57], Chen et al. combine visual and textual features to provide visual explanations for fashion recommendations on platforms including Amazon.
- **Social Media.** In social media such as Facebook and TripAdvisor, recommendations with certain explanations can be used to strengthen people's understanding of the recommended content, serving people's life and friendship. For example, explainable Point-of-Interest (POI) recommendations are designed to discover places of interest and provide explanations for people [362, 439]. Zhao et al. [439] propose a joint sentiment-aspect-region model based on the Yelp dataset. This model could investigate whether people like a certain aspect of a place. Wang et al. [362] propose a tree-enhanced embedding method for transparent and explainable recommendations in tourist attractions and restaurant recommendations.

4.5 Surveys and Tools

In this subsection, we describe the existing surveys on explainability in recommender systems and tools on explainability in AI to facilitate researchers in this field.

4.5.1 Surveys. There have been growing concerns regarding explainability in the recommender system in recent years. Zhang et al. [435] provide a detailed survey of explainable recommendations, and distinguish explainable recommendation models with different algorithmic techniques and explanatory forms. Tintare and Masthoff et al. [344, 345] provide seven useful perspectives for evaluating explanations in the recommender system. Balog and Radlinski [17] discuss the relevance of the seven perspectives and develop two novel explainable evaluation metrics. Moreover, Chen et al. [59] focus on the evaluation part of explainability in the recommender system, and provide main evaluation perspectives for different evaluation methods.

4.5.2 Tools. In this section, we provide commonly-used tools for explainable models. AIX360 [12] is a useful open-source toolkit for explainable models and evaluation metrics in Python environment. In addition, Quantus [148] provides guidance for the evaluation of explainable methods and a comprehensive set of evaluation metrics. For deployment, XAITK [156] provides an open-source collection of explainable AI tools and resources to meet the critical needs of deploying and testing explainable AI systems. It is worth noting that the above tools are mostly oriented to general explainable AI models and methods. Currently, the toolkit for explainable recommendations is still limited, which is also one of the future research directions.

4.6 Future Directions

In this subsection, we discuss future directions for explainable recommendations.

- **Natural Language Generation for Explanation.** Most existing explainable recommendations aim to generate sentence explanations based on predefined templates. However, a more user-friendly explanation could be a natural sentence that is automatically generated. Recently, there have been several studies that have attempted to provide explanations using natural language generation methods [64, 211]. Although they have achieved considerable results, there is still a gap between them and the ideal detailed and natural language explanation, which is one of the important directions of explainable recommendations in the future.
- **Explainable recommendations in more fields.** In addition to methodology, another potential future direction for explainable recommendations is development in more fields, such as medical care and education. At present, explainable recommendations are concentrated mainly in e-commerce and social media, while relevant research in academic support, medical care, education and other fields is still limited. Fortunately, some researchers have noticed the potential value of explainable recommendations in these fields [174, 337]. In the future, the promotion and development of explainable recommendations in these fields will better benefit people's lives.

5 PRIVACY

Human beings have entered the era of big data, in which data can effectively characterize users' profiles via online behaviors (e.g., browsing history and online shopping) and inevitably contain users' private information (e.g., email addresses and gender) [161, 168]. Modern recommender systems, especially deep learning-based strategies, heavily rely on big data and even private data to train algorithms for obtaining high-quality recommendation performance [6, 236]. This raises huge concerns about the safety of private and sensitive data when recommendation algorithms are applied to safety-critical tasks such as finance and healthcare. To build trustworthy recommender systems, protecting data privacy has become increasingly important. Therefore, it is necessary to investigate how to perform privacy attack methods to steal knowledge from the target recommender systems, and then develop privacy-preserving countermeasures to protect data privacy.

In this section, we will focus on privacy attacks and their corresponding strategies with regard to privacy protection for the trustworthiness of recommender systems. We first give some brief concepts and the taxonomy of privacy attacks and privacy protection. Then, representative methods of attacks and privacy-preserving methods on recommender systems are detailed, followed by some surveys and tools related to the privacy of recommender systems. At last, we introduce some real-world applications and describe future directions to explore in achieving trustworthy recommender systems.

5.1 Concepts and Taxonomy

This subsection briefly introduces the widely-received concepts in trustworthy recommender systems, specifically focusing on privacy attacks and privacy-preserving.

5.1.1 Privacy Attacks. In recommender systems, privacy attacks aim to steal knowledge that is not intended to be shared, such as the sensitive information of users and model parameters. Depending on how much knowledge attackers know about target victim recommender systems, privacy attack methods can be classified into white-box attacks, black-box attacks, and grey-box attacks. More specifically, in the white-box setting, attackers are allowed to access all the information about recommender systems, such as the model's architecture, parameters, gradients, training data, etc.

In contrast to the white-box attack, attackers in the black-box setting could have access to minimal knowledge about the victim model. Grey-box attacks, a combination of white-box and black-box attacks, are able to access partial knowledge about the target recommender system, such as users' public reviews and attributes on Amazon.

Moreover, the attack methods can also be divided into four categories based on the information stolen by the attacker: membership inference attacks, property inference attacks, reconstruction attacks, and model extraction attacks.

- **Membership Inference Attacks (MIA)** aim to identify whether the target user is used to *train* the target recommender system. When certain training data is known to the attacker, it can result in a privacy breach. For example, an intelligent system in the healthcare domain recommends treatments for patients with schizophrenia. If an attacker knows a certain person in the training set for building the intelligent system, it is likely to infer that this person is suffering from schizophrenia [316], where the patient's information is sensitive and may not intend to be published.
- **Property Inference Attacks (PIA)** aim at stealing global properties of the training data in the target recommender system. These properties are usually not directly reflected in any specific user but in the global features of the training set, such as the gender ratio and the total amount of Tiktok users [147].
- **Reconstruction Attacks (RA)**, or Model Inversion Attacks, aim to infer private information or labels on training data. Unlike property inference attacks, RA focuses on the sensitive properties of users, such as a user's gender or occupation in Facebook ³.
- **Model Extraction Attacks (MEA)**, also known as Model Stealing Attacks [434], aims to steal the parameters and structure of a target model by creating a new replacement model that behaves similarly to the target model [297]. After that, a successful model extraction can transform the setting into a more manageable white-box attack [236].

5.1.2 Privacy Preserving. In order to defend against privacy attacks, privacy-preserving methods have been proposed based on different strategies, which can be broadly divided into five categories: differential privacy, federated learning, adversarial learning, and anonymization & encryption.

- **Differential Privacy (DP)** is a common way to preserve users' privacy, which can provide strict statistical guarantees for data privacy [91]. The main idea is to add random noise into data to protect actual data while preserving recommendation accuracy [92, 264].
- **Federated Learning (FL)** keeps the training data and recommendation models decentralized, which isolates users' data and the cloud server by only transferring the parameters between them. Thus, it avoids privacy leakage during the data transfer [404, 427]. Specifically, every client uses their data to train a decentralized recommender system locally. Then, the server collects these models' parameters and aggregates them into a new set of recommendation parameters for update [179, 217]. The mathematical model can be defined as follows:

$$\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \sum_{k=1}^n p_k \mathcal{L}_k(\theta), \quad (26)$$

where θ is the global recommender system's parameter, n is the number of the decentralized devices, p_k and \mathcal{L}_k represent the weight and the loss function on the k -th device, respectively.

- **Adversarial Learning (AL)** is a relatively general method for privacy-preserving recommendations, which can be formulated as the minimax simultaneous optimization of recommendation and privacy attacker models [159]. The mathematical model can be formulated as

³<https://www.businessinsider.com/stolen-data-of-533-million-facebook-users-leaked-online-2021-4>

Table 6. The categories of privacy attack methods on recommender systems

	Taxonomy	Related methods
Privacy Attacks	Membership Inference Attacks	[79, 431]
	Property Inference Attacks	[14, 115, 277, 437]
	Reconstruction Attacks	[42, 90, 151, 257, 257, 303]
	Model Extraction Attacks	[418]

follows:

$$\min_{\theta} \max_{\Psi} \mathcal{L}_{rec}(\theta) - \alpha \mathcal{L}_{adv}(\Psi), \quad (27)$$

where θ and Ψ represent the recommender system's parameters and adversary model parameters, respectively, maximizing the loss of adversarial could enhance the attack ability. The minimization of learning loss could provide better performance. α is a hyperparameter to trade-off the contribution of adversarial to the training.

- **Anonymization & Encryption.** Both of them protect the user's privacy by adding noise. Anonymization [248, 336] obscures the privacy attributes of users. Then make it impossible to correlate the privacy attributes with individual identities of people. Encryption techniques prevent people who do not have the authorization from any useful information [125, 413].

5.2 Methods

5.2.1 Privacy Attacks. In this subsection, we will introduce some representative methods of privacy attacks in recommender systems, which are summarized in Table 6. It is worth mentioning that many attack methods utilize shadow training (i.e., building a surrogate system) [316] to generate the training data for the attacker. Shadow training can be divided into two steps: training shadow models and using the predictions of the models to train the attacker [157]. First, use the shadow data (e.g., public auxiliary data) to train a series of shadow models that could mimic the behavior of the target recommender system. Then, the attack model can be well trained using the predictions of the shadow models. The framework of shadow training in recommender systems is shown in Figure 8.

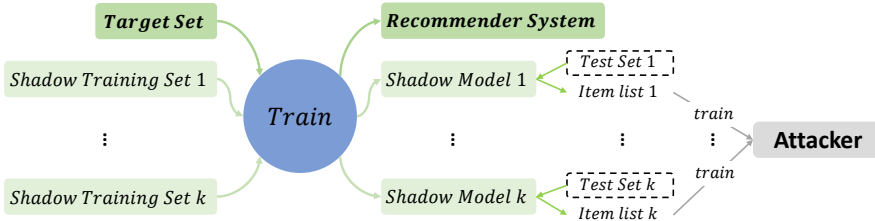


Fig. 8. Shadow training strategy in privacy attacks for recommender systems. The training process consists of two steps: training shadow recommendation models with public accessible auxiliary data and using the recommendation predictions to train the attacker.

- **Membership Inference Attacks (MIA).** The goal of MIA is to identify whether the target user is utilized to *train* the target recommendation model. The general procedure of membership inference attacks is illustrated in Figure 9. By querying the recommender system, the attacker will obtain a recommendation list for the target user. Then the attacker can utilize

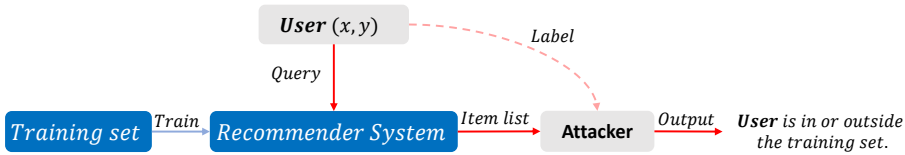


Fig. 9. Membership inference attack. The attacker queries the recommender with a user and obtains the corresponding item list. Afterward, the attacker can infer if the user is in the training set based on the discrepancy between recommended item list and label (i.e., the historically interacted items).

the distribution difference between the recommended item list and the user's historically interacted items to infer whether the target user is used in the training process. Furthermore, some attacking models utilize the temporal differences of the items ranking list to conduct MIA without the label [42]. For instance, to quantify the privacy leakage in recommender systems, Zhang et al. [431] attempt to conduct MIA by measuring the similarity between recommended item list and the user historically interacted item list to infer whether the data of the target user is used by the target recommender system.

- Property Inference Attacks (PIA).** Different from the MIA aiming to infer whether the target user is in the training set or not, property inference attacks mainly focus on the global sensitive information in the training set. To get useful information from machine learning classifier, Ateniese et al. [14] train a series of classifiers to form a meta-classifier, which could recognize the unexpected but useful properties of the target training set. To attack graph-structured data, the work of [437] conducts PIA on graph neural networks to infer the global information of the input graph, such as the number of nodes and links. In addition, the attacker training also adopts the same strategies as shadow training to build shadow models by utilizing auxiliary graphs. Recent works [115, 277] adopt shadow training to conduct property inference attacks against the full connected neural network and convolutional neural networks classification tasks, respectively.
- Reconstruction Attacks.** Recommender systems also face a privacy risk similar to that exists in statistical database queries, i.e., using the user's publicly available attributes and multiple queries to reconstruct their sensitive information [288]. For instance, reconstruction attacks can use the non-sensitive rating information to reconstruct the sensitive features of a certain user by matrix factorization [257]. The work of [42] utilizes the temporal changes in the recommendation item list to conduct a reconstruction attack. With the help of shadow training, recent works [90, 303] leverage the differences in predictions when a new user is added to train the target recommender system to reconstruct the information of the new user. Moreover, the work of [151] reconstructs recommended sensitive items by using data poisoning attacks, in which malicious users with certain items are used to link with the user's private information.
- Model Extraction Attacks.** Another research field close to this attack is Knowledge Distillation (KD) [152]. KD uses a complex "teacher" model to compress into a "student" model that can mimic the behavior of "teacher" but with a simpler structure. KD can be considered as a special case of model extraction attack with strong assumptions, such as accessing the parameters of the "teacher" model and the training data [272]. In contrast, the model extraction attack is usually conducted under the black-box setting. Differing from the other three privacy attack methods that seek to infer the users' information, the model extraction attack targets the parameters and structure of recommender systems. The work of [418]

attempts to steal parameters of sequential recommender systems by utilizing the specific autoregressive regimes of sequential recommender systems. Specifically, the attacker uses the recommended items to generate a series of item sequences, which act as the training data for the attacker. These sequential data can be used to train a victim recommender system which can imitate the behavior of the target recommender system.

5.2.2 Privacy-preserving Methods. In this subsection, we investigate the privacy-preserving approaches in recommender systems. The structure will follow the taxonomy in section 5.1.2, namely differential privacy, federal learning, adversarial learning, and anonymization & encryption. The representative methods for each part are summarized in Table 7.

Table 7. The categories of the privacy-preserving in recommendation systems.

	Taxonomy	Representative Methods
Privacy-preserving Methods	Differential Privacy	[45, 46, 395, 429, 432, 459]
	Federated Learning	[111, 138, 160, 218, 284, 376, 378]
	Adversarial Learning	[22, 208, 229, 295, 352]
	Anonymization & Encryption	[53, 163, 281, 302, 360, 402, 413, 430]

- Differential Privacy.** These privacy-preserving methods can be effective in resisting membership inference attacks by adding random noise [51]. Xue and Sun [459] propose two differential privacy-based methods to preserve users' privacy information, namely differentially private item-based recommendation and differentially private user-based recommendation. To perform privacy-preserving cross-domain recommendation, a two-stage based method (PriCDR) is proposed to firstly preserve the data privacy in the source domain via differential privacy techniques and then transfer cross-domain knowledge to enhance recommendation performance in the target domain [45]. To defend against privacy attacks in various scenario, researchers have successfully applied differential privacy to point-of-interest recommendations [46] and GNNs based recommendations [432]. In view of the trade-off between recommendation performance and privacy protection, Xiao et al. [395] propose a deep reinforcement learning (RL)-based privacy protection method. Differential privacy is utilized to protect users' privacy of the recommender system, and RL is responsible for the choice of privacy budget, which can optimize the privacy budget over time based on the estimated privacy loss. The work of [429] proposes a differential privacy-based privacy-preserving framework (PLORE) to protect the users' privacy in location recommendation systems. PLORE can balance the recommendation performance and the privacy-preserving using the probabilistic differential privacy mechanism.
- Federated Learning.** For federated learning-based privacy-preserving methods, the user's data can be stored in the local devices for training a local recommendation model, while a server in the cloud is responsible for aggregating the distributed model parameters. With such special architecture, federated learning-based recommender systems can avoid the leakage of users' privacy data naturally. For example, FedGNN is proposed to combine federal learning and GNNs based recommendation systems to protect the users' privacy [376], where each client device stores a user-item graph and a local GNN recommendation model. In addition, differential privacy techniques are further used on the local gradients to protect user privacy. In addition, recent works [284, 378] also incorporate local differential privacy for the gradients to ensure the privacy protection on news recommendation. Since these methods may balance the trade-off between privacy protection and the recommendation

performance, Li et al. [218] theoretically analyze such trade-off and provide a privacy error bound. Guo et al. [138] propose an edge-accelerated framework PREFER, which aggregates the decentralized parameters on the edge server (e.g., base station) rather than the cloud server, so as to meet the real-time recommendation needs. The works of [111, 160] introduce federated learning-based multi-view recommendation frameworks to address the cold-start issue in recommender systems.

- **Adversarial Learning.** Adversarial learning is also an effective technique to protect privacy with an attacker discriminator [159, 193, 258]. The work of [22] proposes an adversarial learning-based recommendation framework to defend against reconstruction attacks, which consists of two main components, bayesian personalized ranking recommendation systems, and a reconstruction attacker. In general, the recommendation task is proposed to model the user's preferences, while attackers' gain is minimized to protect users' privacy, which can be formulated a min-max game. On top of that, the user representations themselves may be used together with external data to recover users' sensitive information. To address this, Resheff et al. [295] propose an adversarial recommendation method, which adds the loss of demographic prediction tasks together with the recommendation tasks to the parameters optimization process. Then, the elimination of demographic information about the user can be controlled by hyperparameters. Additionally, increasing attention has been paid to develop adversarial privacy-preserving methods on GNNs [208, 229, 352], which can contribute to enhance the design of privacy-preserving recommender systems.
- **Anonymization & Encryption.** Similar to differential privacy, the goal of anonymization and encryption is to protect sensitive information by adding noise or mapping data to another feature space. In general, anonymization techniques aim to prevent the public data from being linked to individual identities of people, while encryption techniques make data unreadable to those who do not have the key to decrypt it⁴. In [53], Chen et al. propose a suppression and permutation-based anonymous method for collaborative filtering, which can reserve users' privacy by limiting the probability of a successful passive privacy attack. More specifically, the suppression operation is used to suppress an item from a related item list, while the permutation operation aims to permute an item that has climbed in a related item list to a lower position.

Besides, for the encryption techniques in recommender systems, PLAS [360] and EPRT [281] propose homomorphic encryption technologies to protect users' sensitive information. Zhang et al. [430] utilize the BGN Cryptosystem to conduct privacy protection. By adding noise to the original data against privacy attacks, Yu et al. [413] combine knowledge graph enhancement techniques with multi-task learning to improve recommendation performance while adding Gaussian noise to sensitive data to protect privacy. Huo et al. [163] add Laplacian distributed noise to fuse the users' social relationships.

In addition, there are also distortion-based methods [302, 402], where they map the original data into a new feature space through a probabilistic mapping function and achieve data privacy protection under a certain distortion budget constraint.

5.3 Applications

This subsection presents some representative examples of using privacy-preserving techniques to protect sensitive information in real systems.

- **Private medical recommender systems.** In view of the fast development of eHealthcare, the privacy protection of the users' sensitive information attracts great attention. Xu et

⁴https://media13.connectedsocialmedia.com/intel/01/9768/Using_Data_Anonymization_Enhance_Cloud_Security.pdf

al. [398] propose a privacy-preserving medical service recommendation method based on the modified Paillier cryptosystem, truth discovery technology, and the Dirichlet distribution. The work of [188] employs cryptographic privacy-enhancing protocols to deal with the privacy challenges in health service recommender systems. To recommend the trusted physician for patients without privacy leakage, Hoens et al. [153] propose two privacy-friendly recommender frameworks SPA and ACA, where they use secure multiparty computation techniques and anonymization to protect the users' private data, respectively.

- **Location-private recommender systems.** Recent years have witnessed the increased need for location-based services due to the prevalence of smart mobile devices. For example, Facebook and Google Map will collect the location information of the users to conduct more accurate recommendations. It is important to protect sensitive location information. In [447], Zhao et al. obfuscate the rating vectors of the recommender to keep the location information secure. In [429], a recommendation framework PLORE is proposed to address the location privacy challenges, where PLORE uses differential privacy to give sensitive location data privacy guarantees. To hide the real location data, Gao et al. [116] adopt the differential privacy to obfuscate the historically visited locations data.

5.4 Survey and Tools

This subsection collects some existing surveys and tools regarding privacy in recommender systems to facilitate researchers in this field.

5.4.1 Surveys. Recent comprehensive surveys on privacy-preserving recommendations are summarized in [6, 161]. Specifically, [6] investigates privacy attacks and privacy-preserving methods in large-scale social recommendation systems and discusses the main issues in the privacy protection of the online social network. The work of [161] gives the taxonomy of recommender systems and the definition of privacy leakage, in which the characteristics and specific measures of different privacy-preserving methods are compared. On top of that, privacy in machine learning and deep learning is comprehensively reviewed in [259, 297].

5.4.2 Tools. For differential privacy, there are some popular tools, such as Facebook Opacus⁵, TensorFlow-Privacy⁶, OpenDP⁷, Diffpriv [299] and Diffprivlib [154]. For federated learning, popular tools include TFF⁸, FATE⁹, FedML [143], and LEAF [43]. Popular tools in Homomorphic Encryption are Awesome HE¹⁰ and TF Encrypted¹¹.

5.5 Future Directions

Regarding privacy attacks and preservation in recommender systems, we have introduced many methods above. However, the need for privacy-preserving always comes and goes. There are many unresolved privacy issues in the field of recommender systems, including the following three points.

- **Privacy and performance trade-off.** Whether it is differential privacy, anonymization, or encryption, the most prominent means of countering privacy attacks is adding noise to the original data or distorting it. These methods protect privacy while reducing the utility. Because the sensitive information obscured is often the critical information that

⁵<https://opacus.ai/>

⁶<https://github.com/tensorflow/privacy>

⁷<https://opendp.org/home>

⁸<https://github.com/tensorflow/federated>

⁹<https://github.com/FederatedAI/FATE>

¹⁰<https://github.com/jonaschn/awesome-he>

¹¹<https://github.com/tf-encrypted/tf-encrypted>

affects the performance of the recommendation. Therefore, depending on different task requirements, how to protect privacy with minimal performance cost may be a continuous research direction.

- **Comprehensive privacy protection.** The privacy protection methods usually protect against only one kind of privacy attack. However, the actual usage scenario is that a recommender system is subject to multiple privacy attacks. It is still challenging research to combine these privacy protection approaches without degrading the recommendation performance or to propose a comprehensive privacy protection framework.
- **Defence against shadow training.** Among the four privacy attack methods mentioned in this paper, membership inference attack, attribute inference attack, and reconstruction attack all use shadow training methods to train attackers. The training method provides vital support to the privacy attacks but is indeed trained under reasonable assumptions. Therefore, investigating how to defend against such training methods is crucial for privacy protection.

6 ENVIRONMENTAL WELL-BEING

The application of deep learning brings great success to recommender systems [381], which elevates recommending precision to a high level. However, more sophisticated models also result in longer training time and larger energy consumption. For example, Alibaba, an e-commerce site, has to consume several hours to train its model, which contains tens of billions of parameters on hundreds of servers [172]. Adnan et al. [4] study that training a model on the Taobao dataset needs 621 minutes with 4 GPUs, whose average GPU power consumption is 56.39W per hour. As recommender systems have been adopted in many aspects of society, demands for large recommendation models will constantly increase. If we do not control the huge source consumption by recommender systems, the environmental damage will force humans to distrust and even give up the recommendation technique. Therefore, how to build environmental-friendly recommendation models is one essential component of trustworthy recommendation.

In this section, we will conclude existing works on energy saving in the recommendation research field. Firstly, we introduce the concepts and taxonomy of environmental well-being techniques. Next, we summarize some methods to reduce the storage and energy consumption of recommender systems, including model compression and acceleration methods. Then, some applications in real systems are listed. Finally, we summarize some surveys and tools of this topic and give some promising directions.

6.1 Concepts and Taxonomy

In this subsection, we will refer to the concepts and taxonomy of techniques that benefit environmental well-being.

6.1.1 Concepts. As the development of recommender systems, the requests for storage and computation resources increase rapidly. To tackle the problem of intensive natural resources, model compression and acceleration techniques are proposed. The model compression [67] aims to shrink the size of recommendation models to save storage resources. The acceleration techniques [62, 260] focus on reducing training or inference time to save computation resources.

6.1.2 Taxonomy. According to the characteristics of recommendation models, model compression methods are devised for embedding layers or middle layers specifically, while acceleration techniques focus on the training or inference stage.

- **Model Compression.** Different from general deep models, the embedding layers always account for most of the parameters in recommendation models [246, 350, 410], so we categorize model compression techniques for recommendation into two types: (1) *Embedding Layer*, which is always used to map discrete feature, such as ID feature and categorized feature, into a more expressive dense vector. (2) *Middle Layer*, which extracts the user’s preference or relations between features, such as the self-attention layer.
- **Acceleration Techniques.** As for acceleration, those techniques can be divided for *training* and *inference* purposes because models are always placed on different platforms and devices when training or inference.

6.2 Methods

In this subsection, we summarize storage-saving and energy-saving methods for recommendation models, i.e., model compression and acceleration techniques.

Table 8. Classification of Model Compression Methods for Recommendation

	Embedding Layer	Middle Layer
Hash	[80, 209, 307, 438, 456], [184, 227, 313, 355, 422]	[307, 355]
Quantization	[173, 226, 228, 234, 385, 394], [56, 142, 222, 241, 312, 354, 428]	[222, 354, 385]
Knowledge Distillation	[60, 182, 203, 342, 358], [52, 183, 194, 388, 457]	[60, 182, 203, 342, 358], [52, 183, 194, 388, 457]
Neural Architecture Search	[66, 237, 242, 401, 445, 448], [56, 175, 232, 239, 366]	[52, 326]
Others	[128, 311, 332]	[55, 311, 332]

6.2.1 Model Compression. Model compression [67] aims to cut down model size for more efficient training and inference, which is environmental-friendly. While general model compression techniques are mainly designed for various neural network layers in a model, we categorize the recommendation model only into two parts: the embedding layer and the middle layer. Meanwhile, there are five types of methods to achieve model compression in recommendation: hash, quantization, knowledge distillation, neural architecture search, and others. Representative works are listed in Table 8 for a clear description.

- **Hash.** The embedding table is always extremely large in the recommendation model because of plenty of items and users. Therefore, how to compress the embedding layer is a focus in the academic and industrial fields. Hash has been proved an efficient method. The hash function $h(\cdot)$ maps original concrete features $\mathbf{x} \in \{0, 1\}^n$, such as ID, into relatively short binary codes $\mathbf{y} \in \{0, 1\}^m$, where n and m are the vocabulary size of original and hashed feature respectively [368]. The embedding table is compressed much due to $m \ll n$, and thus the memory cost of the embedding table decreases. We cluster hash methods into two groups as follows:
 - *Data-independent methods.* In this thread of methods, the hash function $h(\cdot)$ is pre-defined without considering the dataset. Locality Sensitive Hashing (LSH) [129] is one representative method that can generate pairs of similar items. Two news recommendation works [80, 209] make use of LSH to cluster similar news items and then find similar users to recommend relevant news.

- *Data-dependent methods.* Data-dependent methods always learn hash functions $h(\cdot)$ for specific dataset. For example, some works highlight the importance of users' preference over items and propose to preserve such information when learning hash function [313, 438, 456]. Lian et al. [227] combine with auxiliary information to learn better binary codes. Works of [307, 422] propose to adopt multiple hash functions to tackle collision problem in the hash. Besides, DHE [184] gives a new idea to get embedding without using an embedding table, which adopts multiple hash functions and DNN to get embedding directly. Though most hash methods are designed for a lightweight embedding layer, Bloom Embeddings [307], and ABinCF [355] also optimize the middle layer.
- **Quantization.** Quantization is also an efficient technique to compress the embedding layer. In this type of method, the embedding of one feature will be clustered into several classes, and each embedding can be represented by the center of its cluster, named codeword, which greatly decreases the number of embeddings. To enhance the ability of representation, representation space is always decomposed, and an embedding is quantized to subvectors by several codebooks. An item's quantized representation vector $\mathbf{q}_i \in \mathbb{R}^D$ can be computed as follow:

$$\mathbf{q}_i = f(c_{w_i^1}^1, c_{w_i^2}^2, \dots, c_{w_i^B}^B), \quad (28)$$
 where $c_{w_i^b}^b \in \mathbb{R}^D$ is the w -th codeword in the b -th codebook. $f(\cdot)$ is the composing function of each subvector. According to the types of composing functions, quantization methods can be categorized into product quantization [170], additive quantization [15] and residual quantization [61]. It is worth mentioning that product quantization and additive quantization are often adopted in recommendation field. Besides, some recent works using compositional embedding also belong to quantization.
 - *Product Quantization (PQ).* PQ is a type of quantization method that composes quantized vectors by product. LightRec [226] and pQCF [228] adopt PQ to compress user and item embedding size for matrix factorization methods. Furthermore, Jiang et al. [173] propose a memory-efficient factorization machine based on PQ. Compared with traditional PQ methods, recent works [173, 226, 394] integrate quantization into training process for optimal models. Liu et al. [234] design an online optimized product quantization for the online recommendation model specifically. Besides, LISA [385] and MDQE [354] even utilize PQ to lightweight and accelerate self-attention layer for sequential recommendation models.
 - *Additive Quantization (AQ).* AQ uses add operation to compose vectors. For example, Liu et al. [241] propose an online AQ, which is more efficient than online PQ [234]. Zhang et al. [428] design a new loss for AQ and achieve a lower approximation error in contrast to PQ.
 - *Compositional Embedding.* Recently, another special thread of quantization methods prevailed, named compositional embedding. The main idea of these methods is to generate meta embedding for each feature based on their characteristics [312]. Chen et al. [56] and Hang et al. [142] propose a lightweight compositional embedding for on-device and online recommendation models respectively. Besides, compositional embedding also can be used to create a lightweight self-attention layer in sequential recommendation [222].
- **Knowledge Distillation.** Knowledge Distillation (KD) is one of the most crucial techniques in model compression [131], and is also adopted in the recommendation field for lightweight models. KD aims to use a smaller model (student model) to approximate the capacity of the original big model (teacher model). The key is distillation loss. Similar to [131], existing KD methods for recommendation models can be categorized into two groups according to the distillation loss: response-based and feature-based.

- *Response-based Methods.* Response-based methods transfer knowledge via the output layer of the teacher models, and the distillation loss can be formulated as:

$$\mathcal{L}_{res} = \mathcal{L}_R(z_t, z_s), \quad (29)$$

where z_t and z_s are the logits of teacher and student models, respectively, and $\mathcal{L}(\cdot)$ refers to the divergence loss function. Tang et al. [342] firstly adopt KD to ranking problems and propose the ranking distillation (RD) method. Then, Lee et al. [203] design a novel sampling technique and a new distillation loss function to improve RD. Based on these two methods, researchers further propose KD methods for some specific recommendation tasks, such as POI recommendation [358], sequential recommendation [388] and CTR prediction [457]. Besides, Kweon et al. [194] propose a bidirectional distillation method to elevate the accuracy of teacher and student recommendation models simultaneously.

- *Feature-based Methods.* Compared to response-based methods, feature-based methods transfer the knowledge in intermediate layers of teacher models. The distillation loss of this thread is as follows:

$$\mathcal{L}_{feat} = \mathcal{L}_F(f_t(x), f_s(x)), \quad (30)$$

where $\mathcal{L}_F(\cdot)$ is the similarity function. $f_t(x)$ and $f_s(x)$ are the output from the middle layers of teacher and student models. Recent studies [52, 183] propose to transfer the structure information of teacher models to student models. To fuse external knowledge into models, Chen et al. [60] design a novel training scheme with feature-based KD. Furthermore, Kang et al. [182] combine response-based and feature-based KD, and propose a DE-RRD method.

- **Neural Architecture Search.** Recently, applying automated machine learning (AutoML) techniques to design neural architecture for deep recommendation models has become a hotspot [30, 451]. Neural Architecture Search (NAS) aims to search for the optimal architecture for deep models, which can prune the redundant parameters. The general idea of NAS is utilizing the validation loss to adjust the model architectures. Therefore, the objective of NAS can always be formulated into a bi-level optimization problem:

$$\min_{\mathcal{A}} \mathcal{L}_{valid}(\mathcal{W}^*(\mathcal{A}), \mathcal{A}), \quad (31)$$

$$s.t. \mathcal{W}^*(\mathcal{A}) = \arg \min_{\mathcal{W}} \mathcal{L}_{train}(\mathcal{W}, \mathcal{A}), \quad (32)$$

where \mathcal{L}_{train} and \mathcal{L}_{valid} are training and validation loss respectively. \mathcal{A} and \mathcal{W} are parameters and architectural weights of recommendation models. Most of the NAS works are beneficial to the embedding layer compression, but Song et al. [326], and Chen et al. [52] also utilize neural architecture search to compress the middle layers of recommendation models. The works that aim at the embedding layer can be categorized into two groups: embedding dimension search and automated feature selection.

- *Embedding Dimension Search.* Some related studies focus on searching for optimal and minimal embedding size for each feature, which can compress the embedding layer efficiently. For example, ATML [401] and PEP [242] propose a gradient-based and a pruning-based solution to search optimal embedding size for each feature. However, these two methods face the problem of a vast search space of embedding size. To tackle this problem, AutoEmb [448] and ESAPN [237] cut the embedding into several segments to reduce the search space, which can be called row-wise methods. In detail, AutoEmb designs a soft selection strategy to combine different segments with learnable weights. By contrast, ESAPN proposes a frequency-based hard selection strategy. Compared to AutoEmb and ESAPN, some other works group embedding with different values of a feature field to shrink search space named column-wise methods. AutoDim [445] is a special case of the column-wise method

- because it searches for a unified dimension for each feature field. In detail, AutoDim fuses the embedding of different dimensions by learnable weights, which is similar to AutoEmb. Cheng et al. [66] group values of a feature field based on frequency. Furthermore, NIS [175] and RULE [56] propose to combine both row-wise and column-wise methods.
- *Automated Feature Selection.* In addition to searching embedding dimensions, some works of automated feature selection [232, 239, 366] are also useful for lightweight embedding layers because they decrease the number of input features. Liu et al. [239] propose a method based on reinforcement learning, which regards feature selection as a multi-agent problem. AutoField [366] equips with a controlling architecture to calculate the drop and select probability of each feature field and retrain the recommendation models after selection. Compared with the previous works selecting a fixed set of features for a dataset, Lin et al. [232] propose to select various features for each user-item interaction to capture the dynamics of the practical recommender system.
 - **Others.** There are also some other novel techniques used to compress recommendation models. Shen et al. [311] utilize the alternating direction method of multipliers (ADMM) to jointly optimize feature selection and model compression. Sun et al. [332] propose an adaptive decomposition method for lightweight input and output layers and a parameter sharing scheme to compress middle layers. QFM and QNFM [55] adopt quaternion representations to decrease parameters of factorization machine models. Ginart et al. [128] design a mixed dimension embedding scheme to shrink the memory of the embedding layer and simplify neural architecture search into tuning hyper-parameters.

Table 9. Classification of Acceleration Techniques for Recommendation

		Training	Inference
Hardware-related	Near/In Memory Computing	[196]	[78, 164, 190, 195, 367, 371]
	Cache Optimization	[135, 165, 403, 442]	[93, 397]
	CPU-GPU Co-design	[4, 5, 197, 308, 441, 450]	-
Software-related	Optimization	[128, 137, 146, 411, 454]	[140, 141]
	Efficient Retrieval	-	[81, 113, 191, 287], [238, 263, 339, 400]

6.2.2 Acceleration Techniques. In addition to model compression, the acceleration is also a useful way to reduce training and inference time, and thus can save energy and is environmental-friendly. Existing acceleration techniques [62, 260] mainly focus on: (1) memory-based challenges, which is the difficulty of data access by computation units, and (2) computation-based challenges, which means huge and complex computation. In the recommendation field, we conclude that hardware-related methods always aim at the memory-based challenge, and software-related methods are mainly for computation-based challenges. We will introduce hardware- and software-related methods, respectively, and the representative works are listed in Table 9.

- **Hardware-related.** The advancements in computing units and hardware accelerators (e.g., GPUs) are huge. However, the memory techniques improve much slower. Such a growth gap leads to the problem of memory wall and hinders the improvement of acceleration techniques. Google [32] has proved that data movement between memory and computing units accounts for 62.7% of energy consumption across their many applications. As for the recommendation field, though large embedding tables bring high accuracy to recommendation models, they also cause severe storage and communication burdens [246, 410, 442]. Many hardware-related

methods aim to optimize data moving between the storage device and computing units in model training or inference. We categorize hardware-related methods into three lines:

- *Near/In Memory Computing (NMC/IMC)*. The main idea of NMC is to put computing units closer to the memory, which can lower the distance of data moving and thus reduce latency. Kwon et al. [195] firstly adopt NMC to accelerate the embedding layer for recommendation models. In detail, they design an NMC architecture, named TensorDIMM, based on general DRAMs to elevate bandwidth for embedding operations. RecNMP [190] is proposed for co-location operators for embedding vectors. DIMMining [78] is for Graph-based models specifically. TensorDIMM, RecNMP, and DIMMining are all tailored for DIMM, which is always for GPU-based systems. By contrast, Centaur [164] is designed for FPGAs in CPU-based systems. The methods mentioned above can accelerate model inference efficiently but are not fit for training. Tensor Casting [196] accelerates embedding layer training by design for tensor gather-scatter. Different from NMC, IMC puts memory and computing units together. ReRec [367] and RecSSD [371] design processing units on DRAMs and SSD according to embedding access frequency.
- *Cache Optimization*. The cache mechanism is to specifically store some data that are accessed frequently on the memory device. Though cache provides low communication latency, its size is much limited. Therefore, efficient usage of cache is vital for acceleration. For instance, AIBox [442] caches embedding table on SSDs for training CTR models on a centralized system. Compared to AIBox, ScaleFreeCTR [135] is a distributed training system that contains a MixCache mechanism to accelerate embedding synchronization. Xie et al. [397] propose a novel cache query mechanism to speed up embedding lookup on GPUs. Yang et al. [403] design different precision for training on general memory and cache. Ibrahim et al. [165] and Eisenman et al. [93] both propose an embedding placement strategy to elevate the efficiency of cache usage. It is worth noting that all of the methods mentioned are mainly for training recommendation models, except for Fleche [397] and Bandana [93].
- *CPU-GPU Co-design*. Hybrid CPU-GPU mode prevails in industrial recommendation systems due to huge embedding tables. General recommendation models can be partitioned into the embedding part and DNN part. The embedding part is always stored and processed on the CPU, and the DNN part is on GPU [441]. However, the communication of embedding vector between CPU and GPU is a big challenge, so the CPU-GPU co-design is in need. Adnan et al. [4] propose to put highly accessed embeddings on GPU memory to reduce communication time. Similarly, Recshard [308] partition embedding into several parts according to the distribution of training data and store each part on hierarchical memory. BiPS [450] focuses on the parameter update process during the model training. It contains a bi-tier parameter server to accelerate parameter synchronization. Recently, Adnan et al. [5] and Kwon et al. [197] both propose a data pipeline to speed up parallel training.
- **Software-related**. There are many studies [149, 282, 290] on designing accelerators for DNN to tackle the computation challenges, but not all of them can fit recommendation models. Besides, embedding plays a vital role in the recommendation field, so how to accelerate embedding computation needs exploring. Based on these two observations, we group software-related methods into two groups: Optimization and Efficient Retrieval.
 - *Optimization*. In order to accelerate training recommendation models, some works focus on the training process, such as Cowclip [454] and MetaBalance [146]. Cowclip uses a large batch to train recommendation models and contains an adaptive clipping strategy to maintain training accuracy. Metabalance is proposed to adjust gradients of each loss for the multi-task recommendation dynamically and can reduce training epoch efficiently.

Besides, some other works accelerate training by optimizing the embedding layer. Yin et al. [411] adopt tensor train decomposition to recommendation models and optimize the batched matrix multiplication. Ginart et al. [128] design mixed dimension embeddings for recommendation models. In HyperRec [137], the general embedding layer is replaced by high-dimensional binary vectors, whose computation is more efficient on various computing devices. In addition, DeepRecSched [140] and RecPipe [141] are two optimization methods for inference. DeepRecSched is an effective scheduler to optimize parallel data movement and can reduce latency for several recommendation models. RecPip partition recommendation models into multistage and then compute each stage in a parallel mode under the control of the designed scheduler.

- *Efficient Retrieval*. In industrial, it is common that train user and item embeddings offline to represent their preference and attributes, then get recommending list by Embedding-Based Retrieval (EBR) online. Such a schema can reduce inference latency efficiently under the condition of a huge volume of users and items. Many studies [75, 132] aim to refine the embedding for accuracy, but EBR is also important for efficiency. The main idea of efficient retrieval is to build an index for each embedding to achieve sub-linear searching times rather than conducting the inner product directly. According to the type of index, we categorize efficient retrieval methods into four clusters: tree-based, graph-based, hash-based, and quantization-based. Tree-based methods partition high-dimensional space into several sub-space and use leaf nodes on a tree to represent each sub-space. Due to the lower time complexity of the index, it has been used for efficient retrieval for a long time, such as KD-tree [113]. To further improve efficiency, randomized-partition trees (RP-tree) [81, 191] are proposed with minimal loss of accuracy. Ram et al. [287] combine KD-tree with RP-tree to reach a trade-off. Recently, graph-based methods have attracted much attention. This line of methods builds a similarity graph for retrieval, in which each vertex in the graph represents an embedding, and the edge shows the distance between embeddings. Morozov et al. [263] firstly utilize a Delaunay graph to construct an index graph and propose the theory of similarity graph. ip-NSW+ [238] and NAPG [339] further improve the ip-NSW [263] by introducing an angular similarity graph and a hierarchical navigable small world graph respectively. Xu et al. [400] propose the IPGM algorithm to tackle the problem of node deletion in the similarity graph. Besides, some of hash methods [209, 227, 438, 456] and quantization methods [226, 228, 428] can also be utilized to accelerate EBR, which have been introduced in Model Compression part.

6.3 Applications in Real Systems

Many companies and researchers have devoted efforts to the environmental problems when developing modern recommender systems in many aspects of society. In this subsection, we will introduce these actions from the following three aspects: Big Model, Edge Computation, and Embedding-Based Retrieval Systems.

6.3.1 Big Model. In the field of NLP, Pretrained Language Models (PLMs) such as BERT [86], ERNIE [331] and GPT-3 [35] prevail, because PLMs with only low-cost fine-tuning can adapt to various tasks and possess high precision. Such a training paradigm can emit training a whole model for each specific task and thus save energy largely. Based on the success of PLMs in NLP, many studies aim to shift this scheme to the recommendation field. As the natural similarity between recommendation and NLP, several early works [314, 415, 417] consider a user behavior history as a sentence, and then they can train a big model on datasets from different tasks. The pretrained model often generates a universal user representation, which can be used in many downstream tasks and

get rid of the huge energy cost by retraining for each task. However, these works require all tasks use a unique item set because they are trained based on the item id. To tackle this problem, recent works [76, 123, 155, 212] utilize textual information of items to covert recommendation task to a text-to-text task, which is more suitable for PLMs. Besides, TransRec [356] fuses multi-modality of items to avoid using item id. These works urge recommendation to step into the big model period and contribute to environmental well-being.

6.3.2 Edge Computation. In the traditional service framework, most computations are conducted on the cloud, which may cause high latency of service and high cost of communication. To tackle these challenges, the edge-cloud scheme becomes a hotspot [408]. Due to the limited capacity of edge devices, the model on edge must be computing- and storage-efficient. Alibaba [130] firstly designs a recommender system for edge specifically named EdgeRec, and applies it to their e-commerce application-taobao. EdgeRec deploys the rerank module and user behavior module on edge to respond to users' dynamic preferences and capture real-time behaviors. To avoid storing the whole embedding table on edge, it fetches items and corresponding embedding in each request. Furthermore, Alibaba [407] achieves their "thousands of people with thousands models" with edge-cloud collaboration. Unlike EdgeRec, it can also update the user's model on the edge device, which benefits long-tailed users. To further improve the edge-cloud collaboration, Yao et al. [406] design a meta controller to optimize the recommending list from edge and cloud recommenders. Another recent work [65] regards cloud and edge as slow and fast components and designs a bidirectional collaboration mechanism to benefit both.

6.3.3 Embedding-Based Retrieval Systems. EBR has been widely adopted in many aspects, such as question answering [186, 329], web search [98, 206] and recommender system [158, 216]. EBR plays the main role in the recall stage in a recommender system, and a good EBR system should meet the trade-off of three key points: memory, latency, and accuracy. To achieve these goals, many companies build efficient EBR systems for their services. Airbnb [132] introduces their EBR system mainly from the aspect of how to construct informative embedding to get higher accuracy. Facebook [158] proposes a hybrid EBR system with boolean and KNN matching to elevate recall efficiency. Besides, Amazon [267] publishes their semantic product search model for embedding retrieval. Next, Alibaba [216] and JD [426] also propose MGDSPR and DSPR respectively for their e-commerce applications.

6.4 Surveys and Tools

In this subsection, we give out some surveys and tools about environmental well-being for readers who want to investigate this topic further.

6.4.1 Surveys. To our best knowledge, our survey is the first attempt to conclude works about the environmental well-being of the recommender system. Furthermore, a survey on model compression and acceleration technique for recommender systems is also few, so we list some surveys on these techniques in other fields, such as computer vision (CV). As for model compression, Cheng et al. [67] introduce four types of methods to compress DNN and compare compression rate between several methods on the CIFAR dataset. Among techniques of model compression mentioned above, there is a survey [357] about learning to hash, and another survey [247] about deep hash methods. Besides, Gholami et al. [127], and Gou et al. [131] survey the techniques about quantization and KD for DNN, respectively. It is worth noting that two surveys [30, 451] related to NAS in the recommender system came out recently. Chen et al. [30] categorize methods of NAS into four groups, in which feature selection and feature embedding are two types beneficial to shrink models. Another line of work is acceleration. Deng et al. [84] survey many acceleration methods for DNN, especially

hardware-related methods. Le et al. [199] summarize efficient retrieval methods of recommendation. This survey considers the retrieval pipeline and illustrate the methods from the aspects of two consecutive process: candidate generation and candidate ranking.

6.4.2 Tools. Faiss [177] is a popular tool of similarity search, which is published by Facebook. It implements many product quantization-based methods to accelerate the embedding retrieval. Besides, Faiss integrates a GPU k-selection algorithm, which can utilize GPU more efficiently.

6.5 Future Directions

More research has paid attention to the memory and energy cost of recommender systems for environmental purposes. However, one serious problem is the lack of a framework to measure and predict the energy consumption for recommender systems specifically, like SyNERGY [298]. Such an estimation framework will help researchers to study energy-efficient recommender systems. As for model compression techniques, NAS is a promising direction to shrink model sizes. However, most existing works aim to improve accuracy. How to get the trade-off between the size and accuracy of recommendation models is an interesting problem for NAS. For the aspect of acceleration, the design of collaboration between hardware and software may be a future direction. Edge recommendation is a great example.

7 ACCOUNTABILITY & AUDITABILITY

Accountability for recommendation refers to what extent users can trust recommender systems and who is responsible for the devastating effects brought by the recommender systems. Because recommender systems play the role of information messenger, it is vital to equip recommender systems with accountability to be trustworthy. For example, during a crisis such as the COVID-19 pandemic, recommender systems encouraged the spread of many fake news and conspiracy theories [275], which caused distrust in recommender systems and had harmful effects on society. Besides, due to the lack of transparency and explainability of deep recommendation models, not only general users but also professional experts are unable to control recommender systems absolutely [24, 169]. Recently, the emergence of auditability, which refers to the methodology of evaluating recommender systems, has helped build the accountable recommender systems from a new perspective. To further discuss the accountability and auditability of recommender systems, we first introduce the definition of accountability and the taxonomy of auditability. Then, we present some relevant surveys and tools. At the final of this section, we discuss some future directions to inspire the readers who concern about this topic.

7.1 Concepts and Taxonomy

In this subsection, we will introduce the concepts of accountability and the taxonomy of auditability for the recommendation.

7.1.1 Accountability. Accountability has various concepts among different applications of artificial intelligence. Loi et al. [245] indicate that the general concept of accountability at least includes three dimensions: responsibility, answerability, and sanctionability. From the view of recommendation, we interpret the three dimensions: (1) *Responsibility*. If a user accepts an uncomfortable or illegal recommendation, accountability requires recommender systems to know which part of the system should be blamed. (2) *Answerability*. If an recommender system is accountable, it can reveal the reasons when recommender system has a bad effect. (3) *Sanctionability*. Sanctionability refers that recommender systems should punish and mend the parts that cause harmful impacts. According to these dimensions and AI regulations [1] published by European Commission, we summarize four roles for an accountable recommender systems: (1) **Content Governors**. Content governors

are responsible for examining the facticity and noxiousness of "items" in a recommender system. When a malignant event is reported, they should give an explanation from the aspect of content and decide to remove which item, so they undertake the answerability and sanctionability. (2) **Model Designers**. Model designers build the recommendation models for service. On the one hand, they can design explainable models for answerability. On the other hand, they should make the models reproducible, which benefits the dimension of sanctionability. (3) **System Deployers**. System deployers not only need to deploy recommendation models online but also check the possible trustworthy problems brought by recommender systems to avoid detriments in advance, so they take on the task of Responsibility. (4) **Third-party Auditors**. Third-party auditors are vital to guarantee accountability for general AI systems [286], not excluding recommender systems. They play the role of Responsibility for pointing out existing and potential problems in recommender systems. Besides, they will also reveal the reasons which refer to the role of answerability.

7.1.2 Auditability. Algorithm audits indicate a class of methods to analyze the existence and reasons for harmful problems for AI systems. It can help implement an accountable recommender systems. The methods of recommendation algorithm audits can be categorized into two classes: external and internal audits.

- **External Audits.** External audits regard recommendation models as a black box, and utilize input and output data from recommender systems to evaluate the algorithm [200]. Based on the concepts, we know two roles for accountable recommender systems that are relevant to external audits: Content Governors and Third-party auditors. Content governors aim to identify and remove the items that contain harmful contents regardless of the algorithm itself [11, 41]. However, their works are always complex and burdensome, so that it is impossible to eliminate problematic issues of recommender systems only by content governors. Recently, the methods by third-party auditors became popular. Many works focus on recommendations of YouTube, one of the most popular video platforms, to examine malicious problems, such as inappropriate videos for children [274], user radicalization [296] and pseudoscientific misinformation [275]. These three works conduct three procedures for audits similarly. Firstly, they collect publicly available data from YouTube. Then, they classify normal and bad videos (such as radicalized videos) by manual annotations or well-trained classifiers. At last, they analyze the annotated data to probe problems. In conclusion, external audits can avoid the problem of subjectivity because no system developers and deployers join. However, it cannot be conducted before a recommender system is deployed.
- **Internal Audits.** Internal audits examine the problems with access to training data by the other two roles for accountable recommender systems, i.e., Model Designers and System Deployers. One of the most useful audits means for model designers is to enhance explainability for recommendation models [89], which can output reasons for a bad case. Another efficient way is to achieve reproducibility of recommendation models [24] because a reproducible environment gives auditors more chances to evaluate recommender systems with different strategies. To achieve reproducible recommendation models, many researchers propose recommendation benchmarks, such as Recbole [443, 444], FuxiCTR [458], DaisyRec [333, 334] and ELLIOT [9]. As for System deployers, they audit the system thoroughly based on the designed models before deployment. For example, Wilson et al. [374] propose a five-step (scoping, mapping, artifact collection, testing, and reflection) audit method to explore fairness problems in job recommendation systems. In detail, they analyze the source code and the data to explore some questions relevant to fairness. According to the methods mentioned above, we find that internal audits can minimize the probability of harmful impact before

deployment, but it may cause the problem of subjectivity because designers and auditors are the same groups of people.

7.2 Surveys and Tools

In this subsection, we will summarize the surveys and tools relevant to accountability and auditability.

7.2.1 Surveys. A survey [370] on algorithmic accountability summarizes the theory of accountability and organizes existing works according to five necessary parts of an accountable system [34]. As for auditability, one recent work [19] surveys algorithm audits from the aspects of behavior, domain, organization, and audit methods and focuses on the four problematic behaviors (discrimination, distortion, exploitation, and misjudgment) that auditors should pay attention to. Another survey [200] focuses on external audits. It first formulates the external audit process and then organizes existing works based on two proposed canonical audit forms. However, there is no specific survey about accountability and auditability for recommendations.

7.2.2 Tools. As mentioned above, annotation is one of the most important steps for external audits. Many researchers choose manual annotation for high accuracy, but it is a hard job. A crowdsourced platform is a good tool for this task, such as Amazon Mechanical Turk (AMT)¹. Besides, there are many benchmarking code library that can give convenience to auditors for reproducibility, e.g. Recbole² [443, 444], FuxiCTR³ [458], DaisyRec⁴ [333, 334] and ELLIOT⁵ [9].

7.3 Future Directions

As we know, the accountability of recommender systems is related to many aspects, such as explainability, fairness, and so on. However, most of the existing studies only focus on one aspect, which may lead to inadequate accountability. Therefore, the combination of many aspects for accountable recommender systems should be further considered. As for auditability, many works focus on external audits, but few on internal audits. To minimize the risk before deploying recommendation models, we can explore automated internal audits in the future. Furthermore, the collaboration between external and internal audits can be a promising direction because it can benefit from both merits.

8 INTERACTIONS AMONG DIFFERENT DIMENSIONS

The ideal trustworthy recommender systems would possess all of six features and advantages. However, in real-world context, it is challenging to consider the modeling of multiple features simultaneously, as these features may have many varying levels of interdependence, and even conflict in some aspects.

Despite that a number of studies have investigated the interactions between dimensions of trustworthy AI [101, 236, 399], research on trustworthy recommender systems is still limited. Fortunately, some researchers have recognized the importance of potential interactions among different dimensions, and attempted to explore and utilize them. In this section, we focus on the interactions between dimensions with extensive and close ties to other dimensions.

¹<https://www.mturk.com/>

²<https://github.com/RUCAIBox/RecBole>

³<https://github.com/xue-pai/FuxiCTR>

⁴<https://github.com/AmazingDD/daisyRec>

⁵<https://github.com/sisinflab/elliott>

Interactions with Robustness. Since the robustness of a system is an intrinsic characteristic to ensure its normal operation, it undoubtedly maintains a high degree of close crossover and connection with other dimensions in a recommendation system. Previous research on trustworthy AI [236] shows that the robustness of such systems is positively correlated to their explainability [96, 270], while partly conflicts with their privacy [325] and fairness dimensions [399]. These characteristics are particularly evident in adversarial attacks and robust training. For recommender systems, the issue is comparable. Therefore, how to use positive dimensions to promote robustness, and maintain the balance between robustness and potentially conflicting dimensions, while maintaining system's robustness against adversarial attacks without affecting other dimensions, is a non-trivial issue.

Most recently, Bilge et al. investigate the robustness of four recommendation algorithms based on collaborative filtering with privacy enhancement to determine, which improvement is more appropriate for the interaction between the two sides. Recent research on Trustworthy AI [201, 321] has shifted the focus of researchers to the interaction between robustness and other dimensions of recommendation systems within the context of machine learning. In [433], Zhang et al. design a robust model to combat the attacks and ensure the fairness of the recommender system. In [452], Zheng et al. develop an additive causal model for disentangling user interest and conformity for recommendation with causal embedding. In the meantime, this method ensures the robustness and explainability of the recommendation.

Interactions with Fairness. With the gradual expansion of the research on fairness in terms of trustworthy recommender systems and people's recognition of the importance of fairness, researchers began to take fairness as one of the goals when designing recommender systems, resulting in a number of studies on the interactions between fairness and other dimensions. One of the main directions is the interaction between fairness and explainability. In [49], Chen et al. provide a survey of the research on fairness and analyzes the explainability of the model at the same time. In [114], Fu et al. propose a fairness-aware explainable recommendation model. In [121], Ge et al. provide research on explainable fairness in recommendation. In addition, there are additional interaction studies with fairness, such as the interaction between robustness and fairness [433] mentioned in the previous subsection.

Interactions with Explainability. As mentioned in the previous section, the primary focus of interaction research in this regard is the interaction between explainability and fairness. In addition, there are several additional interaction studies. In [10], Anelli et al. analyze the robustness of their proposed explainable model. In [101], Fan et al. study the interactions between adversarial vulnerability and explainability, and take advantage of explainability to enhance adversarial attacks. In [126], Ghazimatin et al. provide a new counterfactual explanation mechanism for recommendation, which also solved the privacy exposure problem.

9 FUTURE DIRECTIONS

In this survey, we detail the existing works for trustworthy recommendations from six vital dimensions. However, there are also some other potential directions to be explored for the supplement to the definition of trustworthy recommender systems (TRec). In this section, we will illustrate some promising directions for further research on this topic.

Interactions among different dimensions. A trustworthy recommender system should possess six mentioned dimensions simultaneously, but most present researches only focus on one of them. Though few works have paid attention to achieving two dimensions, such as robust-fairness [433], explainability-privacy [126], etc., no one tries to go forward to three or more. Therefore, how to reach more requests of trustworthy dimensions is still an urgent problem for the community of recommender systems. Besides, the conflicts between different dimensions should not be ignored. For example, some recommendation works [320, 361] add an additional module for explainability,

which may bring the risk of violating the aspect of environmental well-being. Such conflicts may ruin the efforts for trustworthiness, so how to resolve the conflicts and get a trade-off is an important direction for TRec. In a word, both positive and negative interactions among different dimensions should be a spotlight in future directions.

Other Dimensions to achieve trustworthy recommender systems. Though we have issued six essential dimensions of TRec, there also exist some other dimensions worthy of being noted. For instance, security is a necessary dimension in many scenes, such as medication recommendation [340] and industrial recommendation [97]. In these scenes, the recommender systems will affect human decisions directly, and any improper decision can cause uncountable losses to life and property. Therefore, the characteristic of security is needed. Another aspect that TRec should possess is controllability. When a recommender system causes a devastating effect, accountability can only give out whether and who should be blamed. By comparison, controllability can help stop harmful recommendations and minimize the horrible effects. The two dimensions mentioned above are still far from the full content of TRec, and more other dimensions of TRec should be explored in the future.

Technology Ecosystem for trustworthy recommender systems. With the increasing demands of TRec, many researchers have devoted themselves to this field. However, the lack of a technology ecosystem causes huge inconvenience to the developments and experiments. A TRec technology ecosystem should contain datasets, metrics, toolkits, etc., but none of these parts have been well developed. For example, a few recent works aim to build up a standard dataset, such as KuaiRec [117] for non-discrimination, but no related work for environmental well-being and accountability specifically. Besides, there is no toolkit for evaluating various dimensions of the recommender system, which is one of the most important reasons why few research focuses on the interaction between various dimensions. Therefore, an integrated technology ecosystem is a vital procedure for achieving trustworthy recommender systems.

10 CONCLUSION

In this survey, we provide a comprehensive overview of trustworthy recommender systems (**TRec**) from a computational perspective. More specifically, we elaborate on six of the most critical dimensions for the trustworthiness of recommender systems: safety & robustness, non-discrimination & fairness, explainability, privacy, environmental well-being, and accountability & auditability. For each dimension, we provide the basic concepts and taxonomy for readers to have a better understanding of this topic, as well as summarize the representative methods in achieving trustworthy recommender systems. In addition, we introduce widely-used applications in real-world systems for achieving trustworthy recommender systems from multiple dimensions. Surveys and tools are also provided for readers' further exploration in this demanding topic. Finally, we also analyze the potential interactions among different dimensions and possible future research directions for trustworthy recommender systems.

REFERENCES

- [1] 2021. Europe fit for the Digital Age: Commission proposes new rules and actions for excellence and trust in Artificial Intelligence. https://ec.europa.eu/commission/presscorner/detail/en/IP_21_1682.
- [2] Himan Abdollahpouri and Masoud Mansoury. 2020. Multi-sided exposure bias in recommendation. *arXiv preprint arXiv:2006.15772* (2020).
- [3] Himan Abdollahpouri, Masoud Mansoury, Robin Burke, Bamshad Mobasher, and Edward Malthouse. 2021. User-centered evaluation of popularity bias in recommender systems. In *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization*. 119–129.
- [4] Muhammad Adnan, Yassaman Ebrahimzadeh Maboud, Divya Mahajan, and Prashant J Nair. 2021. Accelerating recommendation system training by leveraging popular choices. *Proceedings of the VLDB Endowment* 15, 1 (2021),

127–140.

- [5] Muhammad Adnan, Yassaman Ebrahimzadeh Maboud, Divya Mahajan, and Prashant J Nair. 2022. Heterogeneous Acceleration Pipeline for Recommendation System Training. *arXiv preprint arXiv:2204.05436* (2022).
- [6] Erfan Aghasian, Saurabh Garg, and James Montgomery. 2018. User’s Privacy in Recommendation Systems Applying Online Social Network Data, A Survey and Taxonomy. *arXiv preprint arXiv:1806.07629* (2018).
- [7] Qingyao Ai, Vahid Azizi, Xu Chen, and Yongfeng Zhang. 2018. Learning heterogeneous knowledge base embeddings for explainable recommendation. *Algorithms* 11, 9 (2018), 137.
- [8] Qingyao Ai and Lakshmi Narayanan. R. 2021. Model-agnostic vs. Model-intrinsic Interpretability for Explainable Product Search. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 5–15.
- [9] Vito Walter Anelli, Alejandro Bellogín, Antonio Ferrara, Daniele Malitesta, Felice Antonio Merra, Claudio Pomo, Francesco Maria Donini, and Tommaso Di Noia. 2021. Elliot: a comprehensive and rigorous framework for reproducible recommender systems evaluation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2405–2414.
- [10] Vito Walter Anelli, Tommaso Di Noia, Eugenio Di Sciascio, Azzurra Ragone, and Joseph Trotta. 2019. How to make latent factors interpretable by feeding factorization machines with knowledge graphs. In *International Semantic Web Conference*. Springer, 38–56.
- [11] Camila Souza Araújo, Gabriel Magno, Wagner Meira, Virgilio Almeida, Pedro Hartung, and Danilo Doneda. 2017. Characterizing videos, audience and advertising in Youtube channels for kids. In *International Conference on Social Informatics*. Springer, 341–359.
- [12] Vijay Arya, Rachel KE Bellamy, Pin-Yu Chen, Amit Dhurandhar, Michael Hind, Samuel C Hoffman, Stephanie Houde, Q Vera Liao, Ronny Luss, Aleksandra Mojsilovic, et al. 2020. Ai explainability 360: An extensible toolkit for understanding data and machine learning models. *Journal of Machine Learning Research* 21, 130 (2020), 1–6.
- [13] Ashwathy Ashokan and Christian Haas. 2021. Fairness metrics and bias mitigation strategies for rating predictions. *Information Processing & Management* 58, 5 (2021), 102646.
- [14] Giuseppe Ateniese, Giovanni Felici, Luigi V Mancini, Angelo Spognardi, Antonio Villani, and Domenico Vitali. 2013. Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers. *arXiv preprint arXiv:1306.4447* (2013).
- [15] Artem Babenko and Victor Lempitsky. 2014. Additive Quantization for Extreme Vector Compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [16] Ricardo Baeza-Yates. 2018. Bias on the web. *Commun. ACM* 61, 6 (2018), 54–61.
- [17] Krisztian Balog and Filip Radlinski. 2020. Measuring recommendation explanation quality: The conflicting goals of explanations. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 329–338.
- [18] Krisztian Balog, Filip Radlinski, and Shushan Arakelyan. 2019. Transparent, scrutible and explainable user models for personalized recommendation. In *Proceedings of the 42nd international acm sigir conference on research and development in information retrieval*. 265–274.
- [19] Jack Bandy. 2021. Problematic machine behavior: A systematic literature review of algorithm audits. *Proceedings of the acm on human-computer interaction* 5, CSCW1 (2021), 1–34.
- [20] Youjun Bao and Xiaohong Jiang. 2016. An intelligent medicine recommender system framework. In *2016 IEEE 11th conference on industrial electronics and applications (ICIEA)*. IEEE, 1383–1388.
- [21] Baptiste Barreau and Laurent Carlier. 2020. History-Augmented Collaborative Filtering for Financial Recommendations. In *Fourteenth ACM Conference on Recommender Systems*. 492–497.
- [22] Ghazaleh Beigi, Ahmadrza Mosallanezhad, Ruocheng Guo, Hamidreza Alvari, Alexander Nou, and Huan Liu. 2020. Privacy-aware recommendation with private-attribute protection using adversarial learning. In *Proceedings of the 13th International Conference on Web Search and Data Mining*.
- [23] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilović, et al. 2019. AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development* 63, 4/5 (2019), 4–1.
- [24] Alejandro Bellogín and Alan Said. 2021. Improving accountability in recommender systems research through reproducibility. *User Modeling and User-Adapted Interaction* (2021), 1–37.
- [25] Izak Benbasat and Weiquan Wang. 2005. Trust in and adoption of online recommendation agents. *Journal of the association for information systems* 6, 3 (2005), 4.
- [26] Alex Beutel, Jilin Chen, Tulsee Doshi, Hai Qian, Li Wei, Yi Wu, Lukasz Heldt, Zhe Zhao, Lichan Hong, Ed H Chi, et al. 2019. Fairness in recommendation ranking through pairwise comparisons. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2212–2220.

- [27] Runa Bhaumik, Bamshad Mobasher, and Robin Burke. 2011. A clustering approach to unsupervised attack detection in collaborative recommender systems. In *Proceedings of the International Conference on Data Science (ICDATA)*. Citeseer, 1.
- [28] Asia J Biega, Krishna P Gummadi, and Gerhard Weikum. 2018. Equity of attention: Amortizing individual fairness in rankings. In *The 41st international acm sigir conference on research & development in information retrieval*. 405–414.
- [29] Mustafa Bilgic and Raymond J Mooney. 2005. Explaining recommendations: Satisfaction vs. promotion. In *Beyond personalization workshop, IUI*, Vol. 5. 153.
- [30] Chen Bo, Zhao Xiangyu, Wang Yejing, Fan Wenqi, Guo Huifeng, and Ruiming Tang. 2022. Automated Machine Learning for Deep Recommender Systems: A Survey. (2022).
- [31] Rodrigo Borges and Kostas Stefanidis. 2019. Enhancing long term fairness in recommendations with variational autoencoders. In *Proceedings of the 11th international conference on management of digital ecosystems*. 95–102.
- [32] Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, et al. 2018. Google workloads for consumer devices: Mitigating data movement bottlenecks. In *Proceedings of the Twenty-Third International Conference on Architectural Support for Programming Languages and Operating Systems*. 316–331.
- [33] Avishek Bose and William Hamilton. 2019. Compositional fairness constraints for graph embeddings. In *International Conference on Machine Learning*. PMLR, 715–724.
- [34] Mark Bovens. 2007. Analysing and assessing accountability: A conceptual framework 1. *European law journal* 13, 4 (2007), 447–468.
- [35] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [36] Robin Burke. 2017. Multisided fairness for recommendation. *ArXiv preprint abs/1707.00093* (2017). <https://arxiv.org/abs/1707.00093>
- [37] Robin Burke, Bamshad Mobasher, and Runa Bhaumik. 2005. Limited knowledge shilling attacks in collaborative filtering systems. In *Proceedings of 3rd international workshop on intelligent techniques for web personalization (ITWP 2005), 19th international joint conference on artificial intelligence (IJCAI 2005)*. 17–24.
- [38] Robin Burke, Bamshad Mobasher, Runa Bhaumik, and Chad Williams. 2005. Segment-based injection attacks against collaborative filtering recommender systems. In *Fifth IEEE International Conference on Data Mining (ICDM'05)*. IEEE, 4–pp.
- [39] Robin Burke, Bamshad Mobasher, Chad Williams, and Runa Bhaumik. 2006. Classification features for attack detection in collaborative recommender systems. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. 542–547.
- [40] Robin Burke, Michael P O'Mahony, and Neil J Hurley. 2015. Robust collaborative recommendation. In *Recommender systems handbook*. Springer, 961–995.
- [41] Marina Buzzi. 2011. Children and YouTube: access to safe content. In *Proceedings of the 9th ACM SIGCHI Italian Chapter International Conference on Computer-Human Interaction: Facing Complexity*. 125–131.
- [42] Joseph A Calandrino, Ann Kilzer, Arvind Narayanan, Edward W Felten, and Vitaly Shmatikov. 2011. "You might also like:" Privacy risks of collaborative filtering. In *Proc. of SP*.
- [43] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. 2018. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097* (2018).
- [44] Alessandro Castelnovo, Riccardo Crupi, Greta Greco, Daniele Regoli, Ilaria Giuseppina Penco, and Andrea Claudio Cosentini. 2022. A clarification of the nuances in the fairness metrics landscape. *Scientific Reports* 12, 1 (2022), 1–21.
- [45] Chaochao Chen, Huiwen Wu, Jiajie Su, Lingjuan Lyu, Xiaolin Zheng, and Li Wang. 2022. Differential Private Knowledge Transfer for Privacy-Preserving Cross-Domain Recommendation. In *Proceedings of the ACM Web Conference 2022*.
- [46] Chaochao Chen, Jun Zhou, Bingzhe Wu, Wenjing Fang, Li Wang, Yuan Qi, and Xiaolin Zheng. 2020. Practical privacy preserving POI recommendation. *ACM Transactions on Intelligent Systems and Technology (TIST)* (2020).
- [47] Huiyuan Chen and Jing Li. 2019. Adversarial tensor factorization for context-aware recommendation. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 363–367.
- [48] Hongxu Chen, Yicong Li, Xiangguo Sun, Guandong Xu, and Hongzhi Yin. 2021. Temporal meta-path guided explainable recommendation. In *Proceedings of the 14th ACM international conference on web search and data mining*. 1056–1064.
- [49] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. 2020. Bias and debias in recommender system: A survey and future directions. *ArXiv preprint abs/2010.03240* (2020). <https://arxiv.org/abs/2010.03240>
- [50] Jingfan Chen, Wenqi Fan, Guanghui Zhu, Xiangyu Zhao, Chunfeng Yuan, Qing Li, and Yihua Huang. 2022. Knowledge-enhanced Black-box Attacks for Recommendations. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge*

Discovery and Data Mining. 108–117.

- [51] Junjie Chen, Wendy Hui Wang, and Xinghua Shi. 2020. Differential privacy protection against membership inference attack on machine learning for genomic data. In *BIOCOMPUTING 2021: Proceedings of the Pacific Symposium*.
- [52] Lei Chen, Fajie Yuan, Jiayi Yang, Min Yang, and Chengming Li. 2021. Scene-adaptive Knowledge Distillation for Sequential Recommendation via Differentiable Architecture Search. *arXiv preprint arXiv:2107.07173* (2021).
- [53] Rui Chen, Min Xie, and Laks VS Lakshmanan. 2014. Thwarting passive privacy attacks in collaborative filtering. In *Proc. of DASFAA*.
- [54] Tong Chen, Hongzhi Yin, Guanhua Ye, Zi Huang, Yang Wang, and Meng Wang. 2020. Try this instead: Personalized and interpretable substitute recommendation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 891–900.
- [55] Tong Chen, Hongzhi Yin, Xiangliang Zhang, Zi Huang, Yang Wang, and Meng Wang. 2021. Quaternion Factorization Machines: A Lightweight Solution to Intricate Feature Interaction Modeling. *IEEE Transactions on Neural Networks and Learning Systems* (2021).
- [56] Tong Chen, Hongzhi Yin, Yujia Zheng, Zi Huang, Yang Wang, and Meng Wang. 2021. Learning elastic embeddings for customizing on-device recommenders. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 138–147.
- [57] Xu Chen, Hanxiong Chen, Hongteng Xu, Yongfeng Zhang, Yixin Cao, Zheng Qin, and Hongyuan Zha. 2019. Personalized fashion recommendation with visual explanations based on multimodal attention network: Towards visually explainable recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 765–774.
- [58] Xu Chen, Yongfeng Zhang, and Zheng Qin. 2019. Dynamic explainable recommendation based on neural attentive models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 53–60.
- [59] Xu Chen, Yongfeng Zhang, and Ji-Rong Wen. 2022. Measuring "Why" in Recommender Systems: a Comprehensive Survey on the Evaluation of Explainable Recommendation. *arXiv preprint arXiv:2202.06466* (2022).
- [60] Xu Chen, Yongfeng Zhang, Hongteng Xu, Zheng Qin, and Hongyuan Zha. 2018. Adversarial distillation for efficient recommendation with external knowledge. *ACM Transactions on Information Systems (TOIS)* 37, 1 (2018), 1–28.
- [61] Yongjian Chen, Tao Guan, and Cheng Wang. 2010. Approximate nearest neighbor search by residual vector quantization. *Sensors* 10, 12 (2010), 11259–11273.
- [62] Yiran Chen, Yuan Xie, Linghao Song, Fan Chen, and Tianqi Tang. 2020. A survey of accelerator architectures for deep neural networks. *Engineering* 6, 3 (2020), 264–274.
- [63] Zhongxia Chen, Xiting Wang, Xing Xie, Mehul Parsana, Akshay Soni, Xiang Ao, and Enhong Chen. 2021. Towards explainable conversational recommendation. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*. 2994–3000.
- [64] Zhongxia Chen, Xiting Wang, Xing Xie, Tong Wu, Guoqing Bu, Yining Wang, and Enhong Chen. 2019. Co-Attentive Multi-Task Learning for Explainable Recommendation. In *IJCAI*. 2137–2143.
- [65] Zeyuan Chen, Jiangchao Yao, Feng Wang, Kunyang Jia, Bo Han, Wei Zhang, and Hongxia Yang. 2021. MC2-SF: Slow-Fast Learning for Mobile-Cloud Collaborative Recommendation. *arXiv preprint arXiv:2109.12314* (2021).
- [66] Weiyu Cheng, Yanyan Shen, and Linpeng Huang. 2020. Differentiable neural input search for recommender systems. *arXiv preprint arXiv:2006.04466* (2020).
- [67] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. 2017. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:1710.09282* (2017).
- [68] Zhiyong Cheng, Xiaojun Chang, Lei Zhu, Rose C Kanjirathinkal, and Mohan Kankanhalli. 2019. MMALFM: Explainable recommendation by leveraging reviews and images. *ACM Transactions on Information Systems (TOIS)* 37, 2 (2019), 1–28.
- [69] Zhiyong Cheng, Ying Ding, Lei Zhu, and Mohan Kankanhalli. 2018. Aspect-aware latent factor model: Rating prediction with ratings and reviews. In *Proceedings of the 2018 world wide web conference*. 639–648.
- [70] Alexandra Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5, 2 (2017), 153–163.
- [71] Alexandra Chouldechova and Aaron Roth. 2018. The frontiers of fairness in machine learning. *arXiv preprint arXiv:1810.08810* (2018).
- [72] Konstantina Christakopoulou and Arindam Banerjee. 2018. Adversarial recommendation: Attack of the learned fake users. *arXiv preprint arXiv:1809.08336* (2018).
- [73] Konstantina Christakopoulou and Arindam Banerjee. 2019. Adversarial attacks on an oblivious recommender. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 322–330.
- [74] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. 2017. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*. 797–806.

- [75] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*. 191–198.
- [76] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. M6-Rec: Generative Pretrained Language Models are Open-Ended Recommender Systems. *arXiv preprint arXiv:2205.08084* (2022).
- [77] Enyan Dai, Tianxiang Zhao, Huaisheng Zhu, Junjie Xu, Zhimeng Guo, Hui Liu, Jiliang Tang, and Suhang Wang. 2022. A Comprehensive Survey on Trustworthy Graph Neural Networks: Privacy, Robustness, Fairness, and Explainability. *arXiv preprint arXiv:2204.08570* (2022).
- [78] Guohao Dai, Zhenhua Zhu, Tianyu Fu, Chiyue Wei, Bangyan Wang, Xiangyu Li, Yuan Xie, Huazhong Yang, and Yu Wang. 2022. DIMMining: pruning-efficient and parallel graph mining on near-memory-computing. In *Proceedings of the 49th Annual International Symposium on Computer Architecture*. 130–145.
- [79] Pierre Danhier, Clément Massart, and François-Xavier Standaert. 2020. Fidelity leakages: Applying membership inference attacks to preference data. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPs)*.
- [80] Abhinandan S Das, Mayur Datar, Ashutosh Garg, and Shyam Rajaram. 2007. Google news personalization: scalable online collaborative filtering. In *Proceedings of the 16th international conference on World Wide Web*. 271–280.
- [81] Sanjoy Dasgupta and Kaushik Sinha. 2013. Randomized partition trees for exact nearest neighbor search. In *Conference on learning theory*. PMLR, 317–337.
- [82] Yashar Deldjoo, Tommaso Di Noia, and Felice Antonio Merra. 2020. Adversarial machine learning in recommender systems (aml-recsys). In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 869–872.
- [83] Yashar Deldjoo, Tommaso Di Noia, and Felice Antonio Merra. 2021. A survey on adversarial recommender systems: from attack/defense strategies to generative adversarial networks. *ACM Computing Surveys (CSUR)* 54, 2 (2021), 1–38.
- [84] Lei Deng, Guoqi Li, Song Han, Luping Shi, and Yuan Xie. 2020. Model compression and hardware acceleration for neural networks: A comprehensive survey. *Proc. IEEE* 108, 4 (2020), 485–532.
- [85] Tyler Derr, Yao Ma, Wenqi Fan, Xiaorui Liu, Charu Aggarwal, and Jiliang Tang. 2020. Epidemic graph convolutional network. In *Proceedings of the 13th International Conference on Web Search and Data Mining (WSDM)*. 160–168.
- [86] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [87] Virginie Do, Sam Corbett-Davies, Jamal Atif, and Nicolas Usunier. 2021. Two-sided fairness in rankings via Lorenz dominance. *Advances in Neural Information Processing Systems* 34 (2021), 8596–8608.
- [88] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017).
- [89] Finale Doshi-Velez, Mason Kortz, Ryan Budish, Chris Bavitz, Sam Gershman, David O’Brien, Kate Scott, Stuart Schieber, James Waldo, David Weinberger, et al. 2017. Accountability of AI under the law: The role of explanation. *arXiv preprint arXiv:1711.01134* (2017).
- [90] Vasisht Duddu, Antoine Boutet, and Virat Shejwalkar. 2020. Quantifying privacy leakage in graph embedding. In *MobiQuitous 2020-17th EAI International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services*.
- [91] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2016. Calibrating noise to sensitivity in private data analysis. *Journal of Privacy and Confidentiality*.
- [92] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* (2014).
- [93] Assaf Eisenman, Maxim Naumov, Darryl Gardner, Misha Smelyanskiy, Sergey Pupyrev, Kim Hazelwood, Asaf Cidon, and Sachin Katti. 2019. Bandana: Using non-volatile memory for storing deep learning models. *Proceedings of Machine Learning and Systems* 1 (2019), 40–52.
- [94] Michael D Ekstrand, Anubrata Das, Robin Burke, and Fernando Diaz. 2021. Fairness and discrimination in information access systems. *ArXiv preprint abs/2105.05779* (2021). <https://arxiv.org/abs/2105.05779>
- [95] Michael D Ekstrand, Mucun Tian, Ion Madrazo Azpiazu, Jennifer D Ekstrand, Oghenemaro Anuyah, David McNeill, and Maria Soledad Pera. 2018. All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. In *Conference on fairness, accountability and transparency*. PMLR, 172–186.
- [96] Christian Etmann, Sebastian Lunz, Peter Maass, and Carola-Bibiane Schönlieb. 2019. On the connection between adversarial robustness and saliency map interpretability. *arXiv preprint arXiv:1905.04172* (2019).
- [97] Sopuruchukwu Ezenwa, Abhijit D Talpade, Pushkar Ghanekar, Ravi Joshi, Jayachandran Devaraj, Fabio H Ribeiro, and Ray Mentzer. 2022. Toward Improved Safety Culture in Academic and Industrial Chemical Laboratories: An Assessment and Recommendation of Best Practices. *ACS Chemical Health & Safety* 29, 2 (2022), 202–213.
- [98] Miao Fan, Jiacheng Guo, Shuai Zhu, Shuo Miao, Mingming Sun, and Ping Li. 2019. MOBIUS: towards the next generation of query-ad matching in baidu’s sponsored search. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2509–2517.

- [99] Wenqi Fan, Tyler Derr, Yao Ma, Jianping Wang, Jiliang Tang, and Qing Li. 2019. Deep Adversarial Social Recommendation. In *28th International Joint Conference on Artificial Intelligence (IJCAI-19)*. 1351–1357.
- [100] Wenqi Fan, Tyler Derr, Xiangyu Zhao, Yao Ma, Hui Liu, Jianping Wang, Jiliang Tang, and Qing Li. 2021. Attacking black-box recommendations via copying cross-domain user profiles. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 1583–1594.
- [101] Wenqi Fan, Wei Jin, Xiaorui Liu, Han Xu, Xianfeng Tang, Suhang Wang, Qing Li, Jiliang Tang, Jianping Wang, and Charu Aggarwal. 2021. Jointly Attacking Graph Neural Network and its Explanations. *arXiv preprint arXiv:2108.03388* (2021).
- [102] Wenqi Fan, Qing Li, and Min Cheng. 2018. Deep modeling of social relations for recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [103] Wenqi Fan, Xiaorui Liu, Wei Jin, Xiangyu Zhao, Jiliang Tang, and Qing Li. 2022. Graph Trend Filtering Networks for Recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 112–121.
- [104] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. 2019. Graph neural networks for social recommendation. In *The world wide web conference*. 417–426.
- [105] Wenqi Fan, Yao Ma, Qing Li, Jianping Wang, Guoyong Cai, Jiliang Tang, and Dawei Yin. 2020. A graph neural network framework for social recommendations. *IEEE Transactions on Knowledge and Data Engineering* (2020).
- [106] Wenqi Fan, Yao Ma, Dawei Yin, Jianping Wang, Jiliang Tang, and Qing Li. 2019. Deep social collaborative filtering. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 305–313.
- [107] Boli Fang, Miao Jiang, Pei-yi Cheng, Jerry Shen, and Yi Fang. 2020. Achieving Outcome Fairness in Machine Learning Models for Social Decision Problems. In *IJCAI*. 444–450.
- [108] Minghong Fang, Neil Zhenqiang Gong, and Jia Liu. 2020. Influence function based data poisoning attacks to top-n recommender systems. In *Proceedings of The Web Conference 2020*. 3019–3025.
- [109] Minghong Fang, Guolei Yang, Neil Zhenqiang Gong, and Jia Liu. 2018. Poisoning attacks to graph-based recommender systems. In *Proceedings of the 34th annual computer security applications conference*. 381–392.
- [110] Golnoosh Farnadi, Pigi Kouki, Spencer K Thompson, Sriram Srinivasan, and Lise Getoor. 2018. A fairness-aware hybrid recommender system. *arXiv preprint arXiv:1809.09030* (2018).
- [111] Adrian Flanagan, Were Oyomno, Alexander Grigorievskiy, Kuan E Tan, Suleiman A Khan, and Muhammad Ammad-Ud-Din. 2020. Federated multi-view matrix factorization for personalized recommendations. In *Proc. of ECML*.
- [112] Will Fleisher. 2021. What’s Fair about Individual Fairness?. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 480–490.
- [113] Jerome H Friedman, Jon Louis Bentley, and Raphael Ari Finkel. 1977. An algorithm for finding best matches in logarithmic expected time. *ACM Transactions on Mathematical Software (TOMS)* 3, 3 (1977), 209–226.
- [114] Zuohui Fu, Yikun Xian, Ruoyuan Gao, Jieyu Zhao, Qiaoying Huang, Yingqiang Ge, Shuyuan Xu, Shijie Geng, Chirag Shah, Yongfeng Zhang, et al. 2020. Fairness-aware explainable recommendation over knowledge graphs. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 69–78.
- [115] Karan Ganju, Qi Wang, Wei Yang, Carl A Gunter, and Nikita Borisov. 2018. Property inference attacks on fully connected neural networks using permutation invariant representations. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*.
- [116] Chen Gao, Chao Huang, Yue Yu, Huandong Wang, Yong Li, and Depeng Jin. 2019. Privacy-preserving cross-domain location recommendation. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* (2019).
- [117] Chongming Gao, Shijun Li, Wenqiang Lei, Biao Li, Peng Jiang, Jiawei Chen, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. 2022. KuaiRec: A Fully-observed Dataset for Recommender Systems. *arXiv preprint arXiv:2202.10842* (2022).
- [118] Jianling Gao, Lingtao Qi, Haiping Huang, and Chao Sha. 2020. Shilling attack detection scheme in collaborative filtering recommendation system based on recurrent neural network. In *Future of Information and Communication Conference*. Springer, 634–644.
- [119] Yingqiang Ge, Shuchang Liu, Zuohui Fu, Juntao Tan, Zelong Li, Shuyuan Xu, Yunqi Li, Yikun Xian, and Yongfeng Zhang. 2022. A survey on trustworthy recommender systems. *arXiv preprint arXiv:2207.12515* (2022).
- [120] Yingqiang Ge, Shuchang Liu, Ruoyuan Gao, Yikun Xian, Yunqi Li, Xiangyu Zhao, Changhua Pei, Fei Sun, Junfeng Ge, Wenwu Ou, et al. 2021. Towards long-term fairness in recommendation. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 445–453.
- [121] Yingqiang Ge, Juntao Tan, Yan Zhu, Yinglong Xia, Jiebo Luo, Shuchang Liu, Zuohui Fu, Shijie Geng, Zelong Li, and Yongfeng Zhang. 2022. Explainable Fairness in Recommendation. *arXiv preprint arXiv:2204.11159* (2022).
- [122] Yingqiang Ge, Xiaoting Zhao, Lucia Yu, Saurabh Paul, Diane Hu, Chu-Cheng Hsieh, and Yongfeng Zhang. 2022. Toward Pareto Efficient Fairness-Utility Trade-off in Recommendation through Reinforcement Learning. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 316–324.

- [123] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as Language Processing (RLP): A Unified Pretrain, Personalized Prompt & Predict Paradigm (P5). *arXiv preprint arXiv:2203.13366* (2022).
- [124] Sahin Cem Geyik, Stuart Ambler, and Krishnaram Kenthapadi. 2019. Fairness-aware ranking in search & recommendation systems with application to linkedin talent search. In *Proceedings of the 25th acm sigkdd international conference on knowledge discovery & data mining*. 2221–2231.
- [125] Hossein Movafegh Ghadirli, Ali Nodehi, and Rasul Enayatifar. 2019. An Overview of Encryption Algorithms in Color Images. *Signal Processing* (2019).
- [126] Azin Ghazimatin, Oana Balalau, Rishiraj Saha Roy, and Gerhard Weikum. 2020. PRINCE: Provider-side interpretability with counterfactual explanations in recommender systems. In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 196–204.
- [127] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. 2021. A survey of quantization methods for efficient neural network inference. *arXiv preprint arXiv:2103.13630* (2021).
- [128] Antonio A Ginart, Maxim Naumov, Dheevatsa Mudigere, Jiyang Yang, and James Zou. 2021. Mixed dimension embeddings with application to memory-efficient recommendation systems. In *2021 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2786–2791.
- [129] Aristides Gionis, Piotr Indyk, Rajeev Motwani, et al. 1999. Similarity search in high dimensions via hashing. In *Vldb*, Vol. 99. 518–529.
- [130] Yu Gong, Ziwen Jiang, Yufei Feng, Binbin Hu, Kaiqi Zhao, Qingwen Liu, and Wenwu Ou. 2020. EdgeRec: recommender system on edge in Mobile Taobao. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2477–2484.
- [131] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision* 129, 6 (2021), 1789–1819.
- [132] Mihajlo Grbovic and Haibin Cheng. 2018. Real-time personalization using embeddings for search ranking at airbnb. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 311–320.
- [133] Nina Grgić-Hlača, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. 2018. Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [134] Ihsan Gunes, Cihan Kaleli, Alper Bilge, and Huseyin Polat. 2014. Shilling attacks against recommender systems: a comprehensive survey. *Artificial Intelligence Review* 42, 4 (2014), 767–799.
- [135] Huifeng Guo, Wei Guo, Yong Gao, Ruiming Tang, Xiuqiang He, and Wenzhi Liu. 2021. Scalefreectr: Mixcache-based distributed training system for ctr models with huge embedding table. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1269–1278.
- [136] Weiwei Guo, Xiaowei Liu, Sida Wang, Huiji Gao, Ananth Sankar, Zimeng Yang, Qi Guo, Liang Zhang, Bo Long, Bee-Chung Chen, et al. 2020. Detext: A deep text ranking framework with bert. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 2509–2516.
- [137] Yunhui Guo, Mohsen Imani, Jaeyoung Kang, Sahand Salamat, Justin Morris, Baris Aksanli, Yeseong Kim, and Tajana Rosing. 2021. Hyperrec: Efficient recommender systems with hyperdimensional computing. In *2021 26th Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 384–389.
- [138] Yeting Guo, Fang Liu, Zhiping Cai, Hui Zeng, Li Chen, Tongqing Zhou, and Nong Xiao. 2021. PREFER: Point-of-interest REcommendation with efficiency and privacy-preservation via Federated Edge leaRning. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* (2021).
- [139] Ananya Gupta, Eric Johnson, Justin Payan, Aditya Kumar Roy, Ari Kobren, Swetasudha Panda, Jean-Baptiste Tristan, and Michael Wick. 2021. Online post-processing in rankings for fair utility maximization. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 454–462.
- [140] Udit Gupta, Samuel Hsia, Vikram Saraph, Xiaodong Wang, Brandon Reagan, Gu-Yeon Wei, Hsien-Hsin S Lee, David Brooks, and Carole-Jean Wu. 2020. Deeprecsys: A system for optimizing end-to-end at-scale neural recommendation inference. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 982–995.
- [141] Udit Gupta, Samuel Hsia, Jeff Zhang, Mark Wilkening, Javin Pombra, Hsien-Hsin Sean Lee, Gu-Yeon Wei, Carole-Jean Wu, and David Brooks. 2021. RecPipe: Co-designing models and hardware to jointly optimize recommendation quality and performance. In *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture*. 870–884.
- [142] Mengyue Hang, Tobias Schnabel, Longqi Yang, and Jennifer Neville. 2022. Lightweight Compositional Embeddings for Incremental Streaming Recommendation. *arXiv preprint arXiv:2202.02427* (2022).
- [143] Chaoyang He, Songze Li, Jinhyun So, Xiao Zeng, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Praneeth Vepakomma, Abhishek Singh, Hang Qiu, et al. 2020. Fedml: A research library and benchmark for federated machine learning. *arXiv preprint arXiv:2007.13518* (2020).

- [144] Xiangnan He, Tao Chen, Min-Yen Kan, and Xiao Chen. 2015. Trirank: Review-aware explainable recommendation by modeling aspects. In *Proceedings of the 24th ACM international on conference on information and knowledge management*. 1661–1670.
- [145] Xiangnan He, Zhankui He, Xiaoyu Du, and Tat-Seng Chua. 2018. Adversarial personalized ranking for recommendation. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 355–364.
- [146] Yun He, Xue Feng, Cheng Cheng, Geng Ji, Yunsong Guo, and James Caverlee. 2022. MetaBalance: Improving Multi-Task Recommendations via Adapting Gradient Magnitudes of Auxiliary Tasks. In *Proceedings of the ACM Web Conference 2022*. 2205–2215.
- [147] Yingzhe He, Guozhu Meng, Kai Chen, Xingbo Hu, and Jinwen He. 2019. Towards privacy and security of deep learning systems: a survey. *arXiv preprint arXiv:1.12562* (2019).
- [148] Anna Hedström, Leander Weber, Dilyara Bareeva, Franz Motzkus, Wojciech Samek, Sebastian Lapuschkin, and Marina M-C Höhne. 2022. Quantus: an explainable AI toolkit for responsible evaluation of neural network explanations. *arXiv preprint arXiv:2202.06861* (2022).
- [149] Kartik Hegde, Hadi Asghari-Moghaddam, Michael Pellauer, Neal Crago, Aamer Jaleel, Edgar Solomonik, Joel Emer, and Christopher W Fletcher. 2019. Extensor: An accelerator for sparse tensor algebra. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*. 319–333.
- [150] Jonathan L Herlocker, Joseph A Konstan, and John Riedl. 2000. Explaining collaborative filtering recommendations. In *Proceedings of the 2000 ACM conference on Computer supported cooperative work*. 241–250.
- [151] Seira Hidano, Takao Murakami, Shuichi Katsumata, Shinsaku Kiyomoto, and Goichiro Hanaoka. 2020. Exposing Private User Behaviors of Collaborative Filtering via Model Inversion Techniques. *Proceedings on Privacy Enhancing Technologies* (2020).
- [152] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).
- [153] T. Ryan Hoens, Marina Blanton, Aaron Steele, and Nitesh V. Chawla. 2013. Reliable Medical Recommendation Systems with Patient Privacy. *ACM Transactions on Intelligent Systems and Technology* 4 (2013), 67:1–67:31.
- [154] Naoise Holohan, Stefano Braghin, Pól Mac Aonghusa, and Killian Levacher. 2019. Diffprivlib: the IBM differential privacy library. *arXiv preprint arXiv:1907.02444* (2019).
- [155] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. 2022. Towards Universal Sequence Representation Learning for Recommender Systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 585–593.
- [156] Brian Hu, Paul Tunison, Bhavan Vasu, Nitesh Menon, Roddy Collins, and Anthony Hoogs. 2021. XAITK: The explainable AI toolkit. *Applied AI Letters* 2, 4 (2021), e40.
- [157] Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. 2021. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)* (2021).
- [158] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-based retrieval in facebook search. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2553–2561.
- [159] Ling Huang, Anthony D Joseph, Blaine Nelson, Benjamin IP Rubinstein, and J Doug Tygar. 2011. Adversarial machine learning. In *Proceedings of the 4th ACM workshop on Security and artificial intelligence*.
- [160] Mingkai Huang, Hao Li, Bing Bai, Chang Wang, Kun Bai, and Fei Wang. 2020. A federated multi-view deep learning framework for privacy-preserving recommendations. *arXiv preprint arXiv:2008.10808* (2020).
- [161] Weiming Huang, Baisong Liu, and Hao Tang. 2019. Privacy protection for recommendation system: a survey. In *Journal of Physics: Conference Series*.
- [162] Wen Huang, Lu Zhang, and Xintao Wu. 2022. Achieving Counterfactual Fairness for Causal Bandit. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 6952–6959.
- [163] Yongfeng Huo, Bilian Chen, Jing Tang, and Yifeng Zeng. 2021. Privacy-Preserving Point-of-Interest Recommendation Based on Geographical and Social Influence. *Information Sciences* (2021).
- [164] Ranggi Hwang, Taehun Kim, Youngeun Kwon, and Minsoo Rhu. 2020. Centaur: A chiplet-based, hybrid sparse-dense accelerator for personalized recommendations. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 968–981.
- [165] Mohamed Assem Ibrahim, Onur Kayiran, and Shaizeen Aga. 2021. Efficient Cache Utilization via Model-aware Data Placement for Recommendation Models. In *The International Symposium on Memory Systems*. 1–11.
- [166] Eugene Ie, Vihan Jain, Jing Wang, Sanmit Narvekar, Ritesh Agarwal, Rui Wu, Heng-Tze Cheng, Tushar Chandra, and Craig Boutilier. 2019. SlateQ: A tractable decomposition for reinforcement learning with recommendation sets. (2019).

- [167] Rashidul Islam, Kamrun Naher Keya, Ziqian Zeng, Shimei Pan, and James Foulds. 2021. Debiasing career recommendations with neural fair collaborative filtering. In *Proceedings of the Web Conference 2021*. 3779–3790.
- [168] Priyank Jain, Manasi Gyanchandani, and Nilay Khare. 2016. Big Data Privacy: A Technological Perspective and Review. *Journal of Big Data* 3 (2016), 1–25.
- [169] Dietmar Jannach, Sidra Naveed, and Michael Jugovac. 2016. User control in recommender systems: Overview and interaction challenges. In *International Conference on Electronic Commerce and Web Technologies*. Springer, 21–33.
- [170] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2010. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence* 33, 1 (2010), 117–128.
- [171] Sérgio Jesus, Catarina Belém, Vladimir Balayan, João Bento, Pedro Saleiro, Pedro Bizarro, and João Gama. 2021. How can I choose an explainer? An application-grounded evaluation of post-hoc explanations. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 805–815.
- [172] Biye Jiang, Chao Deng, Huimin Yi, Zelin Hu, Guorui Zhou, Yang Zheng, Sui Huang, Xinyang Guo, Dongyue Wang, Yue Song, et al. 2019. XDL: an industrial deep learning framework for high-dimensional sparse data. In *Proceedings of the 1st International Workshop on Deep Learning Practice for High-Dimensional Sparse Data*. 1–9.
- [173] Gangwei Jiang, Hao Wang, Jin Chen, Haoyu Wang, Defu Lian, and Enhong Chen. 2021. xLightFM: Extremely Memory-Efficient Factorization Machine. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 337–346.
- [174] Di Jin, Elena Sergeeva, Wei-Hung Weng, Geeticka Chauhan, and Peter Szolovits. 2022. Explainable deep learning in healthcare: A methodological survey from an attribution view. *WIREs Mechanisms of Disease* 14, 3 (2022), e1548.
- [175] Manas R Joglekar, Cong Li, Mei Chen, Taibai Xu, Xiaoming Wang, Jay K Adams, Pranav Khaitan, Jiahui Liu, and Quoc V Le. 2020. Neural input search for large scale recommendation models. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2387–2397.
- [176] Brittany Johnson and Yuriy Brun. 2022. Fairkit-learn: A Fairness Evaluation and Comparison Toolkit. (2022).
- [177] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [178] Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. 2016. Rawlsian fairness for machine learning. *arXiv preprint arXiv:1610.09559* 1, 2 (2016), 19.
- [179] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. 2021. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning* (2021).
- [180] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Issei Sato. 2016. Model-based approaches for independence-enhanced recommendation. In *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*. IEEE, 860–867.
- [181] Dongyeop Kang, Anusha Balakrishnan, Pararth Shah, Paul Crook, Y-Lan Boureau, and Jason Weston. 2019. Recommendation as a communication game: Self-supervised bot-play for goal-oriented dialogue. *arXiv preprint arXiv:1909.03922* (2019).
- [182] SeongKu Kang, Junyoung Hwang, Wonbin Kweon, and Hwanjo Yu. 2020. DE-RRD: A knowledge distillation framework for recommender system. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 605–614.
- [183] SeongKu Kang, Junyoung Hwang, Wonbin Kweon, and Hwanjo Yu. 2021. Topology distillation for recommender system. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 829–839.
- [184] Wang-Cheng Kang, Derek Zhiyuan Cheng, Tiansheng Yao, Xinyang Yi, Ting Chen, Lichan Hong, and Ed H Chi. 2021. Learning to Embed Categorical Features without Embedding Tables for Recommendation. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 840–850.
- [185] Chen Karako and Putra Manggala. 2018. Using image fairness representations in diversity-based re-ranking for recommendations. In *Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization*. 23–28.
- [186] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6769–6781.
- [187] P Karthikeyan, S Thamarai Selvi, G Neeraja, R Deepika, A Vincent, and V Abinaya. 2017. Prevention of shilling attack in recommender systems using discrete wavelet transform and support vector machine. In *2016 eighth international conference on Advanced Computing (ICoAC)*. IEEE, 99–104.
- [188] Stefan Katzenbeisser and Milan Petkovic. 2008. Privacy-Preserving Recommendation Systems for Consumer Healthcare Services. In *2008 Third International Conference on Availability, Reliability and Security*. 889–895.
- [189] Mesut Kaya, Derek Bridge, and Nava Tintarev. 2020. Ensuring fairness in group recommendations by rank-sensitive balancing of relevance. In *Fourteenth ACM Conference on Recommender Systems*. 101–110.

- [190] Liu Ke, Udit Gupta, Benjamin Youngjae Cho, David Brooks, Vikas Chandra, Utku Diril, Amin Firoozshahian, Kim Hazelwood, Bill Jia, Hsien-Hsin S Lee, et al. 2020. Recnmp: Accelerating personalized recommendation with near-memory processing. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 790–803.
- [191] Omid Keivani, Kaushik Sinha, and Parikshit Ram. 2018. Improved maximum inner product search with better theoretical guarantee using randomized partition trees. *Machine Learning* 107, 6 (2018), 1069–1094.
- [192] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2016. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807* (2016).
- [193] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2016. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236* (2016).
- [194] Wonbin Kweon, SeongKu Kang, and Hwanjo Yu. 2021. Bidirectional distillation for top-K recommender system. In *Proceedings of the Web Conference 2021*. 3861–3871.
- [195] Youngeun Kwon, Yunjae Lee, and Minsoo Rhu. 2019. Tensordimm: A practical near-memory processing architecture for embeddings and tensor operations in deep learning. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*. 740–753.
- [196] Youngeun Kwon, Yunjae Lee, and Minsoo Rhu. 2021. Tensor casting: Co-designing algorithm-architecture for personalized recommendation training. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 235–248.
- [197] Youngeun Kwon and Minsoo Rhu. 2022. Training personalized recommendation systems from (GPU) scratch: look forward not backwards. In *Proceedings of the 49th Annual International Symposium on Computer Architecture*. 860–873.
- [198] Shyong K Lam and John Riedl. 2004. Shilling recommender systems for fun and profit. In *Proceedings of the 13th international conference on World Wide Web*. 393–402.
- [199] Dung D Le and Hady Lauw. 2021. Efficient Retrieval of Matrix Factorization-Based Top-k Recommendations: A Survey of Recent Approaches. *Journal of Artificial Intelligence Research* 70 (2021), 1441–1479.
- [200] Erwan Le Merrer, Ronan Pons, and Gilles Trédan. 2022. Algorithmic audits of algorithms, and the law. (2022).
- [201] Jae-Gil Lee, Yuji Roh, Hwanjun Song, and Steven Euijong Whang. 2021. Machine learning robustness, fairness, and their convergence. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 4046–4047.
- [202] Jong-Seok Lee and Dan Zhu. 2012. Shilling attack detection—a new approach for a trustworthy recommender system. *INFORMS Journal on Computing* 24, 1 (2012), 117–131.
- [203] Jae-woong Lee, Minjin Choi, Jongwuk Lee, and Hyunjung Shim. 2019. Collaborative distillation for top-N recommendation. In *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 369–378.
- [204] Min Kyung Lee, Anuraag Jain, Hea Jin Cha, Shashank Ojha, and Daniel Kusbit. 2019. Procedural justice in algorithmic fairness: Leveraging transparency and outcome control for fair algorithmic mediation. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–26.
- [205] Bo Li, Yining Wang, Aarti Singh, and Yevgeniy Vorobeychik. 2016. Data poisoning attacks on factorization-based collaborative filtering. *Advances in neural information processing systems* 29 (2016).
- [206] Chaozhuo Li, Bochen Pang, Yuming Liu, Hao Sun, Zheng Liu, Xing Xie, Tianqi Yang, Yanling Cui, Liangjie Zhang, and Qi Zhang. 2021. Adsgnn: Behavior-graph augmented relevance modeling in sponsored search. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 223–232.
- [207] Jie Li, Yongli Ren, and Ke Deng. 2022. FairGAN: GANs-based Fairness-aware Learning for Recommendations with Implicit Feedback. In *Proceedings of the ACM Web Conference 2022*. 297–307.
- [208] Kaiyang Li, Guangchun Luo, Yang Ye, Wei Li, Shihao Ji, and Zhipeng Cai. 2020. Adversarial privacy-preserving graph embedding against inference attack. *IEEE Internet of Things Journal* (2020).
- [209] Lei Li, Dingding Wang, Tao Li, Daniel Knox, and Balaji Padmanabhan. 2011. Scene: a scalable two-stage personalized news recommendation system. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. 125–134.
- [210] Lei Li, Yongfeng Zhang, and Li Chen. 2020. Generate neural template explanations for recommendation. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 755–764.
- [211] Lei Li, Yongfeng Zhang, and Li Chen. 2021. Personalized transformer for explainable recommendation. *arXiv preprint arXiv:2105.11601* (2021).
- [212] Lei Li, Yongfeng Zhang, and Li Chen. 2022. Personalized prompt learning for explainable recommendation. *arXiv preprint arXiv:2202.07371* (2022).
- [213] Muyang Li, Xiangyu Zhao, Chuan Lyu, Minghao Zhao, Runze Wu, and Ruocheng Guo. 2022. MLP4Rec: A Pure MLP Architecture for Sequential Recommendations. *arXiv preprint arXiv:2204.11510* (2022).
- [214] Piji Li, Zihao Wang, Zhaochun Ren, Lidong Bing, and Wai Lam. 2017. Neural rating regression with abstractive tips generation for recommendation. In *Proceedings of the 40th International ACM SIGIR conference on Research and*

Development in Information Retrieval. 345–354.

- [215] Roger Zhe Li, Julián Urbano, and Alan Hanjalic. 2021. Leave no user behind: Towards improving the utility of recommender systems for non-mainstream users. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 103–111.
- [216] Sen Li, Fuyu Lv, Taiwei Jin, Guli Lin, Keping Yang, Xiaoyi Zeng, Xiao-Ming Wu, and Qianli Ma. 2021. Embedding-based product retrieval in taobao search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 3181–3189.
- [217] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. 2020. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine* (2020).
- [218] Tan Li, Linqi Song, and Christina Fragouli. 2020. Federated recommendation system via differential privacy. In *2020 IEEE International Symposium on Information Theory (ISIT)*. 2592–2597.
- [219] Yunqi Li, Hanxiong Chen, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2021. User-oriented fairness in recommendation. In *Proceedings of the Web Conference 2021*. 624–632.
- [220] Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, Juntao Tan, Shuchang Liu, and Yongfeng Zhang. 2022. Fairness in Recommendation: A Survey. *arXiv preprint arXiv:2205.13619* (2022).
- [221] Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. 2021. Towards personalized fairness based on causal notion. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1054–1063.
- [222] Yang Li, Tong Chen, Peng-Fei Zhang, and Hongzhi Yin. 2021. Lightweight self-attentive sequential recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 967–977.
- [223] Yunqi Li, Yingqiang Ge, and Yongfeng Zhang. 2021. Tutorial on fairness of machine learning in recommender systems. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2654–2657.
- [224] Yangkun Li, Mohamed-Laid Hedia, Weizhi Ma, Hongyu Lu, Min Zhang, Yiqun Liu, and Shaoping Ma. 2022. Contextualized Fairness for Recommender Systems in Premium Scenarios. *Big Data Research* 27 (2022), 100300.
- [225] Zeyu Li, Wei Cheng, Haiqi Xiao, Wenchao Yu, Haifeng Chen, and Wei Wang. 2021. You Are What and Where You Are: Graph Enhanced Attention Network for Explainable POI Recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3945–3954.
- [226] Defu Lian, Haoyu Wang, Zheng Liu, Jianxun Lian, Enhong Chen, and Xing Xie. 2020. Lightrec: A memory and search-efficient recommender system. In *Proceedings of The Web Conference 2020*. 695–705.
- [227] Defu Lian, Xing Xie, and Enhong Chen. 2019. Discrete matrix factorization and extension for fast item recommendation. *IEEE Transactions on Knowledge and Data Engineering* 33, 5 (2019), 1919–1933.
- [228] Defu Lian, Xing Xie, Enhong Chen, and Hui Xiong. 2020. Product quantized collaborative filtering. *IEEE Transactions on Knowledge and Data Engineering* 33, 9 (2020), 3284–3296.
- [229] Peiyuan Liao, Han Zhao, Keyulu Xu, Tommi Jaakkola, Geoffrey J Gordon, Stefanie Jegelka, and Ruslan Salakhutdinov. 2021. Information obfuscation of graph neural networks. In *Proc. of ICML*.
- [230] Chen Lin, Si Chen, Hui Li, Yanghua Xiao, Lianyun Li, and Qian Yang. 2020. Attacking recommender systems with augmented user profiles. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 855–864.
- [231] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
- [232] Weilin Lin, Xiangyu Zhao, Yejing Wang, Tong Xu, and Xian Wu. 2022. AdaFS: Adaptive Feature Selection in Deep Recommender System. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3309–3317.
- [233] Yujie Lin, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Jun Ma, and Maarten De Rijke. 2019. Explainable outfit recommendation with joint outfit matching and comment generation. *IEEE Transactions on Knowledge and Data Engineering* 32, 8 (2019), 1502–1516.
- [234] Chong Liu, Defu Lian, Min Nie, and Xia Hu. 2020. Online optimized product quantization. In *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 362–371.
- [235] Haochen Liu, Da Tang, Ji Yang, Xiangyu Zhao, Hui Liu, Jiliang Tang, and Youlong Cheng. 2022. Rating Distribution Calibration for Selection Bias Mitigation in Recommendations. In *Proceedings of the ACM Web Conference 2022*. 2048–2057.
- [236] Haochen Liu, Yiqi Wang, Wenqi Fan, Xiaorui Liu, Yaxin Li, Shaili Jain, Yunhao Liu, Anil K Jain, and Jiliang Tang. 2022. Trustworthy ai: A computational perspective. *ACM Transactions on Intelligent Systems and Technology (TIST)* (2022).
- [237] Haochen Liu, Xiangyu Zhao, Chong Wang, Xiaobing Liu, and Jiliang Tang. 2020. Automated embedding size search in deep recommender systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2307–2316.

- [238] Jie Liu, Xiao Yan, Xinyan Dai, Zhirong Li, James Cheng, and Ming-Chang Yang. 2020. Understanding and improving proximity graph based maximum inner product search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 139–146.
- [239] Kumpeng Liu, Yanjie Fu, Le Wu, Xiaolin Li, Charu Aggarwal, and Hui Xiong. 2021. Automated feature selection: A reinforcement learning perspective. *IEEE Transactions on Knowledge and Data Engineering* (2021).
- [240] Ninghao Liu, Yong Ge, Li Li, Xia Hu, Rui Chen, and Soo-Hyun Choi. 2020. Explainable recommender systems via resolving learning representations. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 895–904.
- [241] Qi Liu, Jin Zhang, Defu Lian, Yong Ge, Jianhui Ma, and Enhong Chen. 2021. Online Additive Quantization. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1098–1108.
- [242] Siyi Liu, Chen Gao, Yihong Chen, Depeng Jin, and Yong Li. 2020. Learnable Embedding sizes for Recommender Systems. In *International Conference on Learning Representations*.
- [243] Weiwen Liu, Jun Guo, Nasim Sonboli, Robin Burke, and Shengyu Zhang. 2019. Personalized fairness-aware re-ranking for microlending. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 467–471.
- [244] Weiwen Liu, Feng Liu, Ruiming Tang, Ben Liao, Guangyong Chen, and Pheng Ann Heng. 2020. Balancing between accuracy and fairness for interactive recommendation with reinforcement learning. In *Pacific-asia conference on knowledge discovery and data mining*. Springer, 155–167.
- [245] Michele Loi and Matthias Spielkamp. 2021. Towards accountability in the use of artificial intelligence for public administrations. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 757–766.
- [246] Michael Lui, Yavuz Yetim, Özgür Özkan, Zhuoran Zhao, Shin-Yeh Tsai, Carole-Jean Wu, and Mark Hempstead. 2021. Understanding capacity-driven scale-out neural recommendation inference. In *2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 162–171.
- [247] Xiao Luo, Haixin Wang, Daqing Wu, Chong Chen, Minghua Deng, Jianqiang Huang, and Xian-Sheng Hua. 2020. A survey on deep hashing methods. *ACM Transactions on Knowledge Discovery from Data (TKDD)* (2020).
- [248] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkitasubramaniam. 2006. L-Diversity: Privacy beyond k-Anonymity. In *22nd International Conference on Data Engineering (ICDE'06)*.
- [249] Masoud Mansoury, Himan Abdollahpour, Mykola Pechenizkiy, Bamshad Mobasher, and Robin Burke. 2020. Feedback Loop and Bias Amplification in Recommender Systems. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, Mathieu d'Aquin, Stefan Dietze, Claudia Hauff, Edward Curry, and Philippe Cudré-Mauroux (Eds.). ACM, 2145–2148. <https://doi.org/10.1145/3340531.3412152>
- [250] Masoud Mansoury, Himan Abdollahpour, Mykola Pechenizkiy, Bamshad Mobasher, and Robin Burke. 2021. A graph-based approach for mitigating multi-sided exposure bias in recommender systems. *ACM Transactions on Information Systems (TOIS)* 40, 2 (2021), 1–31.
- [251] Benjamin Marlin, Richard S Zemel, Sam Roweis, and Malcolm Slaney. 2012. Collaborative filtering and the missing at random assumption. *arXiv preprint arXiv:1206.5267* (2012).
- [252] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* 54, 6 (2021), 1–35.
- [253] Rishabh Mehrotra, James McInerney, Hugues Bouchard, Mounia Lalmas, and Fernando Diaz. 2018. Towards a fair marketplace: Counterfactual evaluation of the trade-off between relevance, fairness & satisfaction in recommendation systems. In *Proceedings of the 27th acm international conference on information and knowledge management*. 2243–2251.
- [254] Bhaskar Mehta. 2007. Unsupervised shilling detection for collaborative filtering. In *AAAI*. 1402–1407.
- [255] Bhaskar Mehta and Thomas Hofmann. 2008. A Survey of Attack-Resistant Collaborative Filtering Algorithms. *IEEE Data Eng. Bull.* 31, 2 (2008), 14–22.
- [256] Bhaskar Mehta and Wolfgang Nejdl. 2009. Unsupervised strategies for shilling detection and robust collaborative filtering. *User Modeling and User-Adapted Interaction* 19, 1 (2009), 65–97.
- [257] Xuying Meng, Suhang Wang, Kai Shu, Jundong Li, Bo Chen, Huan Liu, and Yujun Zhang. 2019. Towards Privacy Preserving Social Recommendation under Personalized Privacy Settings. *World Wide Web* (2019).
- [258] David J Miller, Zhen Xiang, and George Kesidis. 2020. Adversarial learning targeting deep neural network classification: A comprehensive review of defenses against attacks. *Proc. IEEE* (2020).
- [259] Fatemehsadat Mireshghallah, Mohammadkazem Taram, Praneeth Vepakomma, Abhishek Singh, Ramesh Raskar, and Hadi Esmailzadeh. 2020. Privacy in deep learning: A survey. *arXiv preprint arXiv:2004.12254* (2020).
- [260] Sparsh Mittal and Shrayish Vaishay. 2019. A survey of techniques for optimizing deep learning on GPUs. *Journal of Systems Architecture* 99 (2019), 101635.
- [261] Bamshad Mobasher, Robin Burke, Runa Bhaumik, and Chad Williams. 2007. Toward trustworthy recommender systems: An analysis of attack models and algorithm robustness. *ACM Transactions on Internet Technology (TOIT)* 7, 4 (2007), 23–es.

- [262] Marco Morik, Ashudeep Singh, Jessica Hong, and Thorsten Joachims. 2020. Controlling fairness and bias in dynamic learning-to-rank. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 429–438.
- [263] Stanislav Morozov and Artem Babenko. 2018. Non-metric similarity graphs for maximum inner product search. *Advances in Neural Information Processing Systems* 31 (2018).
- [264] PETER MÜLLNER, ELISABETH LEX, and DOMINIK KOWALD. 2022. ReuseKNN: Neighborhood Reuse for Differentially-Private KNN-Based Recommendations. *ACM Trans. Intell. Syst. Technol* 1, 1 (2022).
- [265] Mohammadmehdi Naghiaei, Hossein A Rahmani, and Yashar Deldjoo. 2022. Cpfair: Personalized consumer and producer fairness re-ranking for recommender systems. *arXiv preprint arXiv:2204.08085* (2022).
- [266] Arvind Narayanan. 2018. Translation tutorial: 21 fairness definitions and their politics. In *Proc. Conf. Fairness Accountability Transp., New York, USA*, Vol. 1170. 3.
- [267] Priyanka Nigam, Yiwei Song, Vijai Mohan, Vihan Lakshman, Weitian Ding, Ankit Shingavi, Choon Hui Teo, Hao Gu, and Bing Yin. 2019. Semantic product search. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2876–2885.
- [268] Mehrbakhsh Nilashi, Dietmar Jannach, Othman bin Ibrahim, Mohammad Dalvi Esfahani, and Hossein Ahmadi. 2016. Recommendation quality, transparency, and website quality for trust-building in recommendation agents. *Electronic Commerce Research and Applications* 19 (2016), 70–84.
- [269] Hamed Nilforoshan, Johann D Gaebler, Ravi Shroff, and Sharad Goel. 2022. Causal conceptions of fairness and their consequences. In *International Conference on Machine Learning*. PMLR, 16848–16887.
- [270] Adam Noack, Isaac Ahern, Dejing Dou, and Boyang Li. 2021. An empirical study on the relation between network interpretability and adversarial robustness. *SN Computer Science* 2, 1 (2021), 1–13.
- [271] Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kiciman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in Big Data* 2 (2019), 13.
- [272] Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. 2019. Knockoff nets: Stealing functionality of black-box models. In *Proc. of CVPR*.
- [273] Zohreh Ovaisi, Shelby Heinecke, Jia Li, Yongfeng Zhang, Elena Zheleva, and Caiming Xiong. 2022. RGRecSys: A Toolkit for Robustness Evaluation of Recommender Systems. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1597–1600.
- [274] Kostantinos Papadamou, Antonis Papasavva, Savvas Zannettou, Jeremy Blackburn, Nicolas Kourtellis, Ilias Leontiadis, Gianluca Stringhini, and Michael Sirivianos. 2020. Disturbed YouTube for kids: Characterizing and detecting inappropriate videos targeting young children. In *Proceedings of the international AAAI conference on web and social media*, Vol. 14. 522–533.
- [275] Kostantinos Papadamou, Savvas Zannettou, Jeremy Blackburn, Emiliano De Cristofaro, Gianluca Stringhini, and Michael Sirivianos. 2022. “It is just a flu”: Assessing the Effect of Watch History on YouTube’s Pseudoscientific Video Recommendations. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 16. 723–734.
- [276] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
- [277] Mathias PM Parisot, Balazs Pejo, and Dayana Spagnuolo. 2021. Property Inference Attacks on Convolutional Neural Networks: Influence and Implications of Target Model’s Complexity. *arXiv preprint arXiv:2104.13061* (2021).
- [278] Haekyu Park, Hyunsik Jeon, Junghwan Kim, Beunguk Ahn, and U Kang. 2017. Uniwalk: Explainable and accurate recommendation for rating and network data. *arXiv preprint arXiv:1710.07134* (2017).
- [279] Gourab K Patro, Arpita Biswas, Niloy Ganguly, Krishna P Gummedi, and Abhijnan Chakraborty. 2020. Fairrec: Two-sided fairness for personalized recommendations in two-sided platforms. In *Proceedings of The Web Conference 2020*. 1194–1204.
- [280] Georgina Peake and Jun Wang. 2018. Explanation mining: Post hoc interpretability of latent factor models for recommendation systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2060–2069.
- [281] Cong Peng, Debiao He, Jianhua Chen, Neeraj Kumar, and Muhammad Khurram Khan. 2021. EPRT: an efficient privacy-preserving medical service recommendation and trust discovery scheme for eHealth system. *ACM Transactions on Internet Technology (TOIT)* (2021).
- [282] Lillian Pentecost, Marco Donato, Brandon Reagen, Udit Gupta, Siming Ma, Gu-Yeon Wei, and David Brooks. 2019. Maxnm: Maximizing dnn storage density and inference efficiency with sparse encoding and error mitigation. In *Proceedings of the 52Nd Annual IEEE/ACM International Symposium on Microarchitecture*. 769–781.
- [283] Evaggelia Pitoura, Kostas Stefanidis, and Georgia Koutrika. 2021. Fairness in rankings and recommendations: an overview. *The VLDB Journal* (2021), 1–28.

- [284] Tao Qi, Fangzhao Wu, Chuhan Wu, Yongfeng Huang, and Xing Xie. 2020. Privacy-preserving news recommendation model learning. *arXiv preprint arXiv:2003.09592* (2020).
- [285] Tao Qi, Fangzhao Wu, Chuhan Wu, Peijie Sun, Le Wu, Xiting Wang, Yongfeng Huang, and Xing Xie. 2022. ProFairRec: Provider Fairness-aware News Recommendation. *arXiv preprint arXiv:2204.04724* (2022).
- [286] Inioluwa Deborah Raji, Andrew Smart, Rebecca N White, Margaret Mitchell, Timnit Gebru, Ben Hutchinson, Jamila Smith-Loud, Daniel Theron, and Parker Barnes. 2020. Closing the AI accountability gap: defining an end-to-end framework for internal algorithmic auditing. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 33–44.
- [287] Parikshit Ram and Kaushik Sinha. 2019. Revisiting kd-tree for nearest neighbor search. In *Proceedings of the 25th acm sigkdd international conference on knowledge discovery & data mining*. 1378–1388.
- [288] N. Ramakrishnan, B.J. Keller, B.J. Mirza, A.Y. Grama, and G. Karypis. 2001. Privacy Risks in Recommender Systems. *IEEE Internet Computing* (2001).
- [289] Bashir Rastegarpanah, Krishna P Gummadi, and Mark Crovella. 2019. Fighting fire with fire: Using antidote data to improve polarization and fairness of recommender systems. In *Proceedings of the twelfth ACM international conference on web search and data mining*. 231–239.
- [290] Brandon Reagen, Paul Whatmough, Robert Adolf, Saketh Rama, Hyunkwang Lee, Sae Kyu Lee, José Miguel Hernández-Lobato, Gu-Yeon Wei, and David Brooks. 2016. Minerva: Enabling low-power, highly-accurate deep neural network accelerators. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 267–278.
- [291] Zhaochun Ren, Shangsong Liang, Piji Li, Shuaiqiang Wang, and Maarten de Rijke. 2017. Social collaborative viewpoint regression with explainable recommendations. In *Proceedings of the tenth ACM international conference on web search and data mining*. 485–494.
- [292] Steffen Rendle. 2010. Factorization machines. In *2010 IEEE International conference on data mining*. IEEE, 995–1000.
- [293] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).
- [294] Steffen Rendle and Lars Schmidt-Thieme. 2010. Pairwise interaction tensor factorization for personalized tag recommendation. In *Proceedings of the third ACM international conference on Web search and data mining*. 81–90.
- [295] Yehezkel S Resheff, Yanai Elazar, Moni Shahar, and Oren Sar Shalom. 2018. Privacy and fairness in recommender systems via adversarial training of user representations. *arXiv preprint arXiv:1807.03521* (2018).
- [296] Manoel Horta Ribeiro, Raphael Ottoni, Robert West, Virgílio AF Almeida, and Wagner Meira Jr. 2020. Auditing radicalization pathways on YouTube. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 131–141.
- [297] Maria Rigaki and Sebastian Garcia. 2020. A survey of privacy attacks in machine learning. *arXiv preprint arXiv:2007.07646* (2020).
- [298] Crefeda Faviola Rodrigues, Graham Riley, and Mikel Luján. 2018. SyNERGY: An energy measurement and prediction framework for Convolutional Neural Networks on Jetson TX1. In *Proceedings of the International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA)*. The Steering Committee of The World Congress in Computer Science, Computer . . . , 375–382.
- [299] Benjamin IP Rubinstein and A Francesco. 2017. diffpriv: An R package for easy differential privacy. *Journal of Machine Learning Research* (2017).
- [300] Dimitris Sacharidis. 2019. Top-n group recommendations with fairness. In *Proceedings of the 34th ACM/SIGAPP symposium on applied computing*. 1663–1670.
- [301] Bishal Sainju, Chris Hartwell, and John Edwards. 2021. Job satisfaction and employee turnover determinants in Fortune 50 companies: Insights from employee reviews from Indeed. com. *Decision Support Systems* 148 (2021), 113582.
- [302] Salman Salamatian, Amy Zhang, Flavio du Pin Calmon, Sandilya Bhamidipati, Nadia Fawaz, Branislav Kveton, Pedro Oliveira, and Nina Taft. 2015. Managing your private and public data: Bringing down inference attacks against your privacy. *IEEE Journal of Selected Topics in Signal Processing* (2015).
- [303] Ahmed Salem, Apratim Bhattacharya, Michael Backes, Mario Fritz, and Yang Zhang. 2020. {Updates-Leak}: Data Set Inference and Reconstruction Attacks in Online Learning. In *Proc. of USENIX Security*.
- [304] Fatemeh Sarvi, Maria Heuss, Mohammad Aliannejadi, Sebastian Schelter, and Maarten de Rijke. 2022. Understanding and Mitigating the Effect of Outliers in Fair Ranking. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 861–869.
- [305] Ryoma Sato. 2022. Enumerating Fair Packages for Group Recommendations. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 870–878.
- [306] Dimitris Serbos, Shuyao Qi, Nikos Mamoulis, Evaggelia Pitoura, and Panayiotis Tsaparas. 2017. Fairness in package-to-group recommendations. In *Proceedings of the 26th international conference on world wide web*. 371–379.

- [307] Joan Serrà and Alexandros Karatzoglou. 2017. Getting deep recommenders fit: Bloom embeddings for sparse binary input/output networks. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. 279–287.
- [308] Geet Sethi, Bilge Acun, Niket Agarwal, Christos Kozyrakis, Caroline Trippel, and Carole-Jean Wu. 2022. RecShard: statistical feature-based memory optimization for industry-scale neural recommendation. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*. 344–358.
- [309] Behzad Shahrabi, Venugopal Mani, Apoorv Reddy Arrabothu, Deepthi Sharma, Kannan Achan, and Sushant Kumar. 2020. On Detecting Data Pollution Attacks On Recommender Systems Using Sequential GANs. *arXiv preprint arXiv:2012.02509* (2020).
- [310] Amit Sharma and Dan Cosley. 2013. Do social explanations work? Studying and modeling the effects of social explanations in recommender systems. In *Proceedings of the 22nd international conference on World Wide Web*. 1133–1144.
- [311] Jiayi Shen, Haotao Wang, Shupeng Gui, Jianchao Tan, Zhangyang Wang, and Ji Liu. 2020. UMEC: Unified model and embedding compression for efficient recommendation systems. In *International Conference on Learning Representations*.
- [312] Hao-Jun Michael Shi, Dheevatsa Mudigere, Maxim Naumov, and Jiyang Yang. 2020. Compositional embeddings using complementary partitions for memory-efficient recommendation systems. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 165–175.
- [313] Shaoyun Shi, Weizhi Ma, Min Zhang, Yongfeng Zhang, Xinxing Yu, Houzhi Shan, Yiqun Liu, and Shaoping Ma. 2020. Beyond user embedding matrix: Learning to hash for modeling large-scale users in recommendation. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 319–328.
- [314] Kyuyong Shin, Hanock Kwak, Kyung-Min Kim, Minkyu Kim, Young-Jin Park, Jisu Jeong, and Seungjae Jung. 2021. One4all user representation for recommender systems in e-commerce. *arXiv preprint arXiv:2106.00573* (2021).
- [315] Dorin Shmaryahu, Guy Shani, and Bracha Shapira. 2020. Post-hoc Explanations for Complex Model Recommendations using Simple Methods.. In *IntrS@ RecSys*. 26–36.
- [316] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In *Proc. of SP*.
- [317] Mingdan Si and Qingshan Li. 2020. Shilling attacks against collaborative recommender systems: a review. *Artificial Intelligence Review* 53, 1 (2020), 291–319.
- [318] Ashudeep Singh and Thorsten Joachims. 2018. Fairness of exposure in rankings. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2219–2228.
- [319] Jaspreet Singh and Avishek Anand. 2018. Posthoc interpretability of learning to rank models using secondary training data. *arXiv preprint arXiv:1806.11330* (2018).
- [320] Jaspreet Singh and Avishek Anand. 2020. Model agnostic interpretability of rankers via intent modelling. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 618–628.
- [321] Edward Small, Wei Shao, Zeliang Zhang, Peihan Liu, Jeffrey Chan, Kacper Sokol, and Flora Salim. 2022. How Robust is your Fair Model? Exploring the Robustness of Diverse Fairness Strategies. *arXiv preprint arXiv:2207.04581* (2022).
- [322] Nathalie A Smuha. 2019. The EU approach to ethics guidelines for trustworthy artificial intelligence. *Computer Law Review International* 20, 4 (2019), 97–106.
- [323] Nasim Sonboli, Farzad Eskandarian, Robin Burke, Weiwen Liu, and Bamshad Mobasher. 2020. Opportunistic multi-aspect fairness through personalized re-ranking. In *Proceedings of the 28th ACM Conference on User Modeling, Adaptation and Personalization*. 239–247.
- [324] Junshuai Song, Zhao Li, Zehong Hu, Yucheng Wu, Zhenpeng Li, Jian Li, and Jun Gao. 2020. Poisonrec: an adaptive data poisoning framework for attacking black-box recommender systems. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. IEEE, 157–168.
- [325] Liwei Song, Reza Shokri, and Prateek Mittal. 2019. Privacy risks of securing machine learning models against adversarial examples. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 241–257.
- [326] Qingquan Song, Dehua Cheng, Hanning Zhou, Jiyang Yang, Yuandong Tian, and Xia Hu. 2020. Towards automated neural interaction discovery for click-through rate prediction. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 945–955.
- [327] Weiping Song, Zhijian Duan, Ziqing Yang, Hao Zhu, Ming Zhang, and Jian Tang. 2019. Explainable knowledge graph-based recommendation via deep reinforcement learning. *arXiv preprint arXiv:1906.09506* (2019).
- [328] Harald Steck. 2018. Calibrated recommendations. In *Proceedings of the 12th ACM conference on recommender systems*. 154–162.
- [329] Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. 2021. Whitening sentence representations for better semantics and faster retrieval. *arXiv preprint arXiv:2103.15316* (2021).
- [330] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM*

- international conference on information and knowledge management*. 1441–1450.
- [331] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. Ernie: Enhanced representation through knowledge integration. *arXiv preprint arXiv:1904.09223* (2019).
- [332] Yang Sun, Fajie Yuan, Min Yang, Guoao Wei, Zhou Zhao, and Duo Liu. 2020. A generic network compression framework for sequential recommender systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1299–1308.
- [333] Zhu Sun, Hui Fang, Jie Yang, Xinghua Qu, Hongyang Liu, Di Yu, Yew-Soon Ong, and Jie Zhang. 2022. DaisyRec 2.0: Benchmarking Recommendation for Rigorous Evaluation. *arXiv preprint arXiv:2206.10848* (2022).
- [334] Zhu Sun, Di Yu, Hui Fang, Jie Yang, Xinghua Qu, Jie Zhang, and Cong Geng. 2020. Are we evaluating rigorously? benchmarking recommendation for reproducible evaluation and fair comparison. In *Fourteenth ACM conference on recommender systems*. 23–32.
- [335] Özge Sürer, Robin Burke, and Edward C Malthouse. 2018. Multistakeholder recommendation with provider constraints. In *Proceedings of the 12th ACM Conference on Recommender Systems*. 54–62.
- [336] Latanya Sweeney. 2002. k-anonymity: A model for protecting privacy. *International journal of uncertainty, fuzziness and knowledge-based systems* (2002).
- [337] Kyosuke Takami, Yiling Dai, Brendan Flanagan, and Hiroaki Ogata. 2022. Educational Explainable Recommender Usage and its Effectiveness in High School Summer Vacation Assignment. In *LAK22: 12th International Learning Analytics and Knowledge Conference*. 458–464.
- [338] Juntao Tan, Shuyuan Xu, Yingqiang Ge, Yunqi Li, Xu Chen, and Yongfeng Zhang. 2021. Counterfactual explainable recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 1784–1793.
- [339] Shulong Tan, Zhaozhuo Xu, Weijie Zhao, Hongliang Fei, Zhixin Zhou, and Ping Li. 2021. Norm Adjusted Proximity Graph for Fast Inner Product Retrieval. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1552–1560.
- [340] Yanchao Tan, Chengjun Kong, Leisheng Yu, Pan Li, Chaochao Chen, Xiaolin Zheng, Vicki S Hertzberg, and Carl Yang. 2022. 4SDrug: Symptom-based Set-to-set Small and Safe Drug Recommendation. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3970–3980.
- [341] Jinhui Tang, Xiaoyu Du, Xiangnan He, Fajie Yuan, Qi Tian, and Tat-Seng Chua. 2019. Adversarial training towards robust multimedia recommender system. *IEEE Transactions on Knowledge and Data Engineering* 32, 5 (2019), 855–867.
- [342] Jiayi Tang and Ke Wang. 2018. Ranking distillation: Learning compact ranking models with high performance for recommender system. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 2289–2298.
- [343] Jiayi Tang, Hongyi Wen, and Ke Wang. 2020. Revisiting adversarially learned injection attacks against recommender systems. In *Fourteenth ACM conference on recommender systems*. 318–327.
- [344] Nava Tintarev and Judith Masthoff. 2011. Designing and evaluating explanations for recommender systems. In *Recommender systems handbook*. Springer, 479–510.
- [345] Nava Tintarev and Judith Masthoff. 2015. Explaining recommendations: Design and evaluation. In *Recommender systems handbook*. Springer, 353–382.
- [346] Anh Truong, Negar Kiyavash, and Seyed Rasoul Etesami. 2018. Adversarial machine learning: The case of recommendation systems. In *2018 IEEE 19th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 1–5.
- [347] Chun-Hua Tsai and Peter Brusilovsky. 2019. Evaluating visual explanations for similarity-based recommendations: User perception and performance. In *Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization*. 22–30.
- [348] Michael Veale, Max Van Kleek, and Reuben Binns. 2018. Fairness and accountability design needs for algorithmic support in high-stakes public sector decision-making. In *Proceedings of the 2018 chi conference on human factors in computing systems*. 1–14.
- [349] Sahil Verma, John Dickerson, and Keegan Hines. 2020. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596* (2020).
- [350] Hu Wan, Xuan Sun, Yufei Cui, Chia-Lin Yang, Tei-Wei Kuo, and Chun Jason Xue. 2021. FlashEmbedding: storing embedding tables in SSD for large-scale recommender systems. In *Proceedings of the 12th ACM SIGOPS Asia-Pacific Workshop on Systems*. 9–16.
- [351] Mengting Wan, Jianmo Ni, Rishabh Misra, and Julian McAuley. 2020. Addressing marketing bias in product recommendations. In *Proceedings of the 13th international conference on web search and data mining*. 618–626.
- [352] Binghui Wang, Jiayi Guo, Ang Li, Yiran Chen, and Hai Li. 2021. Privacy-preserving representation learning on graphs: A mutual information perspective. In *Proc. of KDD*.

- [353] Chao Wang, Hengshu Zhu, Chen Zhu, Xi Zhang, Enhong Chen, and Hui Xiong. 2020. Personalized employee training course recommendation with career development awareness. In *Proceedings of the Web Conference 2020*. 1648–1659.
- [354] Feng Wang, Miaomiao Dai, Xudong Li, and Liquan Pan. 2021. Compressing Embedding Table via Multi-dimensional Quantization Encoding for Sequential Recommender Model. In *2021 the 7th International Conference on Communication and Information Processing (ICCIP)*. 234–239.
- [355] Jianfang Wang and Pengfei Han. 2019. Adversarial training-based mean Bayesian personalized ranking for recommender system. *IEEE Access* 8 (2019), 7958–7968.
- [356] Jie Wang, Fajie Yuan, Mingyue Cheng, Joemon M Jose, Chenyun Yu, Beibei Kong, Zhijin Wang, Bo Hu, and Zang Li. 2022. TransRec: Learning Transferable Recommendation from Mixture-of-Modality Feedback. *arXiv preprint arXiv:2206.06190* (2022).
- [357] Jingdong Wang, Ting Zhang, Nicu Sebe, Heng Tao Shen, et al. 2017. A survey on learning to hash. *IEEE transactions on pattern analysis and machine intelligence* 40, 4 (2017), 769–790.
- [358] Qinyong Wang, Hongzhi Yin, Tong Chen, Zi Huang, Hao Wang, Yanchang Zhao, and Nguyen Quoc Viet Hung. 2020. Next point-of-interest recommendation on resource-constrained mobile devices. In *Proceedings of the Web conference 2020*. 906–916.
- [359] Shoujin Wang, Xiuzhen Zhang, Yan Wang, Huan Liu, and Francesco Ricci. 2022. Trustworthy Recommender Systems. *arXiv preprint arXiv:2208.06265* (2022).
- [360] Weiqi Wang, An Liu, Zhixu Li, Xiangliang Zhang, Qing Li, and Xiaofang Zhou. 2019. Protecting Multi-Party Privacy in Location-Aware Social Point-of-Interest Recommendation. *World Wide Web* (2019).
- [361] Xiting Wang, Yiru Chen, Jie Yang, Le Wu, Zhengtao Wu, and Xing Xie. 2018. A reinforcement learning framework for explainable recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 587–596.
- [362] Xiang Wang, Xiangnan He, Fuli Feng, Liqiang Nie, and Tat-Seng Chua. 2018. Tem: Tree-enhanced embedding model for explainable recommendation. In *Proceedings of the 2018 World Wide Web Conference*. 1543–1552.
- [363] Xiting Wang, Kunpeng Liu, Dongjie Wang, Le Wu, Yanjie Fu, and Xing Xie. 2022. Multi-level recommendation reasoning over knowledge graphs with reinforcement learning. In *Proceedings of the ACM Web Conference 2022*. 2098–2108.
- [364] Xiang Wang, Dingxian Wang, Canran Xu, Xiangnan He, Yixin Cao, and Tat-Seng Chua. 2019. Explainable reasoning over knowledge graphs for recommendation. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 5329–5336.
- [365] Yifan Wang, Weizhi Ma, Min Zhang, Yiqun Liu, and Shaoping Ma. 2022. A Survey on the Fairness of Recommender Systems. *ArXiv preprint abs/2206.03761* (2022). <https://arxiv.org/abs/2206.03761>
- [366] Yejing Wang, Xiangyu Zhao, Tong Xu, and Xian Wu. 2022. Autofield: Automating feature selection in deep recommender systems. In *Proceedings of the ACM Web Conference 2022*. 1977–1986.
- [367] Yitu Wang, Zhenhua Zhu, Fan Chen, Mingyuan Ma, Guohao Dai, Yu Wang, Hai Li, and Yiran Chen. 2021. REREC: In-ReRAM Acceleration with Access-Aware Mapping for Personalized Recommendation. In *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*. IEEE, 1–9.
- [368] Kilian Weinberger, Anirban Dasgupta, John Langford, Alex Smola, and Josh Attenberg. 2009. Feature hashing for large scale multitask learning. In *Proceedings of the 26th annual international conference on machine learning*. 1113–1120.
- [369] James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viégas, and Jimbo Wilson. 2019. The what-if tool: Interactive probing of machine learning models. *IEEE transactions on visualization and computer graphics* 26, 1 (2019), 56–65.
- [370] Maranke Wieringa. 2020. What to account for when accounting for algorithms: A systematic literature review on algorithmic accountability. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 1–18.
- [371] Mark Wilkening, Udit Gupta, Samuel Hsia, Caroline Trippel, Carole-Jean Wu, David Brooks, and Gu-Yeon Wei. 2021. RecSSD: near data processing for solid state drive based recommendation inference. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*. 717–729.
- [372] Chad Williams and Bamshad Mobasher. 2006. Profile injection attack detection for securing collaborative recommender systems. *DePaul University CTI Technical Report* (2006), 1–47.
- [373] Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* 8, 3 (1992), 229–256.
- [374] Christo Wilson, Avijit Ghosh, Shan Jiang, Alan Mislove, Lewis Baker, Janelle Szary, Kelly Trindel, and Frida Polli. 2021. Building and auditing fair algorithms: A case study in candidate screening. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 666–677.
- [375] Chenwang Wu, Defu Lian, Yong Ge, Zhihao Zhu, and Enhong Chen. 2021. Triple Adversarial Learning for Influence based Poisoning Attack in Recommender Systems. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1830–1840.

- [376] Chuhan Wu, Fangzhao Wu, Yang Cao, Yongfeng Huang, and Xing Xie. 2021. Fedgnn: Federated graph neural network for privacy-preserving recommendation. *arXiv preprint arXiv:2102.04925* (2021).
- [377] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2022. Are Big Recommendation Models Fair to Cold Users? *arXiv preprint arXiv:2202.13607* (2022).
- [378] Chuhan Wu, Fangzhao Wu, Tao Qi, Yongfeng Huang, and Xing Xie. 2022. FedCL: Federated Contrastive Learning for Privacy-Preserving Recommendation. *arXiv preprint arXiv:2204.09850* (2022).
- [379] Chuhan Wu, Fangzhao Wu, Xiting Wang, Yongfeng Huang, and Xing Xie. 2021. Fairness-aware news recommendation with decomposed adversarial learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 4462–4469.
- [380] Le Wu, Lei Chen, Pengyang Shao, Richang Hong, Xiting Wang, and Meng Wang. 2021. Learning fair representations for recommendation: A graph-based perspective. In *Proceedings of the Web Conference 2021*. 2198–2208.
- [381] Le Wu, Xiangnan He, Xiang Wang, Kun Zhang, and Meng Wang. 2021. A survey on neural recommendation: From collaborative filtering to content and context enriched recommendation. *arXiv preprint arXiv:2104.13030* (2021).
- [382] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James Sharpnack. 2020. SSE-PT: Sequential recommendation via personalized transformer. In *Fourteenth ACM Conference on Recommender Systems*. 328–337.
- [383] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. 2020. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys (CSUR)* (2020).
- [384] Yao Wu, Jian Cao, Guandong Xu, and Yudong Tan. 2021. Tffrom: A two-sided fairness-aware recommendation model for both customers and providers. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1013–1022.
- [385] Yongji Wu, Defu Lian, Neil Zhenqiang Gong, Lu Yin, Mingyang Yin, Jingren Zhou, and Hongxia Yang. 2021. Linear-time self attention with codeword histogram for efficient recommendation. In *Proceedings of the Web Conference 2021*. 1262–1273.
- [386] Yiqing Wu, Ruobing Xie, Yongchun Zhu, Fuzhen Zhuang, Xiang Ao, Xu Zhang, Leyu Lin, and Qing He. 2022. Selective Fairness in Recommendation via Prompts. *arXiv preprint arXiv:2205.04682* (2022).
- [387] Yongkai Wu, Lu Zhang, and Xintao Wu. 2018. On discrimination discovery and removal in ranked data using causal graph. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2536–2544.
- [388] Xin Xia, Hongzhi Yin, Junliang Yu, Qinyong Wang, Guandong Xu, and Quoc Viet Hung Nguyen. 2022. On-Device Next-Item Recommendation with Self-Supervised Knowledge Distillation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 546–555.
- [389] Yikun Xian, Zuohui Fu, Qiaoying Huang, Shan Muthukrishnan, and Yongfeng Zhang. 2020. Neural-symbolic reasoning over knowledge graph for multi-stage explainable recommendation. *arXiv preprint arXiv:2007.13207* (2020).
- [390] Yikun Xian, Zuohui Fu, Shan Muthukrishnan, Gerard De Melo, and Yongfeng Zhang. 2019. Reinforcement knowledge graph reasoning for explainable recommendation. In *Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval*. 285–294.
- [391] Yikun Xian, Zuohui Fu, Handong Zhao, Yingqiang Ge, Xu Chen, Qiaoying Huang, Shijie Geng, Zhou Qin, Gerard De Melo, Shan Muthukrishnan, et al. 2020. CAFE: Coarse-to-fine neural symbolic reasoning for explainable recommendation. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 1645–1654.
- [392] Yikun Xian, Tong Zhao, Jin Li, Jim Chan, Andrey Kan, Jun Ma, Xin Luna Dong, Christos Faloutsos, George Karypis, Shan Muthukrishnan, et al. 2021. Ex3: Explainable attribute-aware item-set recommendations. In *Fifteenth ACM Conference on Recommender Systems*. 484–494.
- [393] Lin Xiao, Zhang Min, Zhang Yongfeng, Gu Zhaoquan, Liu Yiqun, and Ma Shaoping. 2017. Fairness-aware group recommendation with pareto-efficiency. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. 107–115.
- [394] Shitao Xiao, Zheng Liu, Yingxia Shao, Defu Lian, and Xing Xie. 2021. Matching-oriented Embedding Quantization For Ad-hoc Retrieval. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 8119–8129.
- [395] Yilin Xiao, Liang Xiao, Xiaozhen Lu, Hailu Zhang, Shui Yu, and H Vincent Poor. 2020. Deep-reinforcement-learning-based user profile perturbation for privacy-aware recommendation. *IEEE Internet of Things Journal* (2020).
- [396] Lijie Xie, Zhaoming Hu, Xingjuan Cai, Wensheng Zhang, and Jinjun Chen. 2021. Explainable recommendation based on knowledge graph and multi-objective optimization. *Complex & Intelligent Systems* 7, 3 (2021), 1241–1252.
- [397] Minhui Xie, Youyou Lu, Jiazhen Lin, Qing Wang, Jian Gao, Kai Ren, and Jiwu Shu. 2022. Fleche: an efficient GPU embedding cache for personalized recommendations. In *Proceedings of the Seventeenth European Conference on Computer Systems*. 402–416.

- [398] Chang Xu, Jiachen Wang, Liehuang Zhu, Chuan Zhang, and Kashif Sharif. 2019. PPMR: A Privacy-Preserving Online Medical Service Recommendation Scheme in eHealthcare System. *IEEE Internet of Things Journal* 6 (2019), 5665–5673.
- [399] Han Xu, Xiaorui Liu, Yaxin Li, Anil Jain, and Jiliang Tang. 2021. To be robust or to be fair: Towards fairness in adversarial training. In *International Conference on Machine Learning*. PMLR, 11492–11501.
- [400] Zhaozhuo Xu, Weijie Zhao, Shulong Tan, Zhixin Zhou, and Ping Li. 2022. Proximity Graph Maintenance for Fast Online Nearest Neighbor Search. *arXiv preprint arXiv:2206.10839* (2022).
- [401] Bencheng Yan, Pengjie Wang, Kai Zhang, Wei Lin, Kuang-Chih Lee, Jian Xu, and Bo Zheng. 2021. Learning Effective and Efficient Embedding via an Adaptively-Masked Twins-based Layer. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3568–3572.
- [402] Dingqi Yang, Bingqing Qu, and Philippe Cudré-Mauroux. 2018. Privacy-preserving social media data publishing for personalized ranking-based recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2018).
- [403] Jie Amy Yang, Jianyu Huang, Jongsoo Park, Ping Tak Peter Tang, and Andrew Tulloch. 2020. Mixed-Precision Embedding Using a Cache. *arXiv preprint arXiv:2010.11305* (2020).
- [404] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)* (2019).
- [405] Tao Yang and Qingyao Ai. 2021. Maximizing marginal fairness for dynamic learning to rank. In *Proceedings of the Web Conference 2021*. 137–145.
- [406] Jiangchao Yao, Feng Wang, Xichen Ding, Shaohu Chen, Bo Han, Jingren Zhou, and Hongxia Yang. 2022. Device-cloud Collaborative Recommendation via Meta Controller. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4353–4362.
- [407] Jiangchao Yao, Feng Wang, Kunyang Jia, Bo Han, Jingren Zhou, and Hongxia Yang. 2021. Device-cloud collaborative learning for recommendation. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 3865–3874.
- [408] Jiangchao Yao, Shengyu Zhang, Yang Yao, Feng Wang, Jianxin Ma, Jianwei Zhang, Yunfei Chu, Luo Ji, Kunyang Jia, Tao Shen, et al. 2022. Edge-Cloud Polarization and Collaboration: A Comprehensive Survey for AI. *IEEE Transactions on Knowledge and Data Engineering* (2022).
- [409] Sirui Yao and Bert Huang. 2017. Beyond parity: Fairness objectives for collaborative filtering. *Advances in neural information processing systems* 30 (2017).
- [410] Xinyang Yi, Yi-Fan Chen, Sukriti Ramesh, Vinu Rajashekhar, Lichan Hong, Noah Fiedel, Nandini Seshadri, Lukasz Heldt, Xiang Wu, and Ed H Chi. 2018. Factorized deep retrieval and distributed tensorflow serving. In *ser. Conference on Machine Learning and Systems*.
- [411] Chunxing Yin, Bilge Acun, Carole-Jean Wu, and Xing Liu. 2021. Tt-rec: Tensor train compression for deep learning recommendation models. *Proceedings of Machine Learning and Systems* 3 (2021), 448–462.
- [412] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 974–983.
- [413] Bin Yu, Chenyu Zhou, Chen Zhang, Guodong Wang, and Yiming Fan. 2020. A privacy-Preserving multi-Task framework for knowledge graph Enhanced recommendation. *IEEE Access* (2020).
- [414] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 31.
- [415] Fajie Yuan, Xiangnan He, Alexandros Karatzoglou, and Liguang Zhang. 2020. Parameter-efficient transfer from sequential behaviors for user modeling and recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 1469–1478.
- [416] Feng Yuan, Lina Yao, and Boualem Benatallah. 2019. Adversarial collaborative neural network for robust recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1065–1068.
- [417] Fajie Yuan, Guoxiao Zhang, Alexandros Karatzoglou, Joemon Jose, Beibei Kong, and Yudong Li. 2021. One person, one model, one world: Learning continual user representation without forgetting. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 696–705.
- [418] Zhenrui Yue, Zhankui He, Huimin Zeng, and Julian McAuley. 2021. Black-Box Attacks on Sequential Recommenders via Data-Free Model Extraction. In *Fifteenth ACM Conference on Recommender Systems*.
- [419] Meike Zehlke, Francesco Bonchi, Carlos Castillo, Sara Hajian, Mohamed Megahed, and Ricardo Baeza-Yates. 2017. Fa* ir: A fair top-k ranking algorithm. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. 1569–1578.
- [420] Meike Zehlke, Ke Yang, and Julia Stoyanovich. 2021. Fairness in ranking: A survey. *arXiv preprint arXiv:2103.14000* (2021).

- [421] Chengjun Zhang, Jin Liu, Yanzhen Qu, Tianqi Han, Xujun Ge, and An Zeng. 2018. Enhancing the robustness of recommender systems against spammers. *PLoS one* 13, 11 (2018), e0206458.
- [422] Caojin Zhang, Yicun Liu, Yuanpu Xie, Sofia Ira Ktena, Alykhan Tejani, Akshay Gupta, Pranay Kumar Myana, Deepak Dilipkumar, Suvadip Paul, Ikuhiro Ihara, et al. 2020. Model size reduction using frequency based double hashing for recommender systems. In *Fourteenth ACM Conference on Recommender Systems*. 521–526.
- [423] Dell Zhang and Jun Wang. 2021. Recommendation Fairness: From Static to Dynamic. *arXiv preprint arXiv:2109.03150* (2021).
- [424] Fuguo Zhang. 2009. A survey of shilling attacks in collaborative filtering recommender systems. In *2009 International Conference on Computational Intelligence and Software Engineering*. IEEE, 1–4.
- [425] Fuzhi Zhang and Quanqiang Zhou. 2014. HHT-SVM: An online method for detecting profile injection attacks in collaborative recommender systems. *Knowledge-Based Systems* 65 (2014), 96–105.
- [426] Han Zhang, Songlin Wang, Kang Zhang, Zhiling Tang, Yunjiang Jiang, Yun Xiao, Weipeng Yan, and Wen-Yun Yang. 2020. Towards personalized and semantic retrieval: An end-to-end solution for e-commerce search via embedding learning. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2407–2416.
- [427] He Zhang, Bang Wu, Xingliang Yuan, Shirui Pan, Hanghang Tong, and Jian Pei. 2022. Trustworthy Graph Neural Networks: Aspects, Methods and Trends. *arXiv preprint arXiv:2205.07424* (2022).
- [428] Jin Zhang, Qi Liu, Defu Lian, Zheng Liu, Le Wu, and Enhong Chen. 2022. Anisotropic Additive Quantization for Fast Inner Product Search. (2022).
- [429] Jia-Dong Zhang and Chi-Yin Chow. 2018. Enabling probabilistic differential privacy protection for location recommendations. *IEEE Transactions on Services Computing* (2018).
- [430] Mingwu Zhang, Yu Chen, and Jingqiang Lin. 2021. A privacy-preserving optimization of neighborhood-based recommendation for medical-aided diagnosis and treatment. *IEEE Internet of Things Journal* (2021).
- [431] Minxing Zhang, Zhaochun Ren, Zihan Wang, Pengjie Ren, Zhunmin Chen, Pengfei Hu, and Yang Zhang. 2021. Membership inference attacks against recommender systems. (2021).
- [432] Shijie Zhang, Hongzhi Yin, Tong Chen, Zi Huang, Lizhen Cui, and Xiangliang Zhang. 2021. Graph embedding for recommendation against attribute inference attacks. In *Proc. of WWW*.
- [433] Shijie Zhang, Hongzhi Yin, Tong Chen, Quoc Viet Nguyen Hung, Zi Huang, and Lizhen Cui. 2020. Gcn-based user representation learning for unifying robust recommendation and fraudster detection. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 689–698.
- [434] Xinyi Zhang, Chengfang Fang, and Jie Shi. 2021. Thief, Beware of What Get You There: Towards Understanding Model Extraction Attack. *arXiv preprint arXiv:2104.05921* (2021).
- [435] Yongfeng Zhang, Xu Chen, et al. 2020. Explainable recommendation: A survey and new perspectives. *Foundations and Trends® in Information Retrieval* 14, 1 (2020), 1–101.
- [436] Yongfeng Zhang, Guokun Lai, Min Zhang, Yi Zhang, Yiqun Liu, and Shaoping Ma. 2014. Explicit factor models for explainable recommendation based on phrase-level sentiment analysis. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. 83–92.
- [437] Zhikun Zhang, Min Chen, Michael Backes, Yun Shen, and Yang Zhang. 2022. Inference attacks against graph neural networks. In *Proc. of USENIX Security*.
- [438] Zhiwei Zhang, Qifan Wang, Lingyun Ruan, and Luo Si. 2014. Preference preserving hashing for efficient recommendation. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. 183–192.
- [439] Kaiqi Zhao, Gao Cong, Quan Yuan, and Kenny Q Zhu. 2015. SAR: A sentiment-aspect-region model for user preference analysis in geo-tagged reviews. In *2015 IEEE 31st international conference on data engineering*. IEEE, 675–686.
- [440] Kangzhi Zhao, Xiting Wang, Yuren Zhang, Li Zhao, Zheng Liu, Chunxiao Xing, and Xing Xie. 2020. Leveraging demonstrations for reinforcement recommendation reasoning over knowledge graphs. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 239–248.
- [441] Weijie Zhao, Deping Xie, Ronglai Jia, Yulei Qian, Ruiquan Ding, Mingming Sun, and Ping Li. 2020. Distributed hierarchical gpu parameter server for massive scale deep learning ads systems. *Proceedings of Machine Learning and Systems* 2 (2020), 412–428.
- [442] Weijie Zhao, Jingyuan Zhang, Deping Xie, Yulei Qian, Ronglai Jia, and Ping Li. 2019. Aibox: Ctr prediction model training on a single node. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 319–328.
- [443] Wayne Xin Zhao, Yupeng Hou, Xingyu Pan, Chen Yang, Zeyu Zhang, Zihan Lin, Jingsen Zhang, Shuqing Bian, Jiakai Tang, Wenqi Sun, et al. 2022. RecBole 2.0: Towards a More Up-to-Date Recommendation Library. *arXiv preprint arXiv:2206.07351* (2022).

- [444] Wayne Xin Zhao, Shanlei Mu, Yupeng Hou, Zihan Lin, Yushuo Chen, Xingyu Pan, Kaiyuan Li, Yujie Lu, Hui Wang, Changxin Tian, et al. 2021. Recbole: Towards a unified, comprehensive and efficient framework for recommendation algorithms. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 4653–4664.
- [445] Xiangyu Zhao, Haochen Liu, Hui Liu, Jiliang Tang, Weiwei Guo, Jun Shi, Sida Wang, Huiji Gao, and Bo Long. 2021. Autodim: Field-aware embedding dimension search in recommender systems. In *Proceedings of the Web Conference 2021*. 3015–3022.
- [446] Xiangyu Zhao, Long Xia, Liang Zhang, Zhuoye Ding, Dawei Yin, and Jiliang Tang. 2018. Deep Reinforcement Learning for Page-wise Recommendations. In *Proceedings of the 12th ACM Conference on Recommender Systems*. ACM, 95–103.
- [447] Yuchen Zhao, Juan Ye, and Tristan Henderson. 2014. Privacy-Aware Location Privacy Preference Recommendations. In *Proceedings of the 11th International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services*. 120–129.
- [448] Xiangyu Zhao, Haochen Liu, Wenqi Fan, Hui Liu, Jiliang Tang, Chong Wang, Ming Chen, Xudong Zheng, Xiaobing Liu, and Xiwang Yang. 2021. Autoemb: Automated embedding dimensionality search in streaming recommendations. In *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 896–905.
- [449] Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang, Nicholas Jing Yuan, Xing Xie, and Zhenhui Li. 2018. DRN: A deep reinforcement learning framework for news recommendation. In *Proceedings of the 2018 world wide web conference*. 167–176.
- [450] Qiming Zheng, Quan Chen, Kaihao Bai, Huifeng Guo, Yong Gao, Xiuqiang He, and Minyi Guo. 2021. BiPS: Hotness-aware Bi-tier Parameter Synchronization for Recommendation Models. In *2021 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. IEEE, 609–618.
- [451] Ruiqi Zheng, Liang Qu, Bin Cui, Yuhui Shi, and Hongzhi Yin. 2022. AutoML for Deep Recommender Systems: A Survey. *arXiv preprint arXiv:2203.13922* (2022).
- [452] Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Yong Li, and Depeng Jin. 2021. Disentangling user interest and conformity for recommendation with causal embedding. In *Proceedings of the Web Conference 2021*. 2980–2991.
- [453] Zhi Zheng, Chao Wang, Tong Xu, Dazhong Shen, Penggang Qin, Xiangyu Zhao, Baoxing Huai, Xian Wu, and Enhong Chen. 2022. Interaction-aware Drug Package Recommendation via Policy Gradient. *ACM Transactions on Information Systems (TOIS)* (2022).
- [454] Zangwei Zheng, Pengtai Xu, Xuan Zou, Da Tang, Zhen Li, Chenguang Xi, Peng Wu, Leqi Zou, Yijie Zhu, Ming Chen, et al. 2022. CowClip: Reducing CTR Prediction Model Training Time from 12 hours to 10 minutes on 1 GPU. *arXiv preprint arXiv:2204.06240* (2022).
- [455] Kun Zhou, Hui Yu, Wayne Xin Zhao, and Ji-Rong Wen. 2022. Filter-enhanced MLP is all you need for sequential recommendation. In *Proceedings of the ACM Web Conference 2022*. 2388–2399.
- [456] Ke Zhou and Hongyuan Zha. 2012. Learning binary codes for collaborative filtering. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 498–506.
- [457] Jieming Zhu, Jinyang Liu, Weiqi Li, Jincai Lai, Xiuqiang He, Liang Chen, and Zibin Zheng. 2020. Ensembled CTR prediction via knowledge distillation. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2941–2958.
- [458] Jieming Zhu, Jinyang Liu, Shuai Yang, Qi Zhang, and Xiuqiang He. 2020. Fuxictr: An open benchmark for click-through rate prediction. *arXiv preprint arXiv:2009.05794* (2020).
- [459] Xue Zhu and Yuqing Sun. 2016. Differential privacy for collaborative filtering recommender algorithm. In *Proceedings of the 2016 ACM on International Workshop on Security and Privacy Analytics*.
- [460] Yaxin Zhu, Yikun Xian, Zuohui Fu, Gerard de Melo, and Yongfeng Zhang. 2021. Faithfully explainable recommendation via neural logic reasoning. *arXiv preprint arXiv:2104.07869* (2021).
- [461] Ziwei Zhu, Xia Hu, and James Caverlee. 2018. Fairness-aware tensor-based recommendation. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 1153–1162.
- [462] Ziwei Zhu, Jingu Kim, Trung Nguyen, Aish Fenton, and James Caverlee. 2021. Fairness among new items in cold start recommender systems. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 767–776.
- [463] Ziwei Zhu, Jianling Wang, and James Caverlee. 2020. Measuring and mitigating item under-recommendation bias in personalized ranking systems. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 449–458.