

Robustness in Fairness against Edge-level Perturbations in GNN-based Recommendation

Ludovico Boratto¹[0000-0002-6053-3015], Francesco Fabbri²[0000-0002-9631-1799],
Gianni Fenu¹[0000-0003-4668-2476], Mirko Marras^{*,1}[0000-0003-1989-6057], and
Giacomo Medda¹[0000-0002-1300-1876]

¹ University of Cagliari, Cagliari, Italy

{ludovico.boratto,mirko.marras}@acm.org, {fenu,giacomo.medda}@unica.it,

² Spotify, Barcelona, Spain

francescof@spotify.com

Abstract. Efforts in the recommendation community are shifting from the sole emphasis on utility to considering beyond-utility factors, such as fairness and robustness. Robustness of recommendation models is typically linked to their ability to maintain the original utility when subjected to attacks. Limited research has explored the robustness of a recommendation model in terms of fairness, e.g., the parity in performance across groups, under attack scenarios. In this paper, we aim to assess the robustness of graph-based recommender systems concerning fairness, when exposed to attacks based on edge-level perturbations. To this end, we considered four different fairness operationalizations, including both consumer and provider perspectives. Experiments on three datasets shed light on the impact of perturbations on the targeted fairness notion, uncovering key shortcomings in existing evaluation protocols for robustness. As an example, we observed perturbations affect consumer fairness on a higher extent than provider fairness, with alarming unfairness for the former. Source code: <https://github.com/jackmedda/CPFairRobust>.

Keywords: Robustness · Fairness · Recommendation · GNN · Perturbation · Multi-Stakeholder · Provider · Consumer.

1 Introduction

Individuals are increasingly interacting with recommender systems, enjoying the benefits of personalized services provided by e-commerce and streaming platforms. These services are designed to adapt to the preferences and interests of consumers about the content they discover, while also meeting the expectations of content providers, who seek visibility and engagement. However, the experiences of these stakeholders can be compromised by specialized attacks targeting recommender systems. These attacks aim to manipulate the recommendations generated by the systems according to the attacker’s objectives [63,42,1].

* Corresponding author.

The effectiveness of attacks against recommender systems has been demonstrated across a diverse range of recommendation models, including those based on k -nearest neighborhood [42], matrix factorization [31], association rules [57], recurrent neural networks [49], and graphs [19,39]. In this context, ‘poisoning’ attacks have become particularly prevalent [63,51]. These attacks primarily involve data perturbation during the training stage [1,51], often through the introduction of fake users, also known as ‘shilling’ attacks [1,42,49]. Concerted efforts have been increasingly devoted to enhancing the robustness of recommender systems against various attacks [63,12,29,65,64]. Fraudster detection [65,10] and adversarial learning [12,29,55] have emerged as the primary defensive strategies. The former seeks to identify and mitigate the influence of fake users, while the latter introduces perturbations to strengthen models against adversarial samples.

Differently from other research fields (e.g., computer sensing [16,15,59,3], code generation [36], and program understanding [62]), attacks against recommendation have prioritized the maximization of model disruption, often at the expense of constructive objectives, without considering their impact on model robustness. The existing methods for evaluating robustness merely compare recommendation utility before/after attacks [63]. This practice is unfortunately limited, given that the overall utility can remain stable even when recommendations are significantly altered [41]. This limitation prevents from detecting the impact of attacks on beyond-accuracy objectives, such as trustworthiness [52], fairness [54], and explainability [66]. Few prior works consider robustness in beyond-accuracy properties, such as bias [46] and sparseness [68], but do not cover fairness [63]. In other domains, this interplay has already been addressed [43,35].

In this paper, we provide a novel comprehensive analysis on the robustness of graph-based recommender systems in terms of fairness, referred to as *robustness in fairness*. Specifically, we investigate the extent to which the system fairness remains stable, from both the consumer [7,13,54] and the provider [47,25,26] sides, under attack scenarios. We address this issue on systems based on graphs due to their state-of-the-art performance and the extensive range of attacks on graph data [51,19,64,39,65,14]. Adding and deleting edges is a popular technique for attacks in graph data [51]. To this end, we extended an approach that perturbs a graph at the edge-level to explain the predictions in several downstream tasks [34,30]. This approach iteratively performs poisoning-like attacks against recommender systems based on Graph Neural Networks (GNNs) and monitors fairness as the user-item interaction graph gets gradually perturbed, encompassing different types of perturbations and fairness operationalizations. Although our experimental evaluation is driven by the employed attack, it is important to note that the attack itself does not constitute the main contribution. Rather, this paper specifically aims to study the robustness in fairness in recommendation and explore the nuances in the GNN models’ outcomes after attacks.

Our study operates within a white-box scenario, simulating the role of an attacker aiming to compromise the group fairness of a recommender system. In this scenario, the perturbation process involves modifying the input graph that feeds into the GNN. Such attacks may have real-world consequences, including

compromising a company’s reputation in the public eye [38], both through media coverage and legal implications that may result in sanctions and other repercussions. Concerning the recent regulations in terms of robustness and fairness of automated systems [40,21], such consequences could illustrate a worrying scenario. Our experimental study showcases an extensive characterization of robustness in fairness against poisoning-like attacks, by employing three datasets and three GNN-based recommender systems. The tested models exhibit a higher sensitivity to attacks tailored for group consumer fairness compared with provider one. Specifically, the unfairness levels across consumer groups can be increased by a restrained amount of perturbations, whereas the impact on provider fairness is limited by the prior unfairness level exhibited in the original recommendations.

2 Related Work

2.1 Attacks and Robustness in Recommendation

The researchers addressing attacks and robustness in recommendation do not necessarily see these properties as interconnected, although most of the literature in robustness regards attacks [63]. In fact, several papers solely focused on identifying attacks and treating them as strategies for achieving maximum disruptions [42,19,14,49,57,31] by injecting fake users to increase the recommendation of specific items [42,19], or adversarially generate unnoticeable fake profiles [14]. Conversely, other works focused on improving robustness without necessarily viewing the adopted attack as the actual contribution [41,65,60,12,61]. For instance, [65] detected suspicious users as fraudsters using neural random forests. Despite these advancements, a comprehensive analysis of the impact of attacks on the models’ robustness, particularly in terms of accuracy and other critical properties, is notably absent. This gap in research is remarkable, especially when compared to analogous studies conducted in other fields [16,15,36,59,3,62,43].

2.2 Fairness in Recommendation

Due to recently issued regulations [40,21], researchers are increasingly prioritizing beyond-accuracy aspects in recommendation, as explainability [66] and fairness [54]. The relevant amount of recent works studying consumer and provider (un)fairness addressed their assessment [7,8,25,26,47], mitigation [32,9,5,22,6], and explanation [23,24,17,37]. Despite calls for unifying the goals of robustness and fairness in recommendation [63], to the best of our knowledge, [58] is the only work that focused on both properties. Specifically, [58] proposed a fair and distributionally robust method to solve the distribution shift problem between the training and testing sets. Unfortunately, no study addressed the assessment of robustness in group consumer/provider fairness against specialized attacks.

2.3 Robustness and Beyond-accuracy Aspects

Some studies on robustness in recommendation considered beyond-accuracy properties, e.g., bias [46] and sparseness [68], as pertaining to a kind of robustness [63].

However, their scope does not cover the fairness property envisioned in our study. On the other hand, the literature in other fields has witnessed the introduction of novel techniques of certified robustness for text classification [43], and novel attacks that target the fairness of classifiers [38,48]. Nevertheless, their works regard classification tasks, where attacks and robustness methods differ from the ones employed in recommendation, as those targeted by this paper.

3 Methodology

3.1 Perturbation Task in Graph-based Recommendation

Our perturbation task tailored for GNN-based recommender systems aims to perturb the adjacency matrix through edge perturbations to alter the predicted recommendation lists, and test the systems’ robustness in fairness. We then distinguish between the recommendation task and the proper perturbation task.

Recommendation task. In a typical recommendation scenario, a model learns the preferences of a set of users U from their past interactions with a catalog of items I . The network of user-item interactions can be represented by means of an undirected bipartite graph $G = (V, E)$, where $V = U \cup I$ is the set of nodes, E is the set of edges between user and item nodes. G can be encoded in a $n \times n$ adjacency matrix A , where $A_{u,i} > 0$ denotes an edge links the user u with the item i , otherwise $A_{u,i} = 0$. We can feed A to a GNN f to predict the probability of missing user-item links. Specifically, f can be parameterized by a weight matrix W and optimized to recommend to each user a list of the top- k items sorted by the predicted linking probability in descending order. Let $q_u@k$ be the top- k list recommended to user u and $Q@k$ the set of all $q_u@k, \forall u \in U$.

Perturbation task. Following [63], robustness can be estimated by the disparity between the performance measured with the original (non-perturbed) data and the perturbed data. A model reporting a disparity lower than a threshold ϵ and a perturbation bounded by a constant γ would be denoted as (γ, ϵ) -robust. In our graph-based scenario, the original data is the adjacency matrix A , and we denote its perturbed version as \tilde{A} . Given our fairness-related task, we define the performance by means of a fairness metric M . We can then formally define the (γ, ϵ) -robustness in fairness of our recommender system f as follows:

$$\Delta = M(f(\tilde{A}, W), A) - M(f(A, W), A), \quad \|\Delta\|_2^2 \leq \epsilon, \quad |\tilde{E}| \leq \gamma \quad (1)$$

where M estimates the fairness level based on the outcome of f and A , and \tilde{E} denotes the set of candidate edges for perturbation. Although [63] does not guide the selection of ϵ and γ , a small ϵ guarantees a greater level of robustness, while a small γ reflects an attack that is harder to detect. Focused on analyzing the robustness in fairness, we do not set a fixed bound on the extent of edge perturbations, i.e. $\gamma = +\infty$. Addressing edge perturbations as deletion and addition, the range of $|\tilde{E}|$ is then $[1, |E|]$ for deletion, and $[1, |U| \times |I| - |E|]$ for addition.

We seek to test the robustness in fairness in recommendation by identifying the edge perturbations that maximize the disparity in (1), prioritizing fewer perturbations. In other words, we aim to optimize the following objective function:

$$\min_{\hat{p}} - \left\| M(f(\tilde{A}, W; \hat{p}), A) - M(f(A, W), A) \right\|_2^2 + \lambda \left\| \Gamma(\tilde{A}, A) \right\|_2^2 \quad (2)$$

where Γ is a distance function [51,34], \hat{p} is a trainable weight used to identify the edges to be perturbed, $\lambda \in \mathbb{R}$ is a hyper-parameter that is used to control the weight between the two terms. The minus sign applied to the first term optimizes the objective function towards maximizing the disparity Δ , resulting in unfair recommendations. However, the original fairness level estimated by M is not affected by the perturbation process. Hence, we can simplify (2) as follows:

$$\min_{\hat{p}} -M(f(\tilde{A}, W; \hat{p}), A) + \lambda \left\| \Gamma(\tilde{A}, A) \right\|_2^2 \quad (3)$$

3.2 Graph Perturbation Mechanism

Following works of explainability in GNNs [34,30], we use a sparsification method to obtain a binary perturbation tensor from a trainable real-valued weight (first described by [50]) and extend it to the recommendation scenario. First, we enlarge the space of the candidate edges to the entire graph. Second, we replace the perturbation matrix P used in [34] with a binary perturbation vector $p \in \{0, 1\}^{|\tilde{E}|}$, such that solely the relevant user-item connections are perturbed instead of affecting also self-loops, user-user and item-item links. p can be derived from the trainable weight \hat{p} (2) by applying a sigmoid transformation and a binarization to \hat{p} , so as to map values lower than 0.5 to 0, otherwise to 1.

The perturbation mechanism can be thought as a substitution process that updates the entries of the adjacency matrix A with the entries of p , resulting in \tilde{A} . A fixed relation between the 2D index of A and the 1D index of p establishes which candidate edge $(u, i) \in \tilde{E}$ will be perturbed by the j -th entry of p . The entries update process resulting in \tilde{A} is formally defined as:

$$\tilde{A} = A \dot{+} p, \quad \tilde{A}_{u,i} = \begin{cases} p_j & \text{if } (u, i) \in \tilde{E} \\ A_{u,i} & \text{otherwise} \end{cases} \quad (4)$$

$\dot{+}$ denotes the perturbation operator for edge deletion $\dot{-}$ *Del* or addition $\dot{+}$ *Add*.

This perturbation mechanism is performed iteratively by gradually modifying A to generate \tilde{A} , until the perturbed edges optimize the targeted task. Specifically, we initialize \hat{p} based on the perturbation type (deletion or addition), such that no edge is affected in A , i.e. $A \dot{+} p = A$. At each iteration, we generate \tilde{A} with \hat{p} , and feed \tilde{A} to the GNN f to produce the corresponding recommendations. The latter are processed by (3) to estimate the fairness level, the distance between \tilde{A} and A , and to update the weight \hat{p} accordingly. Finally, the process stops based on a predefined criterion, e.g., the impact of the last perturbation.

3.3 Fairness Notion and Operationalization

Fairness Notion. We proceed to define the fairness metric M . We follow recent works [7,56,32,23,47,22,20,2] that emphasized the relevance of the group fairness notion of *demographic parity* from both the consumer [7,56,32] and provider side [23,47,22]. For the former side, demographic parity is satisfied if consumer groups experience the same level of recommendation utility. For the latter side, the notion is satisfied if the probability of being recommended is equal across provider groups, proportionally to their representation in the catalog.

We ground our work on a binary setting as previous studies [7,32,23,22], where each stakeholder set Z ($Z \subseteq U$ for consumers, $Z \subseteq I$ for providers³) can be partitioned in two groups, Z_1 and Z_2 . Multiple attributes, e.g., gender and age, of each stakeholder could produce a distinct partition of Z in two groups. Demographic parity (DP) can then be operationalized as the following disparity:

$$DP = \|S(f(A, W), A^{Z_1}) - S(f(A, W), A^{Z_2})\|_2^2 \quad (5)$$

where S represents the metric used to estimate the performance w.r.t. the corresponding stakeholder, e.g., exposure for a provider, A^{Z_1} and A^{Z_2} respectively denote the adjacency sub-matrices with regard to the two partitions Z_1 and Z_2 .

Consumer and Provider Fairness Operationalization. Depending on the adopted metric S , we can define specific operationalizations of DP, which reflect distinct perspectives of the unfairness issue. For each stakeholder, we contemplate two types of operationalization, a rank-aware and a rank-agnostic one.

We underline that DP represents M in (1)-(2)-(3), given that it estimates the fairness performance of a recommender system. Therefore, we define each operationalization of DP, the corresponding metric S used for evaluation, and the differentiable approximation of S to be used in the objective function in (3):

- **Consumer Preference (CP):** it estimates the consumer fairness as the disparity across consumers groups in rank-aware top- k recommendation utility, which can be measured by the Normalized Discounted Cumulative Gain (NDCG@ k). Following [56,44], NDCG (N@ k) is approximated as the differentiable function \widehat{NDCG} , where the rank of an item is defined in terms of the pairwise preference with respect to any other item in the catalog.
- **Consumer Satisfaction (CS):** it estimates the consumer fairness as the disparity across consumers groups in rank-agnostic top- k recommendation utility, which can be measured by the Precision (P@ k). We optimize P@ k by treating the recommendation task as a binary classification task, by using a sigmoid function followed by a binary cross entropy loss. This task aims to include relevant items in the top- k list, regardless of their position.
- **Provider Exposure (PE):** it estimates the provider fairness as the disparity in exposure across providers groups. Following [25,26,23,22], we define the

³ We do not consider other features and associate each item with a distinct provider.

exposure of a generic provider group I_* as the average number of exposures in $Q@k$ (estimated by an indicator function $\mathbb{1}[\cdot]$ as [23,22]), discounted by the importance of their position [47,25,26] (as for DCG), and normalized by the ideal exposure [25,26] (as for NDCG). Formally:

$$\text{Exposure}(I_* | Q@k) = \frac{|I|}{|I_*|} \frac{1}{|U|} \sum_{u \in U} \frac{\sum_{j=1}^k \frac{\mathbb{1}[i_j \in I_*]}{\log_2(j+1)}}{\sum_{j=1}^k \frac{\mathbb{1}[i_j \in I]}{\log_2(j+1)}} \quad (6)$$

We leverage the approximation proposed by [23], where the indicator function is replaced by the predicted linking probability (item relevance).

- **Provider Visibility (PV)**: it estimates the provider fairness as the disparity in visibility across providers groups. Following [18,25,26], we define the visibility of a generic provider group I_* as the average number of exposures in $Q@k$ (estimated by $\mathbb{1}[\cdot]$ as for PE). Formally:

$$\text{Visibility}(I_* | Q@k) = \frac{|I|}{|I_*|} \frac{1}{|U|k} \sum_{u \in U} \sum_{i \in q_u @k} \mathbb{1}[i \in I_*] \quad (7)$$

We approximate PV by replacing $\text{Visibility}(\cdot, \cdot)$ with the loss function used for CS , which would aim to include the items of I_* in the top- k list.

CP and CS $\in [0, 1]$, PE and PV $\in [0, \frac{|I|}{|I_1|}]$, for which 0 denotes fairness.

The approximations of PE and PV are inspired by [23], but the authors measure the disparity only according to the top- k items, limiting the exposure/visibility information, given that k is usually small. It causes the gradient to be computed only for the top- k items, which will not necessarily be included in the final recommendations, due, for instance, to items already enjoyed by some users. Instead, we set k to be 10% of the items catalog size $|I|$ to expand the operational scope of PE and PV. This choice also helps when just one of the groups is represented in $Q@k$, and enables our perturbation task to better optimize the presence of one of the groups in positions closer to the top- k ones.

4 Experimental Evaluation

The following experiments aim to answer the following research questions:

- RQ1:** What is the extent to which edge perturbations impact the robustness in fairness of recommender systems?
- RQ2:** Are the adopted models similarly affected in terms of alterations in robustness in fairness as edge perturbations gradually increase?
- RQ3:** Which consumer or provider group should be more affected by the edge perturbation to nuke the models' robustness in fairness?

Table 1. **Left:** datasets’ statistics (*Repr*: Representation, *O*: Older, *Y*: Younger, *F*: Female, *M*: Males). **Right:** original models’ performance in N@10 (%) and P@10 (%).

	ML1M [27]	LF1K [11]	INS [32]						
# Users	6,040	268	346						
# Items	3,706	51,609	20						
# Interactions	1,000,209	200,586	1,879						
Domain	Movie	Music	Insurance						
Age	O : 43.4%	O : 42.2%	O : 49.4%						
	Y : 56.6%	Y : 57.8%	Y : 50.6%						
Repr.	F : 28.3%	F : 42.2%	F : 23.4%						
	M : 71.7%	M : 57.8%	M : 76.6%						

	INS		LF1K		ML1M	
	N@10	P@10	N@10	P@10	N@10	P@10
GCMC	76.40	9.60	39.66	38.47	12.62	11.35
LGCN	78.07	9.74	39.81	38.36	12.68	11.35
NGCF	78.39	9.74	39.71	38.43	12.94	11.65

4.1 Evaluation Setting

Evaluation Protocol The perturbation process is run for 200 epochs, but early stopped if the increment in Δ after 15 consecutive epochs is lower than 0.001.

Given our objective of testing the robustness in fairness, we do not perform a classic poisoning attack as defined in [63,1,51]. Specifically, the combination of datasets, models, perturbation types, and fairness operationalizations sums up to 108 attacks, but we aim to also address intermediary perturbation stages, which would result in an impracticable amount of models re-training processes. To this end, we estimate the first term of Δ by substituting the original adjacency matrix A with the perturbed one \tilde{A} at the inference stage, and maintaining the models’ parameters constant. If \tilde{A} was generated by external tools or approaches, such poisoning-like attack would reflect a white/grey-box setting, with the unique requirement of having access to the saved representation of A .

Models. We rely on Recbole [67] and select GCMC [4], LightGCN (LGCN) [28], and NGCF [53] as the GNN-based recommender systems for our study. Though the set of employed models is limited, they cover different architectures to learn the users’ preferences: GCMC leverages an auto-encoder structure, LGCN learns from linear relationships between users and items, NGCF adopts features transformation and nonlinear activation on the message-passing step. While the last two models generate the recommendations by the learnt user and item embeddings, GCMC performs a complete forward process during inference. Thus, GCMC represents a more suitable candidate for the attacker, who does not need to force the embeddings re-generation after the graph perturbation.

Datasets. We rely on [7], which includes MovieLens-1M [27] (ML1M) and LFM-1K [11] (LF1K)⁴. We also consider Insurance [33] (INS) and discard consumers with less than 5 interactions. Table 1 reports the datasets’ statistics and the original models’ performance. We use the items’ popularity to partition the provider set in short-head I_1 and long-tail I_2 items such that $\frac{|I_1|}{|I_2|} = \frac{1}{4}$ as in [23,22]. For

⁴ The timestamp of each (u, i) refers to the last interaction between u and a i ’s song.

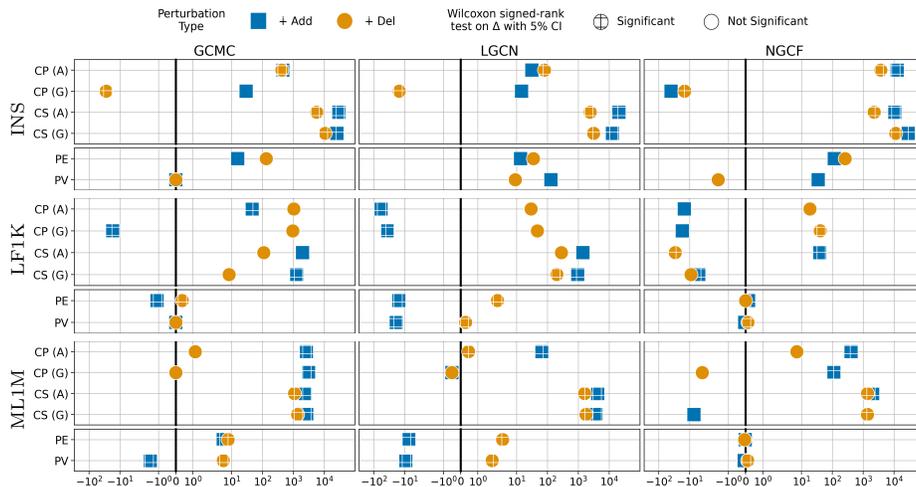


Fig. 1. Impact of edges addition ($\dot{+}$ *Add*) and deletion ($\dot{+}$ *Del*) on the robustness in fairness, reported as the relative difference in M between the non-perturbed and perturbed model, i.e. $\Delta/M(f(A, W), A)$ ((1)). A stands for *Age*, G for *Gender*.

each user, the interactions are sorted in ascending order of recency, split by a ratio 7:1:2, and each split respectively assigned to the training, validation, and testing set. We use the validation set to select the best original model, and the ground truth of the testing set to optimize the fairness operationalizations (Section 3.3). Using the test relevance judgements would result in a powerful attack and support our main concern, i.e. the analysis of robustness in fairness⁵.

4.2 RQ1: Impact on Robustness in Fairness

We first assess the impact of edge-level perturbations on the robustness in fairness, i.e. Δ , where M corresponds to the formulation of DP in (5) and its specialized operationalizations, e.g., PE for provider exposure fairness.

Figure 1 reports the impact on robustness in fairness as the relative difference in M between the original fairness level and after each perturbation type, i.e. Δ divided by the second term in (1). The reported values pertain to the iteration where at least one edge was perturbed and \hat{A} affected the most Δ for each model.

The formulation of Δ in (1) enables us to easily denote positive values as increment in fairness level disparity across groups, i.e. increment in DP, otherwise a decrement of the latter. Therefore, positive values reflect a successful outcome for the attacker, whose goal is to make the system generate unfair recommendations. Negative values derive from instances of \hat{A} that reduce the unfairness level even at the first iterations. The x-axis has been symmetrically log-normalized to highlight the remarkable impact on Δ on the systems. For instance, both $\dot{+}$ *Add* and $\dot{+}$ *Del* significantly caused an impact for CS on LGCN under ML1M,

⁵ This design choice does not constitute the issue highlighted in [45].

precisely increasing DP by more than 1,000%. Also the fairness levels under INS were remarkably affected, reaching values higher than 10,000%.

Successful attacks are especially reported in terms of edge deletions, given that most of the points labeled as $\dot{+}$ *Del* are depicted at the right of the zero line. Conversely, $\dot{+}$ *Add* caused varied results, ranging from reductions of fairness level disparity, e.g., on LGCN under LF1K, to disruptions in robustness in fairness, e.g., on GCMC under ML1M for the consumer side. This observation could be related to the addition of information in the graph due to $\dot{+}$ *Add*, but also to the larger edge sample space for the latter compared with $\dot{+}$ *Del*. The experiments on provider fairness confirm $\dot{+}$ *Del* as a more effective perturbation attack than $\dot{+}$ *Add*. However, several settings report an optimal robustness, e.g., GCMC and NGCF under LF1K, and other ones a negative orientation of Δ , e.g., LGCN under LF1K and ML1M. Across the models, NGCF exhibits the least sensitivity to perturbations, especially on provider fairness and under LF1K, given that most of the points are close to the zero line. This may be attributed to the feature transformation and nonlinear activation applied in the NGCF message-passing scheme, which diminish the impact of the perturbed graph on the predictions.

4.3 RQ2: Robustness in Fairness under Incremental Perturbations

The previous research question aimed to highlight if the edge perturbations could be able to affect the robustness in fairness of the considered models, accounting only for the iteration with the highest impact on Δ . It is then unclear to what extent the robustness is affected by the gradual increment of edge perturbations.

To this end, we report the DP for each iteration of the perturbation process through the points in Figure 2. For each setting, a horizontal dashed line denotes the original DP to highlight the gap caused by the increment of perturbed edges. A robust model would then be described by points close to the dashed line, while far points conceive the perturbed edges significantly affected the fairness level. Even the size of such points is important: big points close to the dashed line denote the model is robust despite a relevant amount of perturbed edges.

Consumer-side. Across the datasets, NGCF is the most robust model, given that in most of the setting the points are distributed closer to the original DP compared with the other two systems. This likely stems from the observation highlighted in the previous section. GCMC and LGCN behave in a similar way, especially in terms of CS, for which the addition of edges ($\dot{+}$ *Add*) has a more significant influence on DP compared with the deletion ($\dot{+}$ *Del*). However, LGCN reports a systematic robustness for the rank-aware operationalization CP with only a few small points reported in the proximity of the original DP.

Some experiments underline the attack could be influenced by the dataset and sensitive attribute. The former influence factor can be observed on GCMC, for which $\dot{+}$ *Del* significantly affected CP under LF1K; a similar result was reported under ML1M, but due to the other perturbation type $\dot{+}$ *Add*. The latter influence factor can be observed under INS on NGCF, where the impact of $\dot{+}$ *Add* on DP is reported only across age groups, regardless of the operationalization.

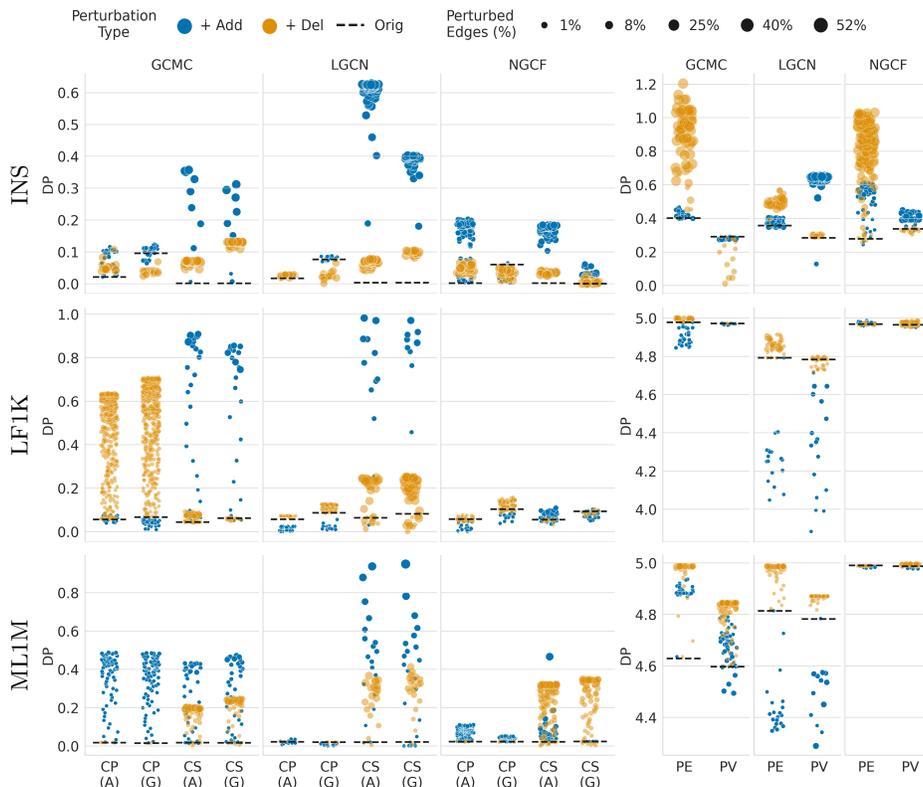


Fig. 2. Trend of the operationalizations of DP across different stages of the perturbation process. Each stage reflects a fraction of perturbed edges w.r.t. \tilde{E} , depicted as points gradually larger as the amount of perturbed edges increases. The horizontal dashed line labeled as *Orig* denotes the DP of the recommendations generated by the non-perturbed system. In the consumer-side results, *A* stands for *Age*, *G* for *Gender*.

Provider-side. Differently from the consumer-side, the original systems report a high degree of unfairness, close to the maximum value $\frac{|I|}{|I_1|} = 5$. Hence, the edge perturbations adopted by a potential attacker cannot significantly affect the robustness in fairness, as highlighted under LF1K on all models, and under ML1M on NGCF, where the original level of PE and PV is remarkable.

In the other settings, deleting edges is more reliable in influencing the robustness in provider fairness compared with \dagger *Add*. This behavior is especially reported under INS, but also under ML1M on GCMC and LGCN.

Other observations contrast with the consumer-side evaluation. First, NGCF exhibits a high degree of sensitivity to the perturbations under INS, as emphasized by the gradual increment in PE as more edges are deleted. This is possibly due to the limited size of INS, which prevents NGCF from robustly learning the users' preferences in a generalized way. Second, the differences in trend between GCMC and LGCN are more remarked, with \dagger *Add* affecting more the robust-

ness of GCMC under ML1M w.r.t. the LGCN’s one. This can be explained by the GCMC encoder architecture, which possibly interprets the added edges as noise, which is instead captured as new information by the linear step in LGCN.

4.4 RQ3: Edge Perturbations Influence on Groups

Group unfairness denotes one of the groups is favored by the system, e.g., with a higher level of exposure. We denote such group as *advantaged* and the other group as *disadvantaged*. We seek to discover whether the impact on Δ is due to a greater amount of edge perturbations applied to the advantaged or the disadvantaged group. To this end, we define *Edge Impact* (EI) as the ratio between the distribution of edge perturbations on the nodes of a group and the representation of the latter. Let $\neg\tilde{E} \subseteq \tilde{E}$ be the subset of edges actually perturbed, formally $EI_* = \frac{|\neg\tilde{E}_*|/|\neg\tilde{E}|}{|Z_*|/|Z|}$, where Z_* is a group of consumers or providers, and $\neg\tilde{E}_*$ is the subset of edges perturbed w.r.t. the nodes in Z_* . We also define the difference $\Delta EI = EI_{Adv} - EI_{Disadv}$ between the advantaged and the disadvantaged group.

In Figure 3, we depict ΔEI on the y-axis and Δ on the x-axis to study their relationship. The values of Δ pertain to our first experiment in Section 4.2.

In the consumer side, most of the settings reporting an impact on Δ were caused by edge perturbations prioritizing the disadvantaged group, i.e. $\Delta EI < 0$. In particular, this behavior regards $\dagger Del$, which removed interactions of the disadvantaged group to, intuitively, reduce its recommendation utility. Conversely,

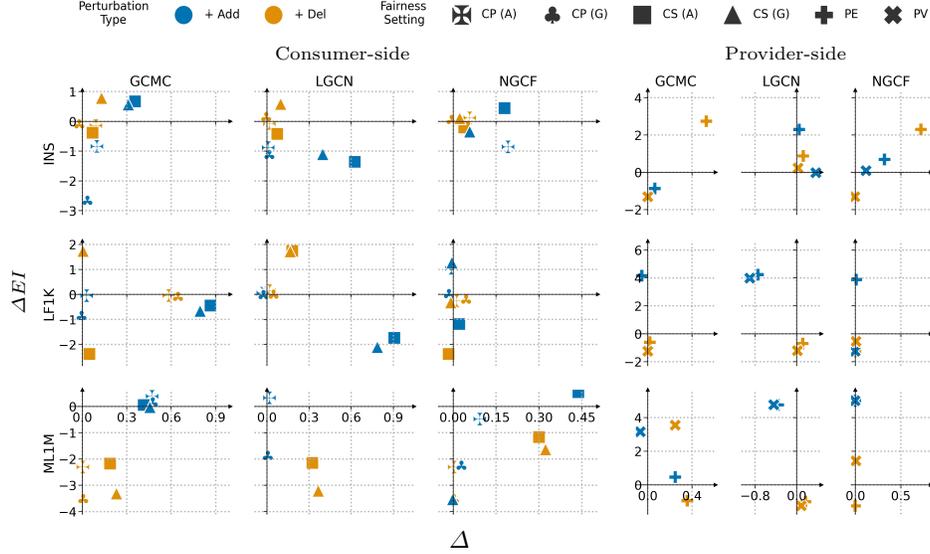


Fig. 3. Relationship between ΔEI (y-axis), i.e. the disparity in edge perturbations distribution between the advantaged and the disadvantaged group, and Δ (x-axis), i.e. the disparity in fairness level before and after the attack. Successful attacks ($\Delta > 0$) are caused by targeting more the advantaged group if $\Delta EI > 0$, otherwise the disadvantaged group was targeted. *A* stands for *Age*, *G* for *Gender*. Some points overlap.

we expect \dagger *Add* to report $\Delta EI > 0$. However, this is actually true only under ML1M on LGCN and NGCF, and slightly under INS on GCMC and NGCF. Other settings, e.g., under LF1K, illustrate Δ was affected by adding more edges to the disadvantaged group, highlighting that a higher number of interactions does not systematically correspond to a higher recommendation utility.

In the provider side, the disparity in edge perturbations distribution is more remarked in comparison with the consumer-side results. However, such disparity does not reflect a relevant impact on Δ in most settings, as already highlighted in other experiments. Successful attacks ($\Delta > 0$) are mostly reported under INS for PE, where ΔEI is systematically above the origin. Hence, the impact on the robustness in provider fairness was caused by deleting or adding more edges to the advantaged group. Other scenarios where Δ was affected, e.g., on GCMC under ML-1M, confirm such observation or report close points to the x-axis.

5 Conclusions

In this paper, we present the concept of robustness in fairness and raised attention towards the issues caused by related attacks in recommendation. Compared to prior work, our analysis aimed to assess the robustness in fairness of GNN-based recommender systems against poisoning-like attacks based on edge-level perturbations, focusing on the models’ robustness and not on the attack itself.

Even though the range of considered models is limited, they represent consistent baselines in the literature, and cover varied GNN architectures. It follows that the adopted models represent perfect candidates for an analysis of robustness in fairness, a topic that is still unexplored in recommendation compared with other fields. Nevertheless, future works will cover a wider set of models.

Moreover, although the adopted attack is focused to test the robustness in fairness, the same approach could be extended to targeted attacks, e.g., against specific consumers or providers. In other words, based on how each stakeholder partitions are established, our attack could be used to modify the model performance to specifically favor one of the tailored groups. We plan to investigate other types of attack, such as those relying on perturbations based on re-wiring.

Additionally, our white-box setting could be simplified to a grey-box one by removing the assumption of having access to the model parameters. Given that the attack solely requires the substitution of the adjacency matrix with its perturbed version, the latter could be generated through a different process. Future work will be focused on exploring grey- or black-box settings.

Acknowledgement. We acknowledge financial support under the National Recovery and Resilience Plan (NRRP), Miss. 4 Comp. 2 Inv. 1.5 - Call for tender No.3277 published on Dec 30, 2021 by the Italian Ministry of University and Research (MUR) funded by the European Union – NextGenerationEU. Prj. Code ECS0000038 eINS Ecosystem of Innovation for Next Generation Sardinia, CUP F53C22000430001, Grant Assignment Decree N. 1056, Jun 23, 2022 by the MUR.

References

1. Anelli, V.W., Deldjoo, Y., Noia, T.D., Merra, F.A.: Adversarial recommender systems: Attack, defense, and advances. In: Ricci, F., Rokach, L., Shapira, B. (eds.) *Recommender Systems Handbook*, pp. 335–379. Springer US (2022)
2. Atzori, A., Fenu, G., Marras, M.: Explaining bias in deep face recognition via image characteristics. In: *Proc. of the IEEE International Joint Conference on Biometrics, IJCB*. pp. 1–10. IEEE (2022)
3. Barbu, A., Mayo, D., Alverio, J., Luo, W., Wang, C., Gutfreund, D., Tenenbaum, J., Katz, B.: Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In: *Proc. of the Annual Conference on Neural Information Processing Systems, NeurIPS*. pp. 9448–9458 (2019)
4. den Berg, R.V., Kipf, T.N., Welling, M.: Graph convolutional matrix completion. *CoRR* **abs/1706.02263** (2017)
5. Boratto, L., Fabbri, F., Fenu, G., Marras, M., Medda, G.: Counterfactual graph augmentation for consumer unfairness mitigation in recommender systems. In: *Proc. of the 32nd ACM International Conference on Information and Knowledge Management, CIKM*. pp. 3753–3757. ACM (2023)
6. Boratto, L., Fenu, G., Marras, M.: Interplay between upsampling and regularization for provider fairness in recommender systems. *User Model. User Adapt. Interact.* **31**(3), 421–455 (2021)
7. Boratto, L., Fenu, G., Marras, M., Medda, G.: Consumer fairness in recommender systems: Contextualizing definitions and mitigations. In: *Proc. of the 44th European Conference on IR Research, ECIR*. LNCS, vol. 13185, pp. 552–566. Springer (2022)
8. Boratto, L., Fenu, G., Marras, M., Medda, G.: Practical perspectives of consumer fairness in recommendation. *Inf. Process. Manag.* **60**(2), 103208 (2023)
9. Burke, R., Sonboli, N., Ordonez-Gauger, A.: Balanced neighborhoods for multi-sided fairness in recommendation. In: *Proc. of the Conference on Fairness, Accountability and Transparency, FAT*. vol. 81, pp. 202–214. PMLR (2018)
10. Cao, Y., Chen, X., Yao, L., Wang, X., Zhang, W.E.: Adversarial attacks and detection on reinforcement learning-based interactive recommender systems. In: *Proc. of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR*. pp. 1669–1672. ACM (2020)
11. Celma, Ò.: *Music Recommendation and Discovery - The Long Tail, Long Fail, and Long Play in the Digital Music Space*. Springer (2010)
12. Chen, H., Zhou, K., Lai, K., Hu, X., Wang, F., Yang, H.: Adversarial graph perturbations for recommendations at scale. In: *Proc. of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR*. pp. 1854–1858. ACM (2022)
13. Chen, J., Dong, H., Wang, X., Feng, F., Wang, M., He, X.: Bias and debias in recommender system: A survey and future directions. *ACM Trans. Inf. Syst.* **41**(3), 67:1–67:39 (2023)
14. Christakopoulou, K., Banerjee, A.: Adversarial attacks on an oblivious recommender. In: Bogers, T., Said, A., Brusilovsky, P., Tik, D. (eds.) *Proc. of the 13th ACM Conference on Recommender Systems, RecSys*. pp. 322–330. ACM (2019)
15. Croce, F., Goyal, S., Brunner, T., Shelhamer, E., Hein, M., Cemgil, A.T.: Evaluating the adversarial robustness of adaptive test-time defenses. In: *Proc. of the International Conference on Machine Learning, ICML*. vol. 162, pp. 4421–4435. PMLR (2022)

16. Croce, F., Hein, M.: Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In: Proc. of the 37th International Conference on Machine Learning, ICML. vol. 119, pp. 2206–2216. PMLR (2020)
17. Deldjoo, Y., Bellogín, A., Noia, T.D.: Explaining recommender systems fairness and accuracy through the lens of data characteristics. *Inf. Process. Manag.* **58**(5), 102662 (2021)
18. Fabbri, F., Croci, M.L., Bonchi, F., Castillo, C.: Exposure inequality in people recommender systems: The long-term effects. In: Proc. of the Sixteenth International AAAI Conference on Web and Social Media, ICWSM. pp. 194–204. AAAI Press (2022)
19. Fang, M., Yang, G., Gong, N.Z., Liu, J.: Poisoning attacks to graph-based recommender systems. In: Proc. of the 34th Annual Computer Security Applications Conference, ACSAC. pp. 381–392. ACM (2018)
20. Fenu, G., Marras, M., Medda, G., Meloni, G.: Fair voice biometrics: Impact of demographic imbalance on group fairness in speaker recognition. In: Proc. of the 22nd Annual Conference of the International Speech Communication Association, Interspeech. pp. 1892–1896. ISCA (2021)
21. Floridi, L., Holweg, M., Taddeo, M., Silva, J., Mokander, J., Wen, Y.: capai - a procedure for conducting conformity assessment of ai systems in line with the eu artificial intelligence act. *SSRN Electronic Journal* (01 2022)
22. Ge, Y., Liu, S., Gao, R., Xian, Y., Li, Y., Zhao, X., Pei, C., Sun, F., Ge, J., Ou, W., Zhang, Y.: Towards long-term fairness in recommendation. In: Proc. of the Fourteenth ACM International Conference on Web Search and Data Mining, WSDM. pp. 445–453. ACM (2021)
23. Ge, Y., Tan, J., Zhu, Y., Xia, Y., Luo, J., Liu, S., Fu, Z., Geng, S., Li, Z., Zhang, Y.: Explainable fairness in recommendation. In: Proc. of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR. pp. 681–691. ACM (2022)
24. Ghazimatin, A., Balalau, O., Roy, R.S., Weikum, G.: PRINCE: provider-side interpretability with counterfactual explanations in recommender systems. In: WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining. pp. 196–204. ACM (2020)
25. Gómez, E., Zhang, C.S., Boratto, L., Salamó, M., Marras, M.: The winner takes it all: Geographic imbalance and provider (un)fairness in educational recommender systems. In: Proc. of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR. pp. 1808–1812. ACM (2021)
26. Gómez, E., Zhang, C.S., Boratto, L., Salamó, M., Ramos, G.: Enabling cross-continent provider fairness in educational recommender systems. *Future Gener. Comput. Syst.* **127**, 435–447 (2022)
27. Harper, F.M., Konstan, J.A.: The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.* **5**(4), 19:1–19:19 (2016)
28. He, X., Deng, K., Wang, X., Li, Y., Zhang, Y., Wang, M.: Lightgcn: Simplifying and powering graph convolution network for recommendation. In: Proc. of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR. pp. 639–648. ACM (2020)
29. He, X., He, Z., Du, X., Chua, T.: Adversarial personalized ranking for recommendation. In: Proc. of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR. pp. 355–364. ACM (2018)
30. Kang, B., Lijffijt, J., Bie, T.D.: Explanations for network embedding-based link predictions. In: Proc. of the International Workshops of the European Confer-

- ence on Machine Learning and Principles and Practice of Knowledge Discovery in Databases, ECML PKDD. vol. 1524, pp. 473–488. Springer (2021)
31. Li, B., Wang, Y., Singh, A., Vorobeychik, Y.: Data poisoning attacks on factorization-based collaborative filtering. In: Proc. of the Annual Conference on Neural Information Processing Systems, NeurIPS. pp. 1885–1893 (2016)
 32. Li, Y., Chen, H., Fu, Z., Ge, Y., Zhang, Y.: User-oriented fairness in recommendation. In: Proc. of the Web Conference, TheWebConf. pp. 624–632. ACM / IW3C2 (2021)
 33. Li, Y., Chen, H., Xu, S., Ge, Y., Zhang, Y.: Towards personalized fairness based on causal notion. In: Proc. of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR. pp. 1054–1063. ACM (2021)
 34. Lucic, A., ter Hoeve, M.A., Tolomei, G., de Rijke, M., Silvestri, F.: Cf-gnnexplainer: Counterfactual explanations for graph neural networks. In: Proc. of the International Conference on Artificial Intelligence and Statistics, AISTATS. vol. 151, pp. 4499–4511. PMLR (2022)
 35. Marras, M., Korus, P., Jain, A., Memon, N.D.: Dictionary attacks on speaker verification. *IEEE Trans. Inf. Forensics Secur.* **18**, 773–788 (2023)
 36. Mastropaolo, A., Pascarella, L., Guglielmi, E., Ciniselli, M., Scalabrino, S., Oliveto, R., Bavota, G.: On the robustness of code generation techniques: An empirical study on github copilot. In: Proc. of the 45th IEEE/ACM International Conference on Software Engineering, ICSE. pp. 2149–2160. IEEE (2023)
 37. Medda, G., Fabbri, F., Marras, M., Boratto, L., Fenu, G.: GNNUERS: fairness explanation in gnn for recommendation via counterfactual reasoning. *CoRR abs/2304.06182* (2023)
 38. Mehrabi, N., Naveed, M., Morstatter, F., Galstyan, A.: Exacerbating algorithmic bias through fairness attacks. In: Proc. of the Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI , Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI. pp. 8930–8938. AAAI Press (2021)
 39. Nguyen, T.T., Quach, N.D.K., Nguyen, T.T., Huynh, T.T., Vu, V.H., Nguyen, P.L., Jo, J., Nguyen, Q.V.H.: Poisoning gnn-based recommender systems with generative surrogate-based attacks. *ACM Trans. Inf. Syst.* **41**(3), 58:1–58:24 (2023)
 40. Noia, T.D., Tintarev, N., Fatourou, P., Schedl, M.: Recommender systems under european AI regulations. *Commun. ACM* **65**(4), 69–73 (2022)
 41. Oh, S., Ustun, B., McAuley, J.J., Kumar, S.: Rank list sensitivity of recommender systems to interaction perturbations. In: Proc. of the 31st ACM International Conference on Information & Knowledge Management, CIKM. pp. 1584–1594. ACM (2022)
 42. O’Mahony, M.P., Hurley, N.J., Silvestre, G.C.M.: Recommender systems: Attack types and strategies. In: Proc. of the Twentieth National Conference on Artificial Intelligence and the Seventeenth Innovative Applications of Artificial Intelligence Conference, AAAI. pp. 334–339. AAAI Press / The MIT Press (2005)
 43. Pruksachatkun, Y., Krishna, S., Dhamala, J., Gupta, R., Chang, K.: Does robustness improve fairness? approaching fairness with word substitution robustness methods for text classification. In: Proc. of the Findings of the Association for Computational Linguistics: ACL/IJCNLP. Findings of ACL, vol. ACL/IJCNLP 2021, pp. 3320–3331. ACL (2021)
 44. Qin, T., Liu, T., Li, H.: A general approximation framework for direct optimization of information retrieval measures. *Inf. Retr.* **13**(4), 375–397 (2010)

45. Rahmani, H.A., Naghiaei, M., Dehghan, M., Aliannejadi, M.: Experiments on generalizability of user-oriented fairness in recommender systems. In: Proc. of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR. pp. 2755–2764. ACM (2022)
46. Sato, M., Takemori, S., Singh, J., Ohkuma, T.: Unbiased learning for the causal effect of recommendation. In: Proc. of the Fourteenth ACM Conference on Recommender Systems, RecSys. pp. 378–387. ACM (2020)
47. Singh, A., Joachims, T.: Fairness of exposure in rankings. In: Proc. of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD. pp. 2219–2228. ACM (2018)
48. Solans, D., Biggio, B., Castillo, C.: Poisoning attacks on algorithmic fairness. In: Proc. of the European Conference on Machine Learning and Knowledge Discovery in Database, ECML. LNCS, vol. 12457, pp. 162–177. Springer (2020)
49. Song, J., Li, Z., Hu, Z., Wu, Y., Li, Z., Li, J., Gao, J.: Poisonrec: An adaptive data poisoning framework for attacking black-box recommender systems. In: Proc. of the 36th IEEE International Conference on Data Engineering, ICDE. pp. 157–168. IEEE (2020)
50. Srinivas, S., Subramanya, A., Babu, R.V.: Training sparse neural networks. In: Proc. of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops. pp. 455–462. IEEE (2017)
51. Sun, L., Dou, Y., Yang, C.J., Zhang, K., Wang, J., Yu, P.S., He, L., Li, B.: Adversarial attack and defense on graph data: A survey. *IEEE Trans. Knowl. Data Eng.* **35**(8), 7693–7711 (2023)
52. Wang, S., Zhang, X., Wang, Y., Liu, H., Ricci, F.: Trustworthy recommender systems. *CoRR* **abs/2208.06265** (2022)
53. Wang, X., He, X., Wang, M., Feng, F., Chua, T.: Neural graph collaborative filtering. In: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21–25, 2019. ACM (2019)
54. Wang, Y., Ma, W., Zhang, M., Liu, Y., Ma, S.: A survey on the fairness of recommender systems. *ACM Trans. Inf. Syst.* (jul 2022)
55. Wu, C., Wu, F., Wang, X., Huang, Y., Xie, X.: Fairness-aware news recommendation with decomposed adversarial learning. In: Proc. of the Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI. pp. 4462–4469. AAAI Press (2021)
56. Wu, H., Ma, C., Mitra, B., Diaz, F., Liu, X.: A multi-objective optimization framework for multi-stakeholder fairness-aware recommendation. *ACM Trans. Inf. Syst.* (aug 2022), just Accepted
57. Yang, G., Gong, N.Z., Cai, Y.: Fake co-visitation injection attacks to recommender systems. In: Proc. of the 24th Annual Network and Distributed System Security Symposium, NDSS. The Internet Society (2017)
58. Yang, H., Liu, Z., Zhang, Z., Zhuang, C., Chen, X.: Towards robust fairness-aware recommendation. In: Proc. of the 17th ACM Conference on Recommender Systems, RecSys. pp. 211–222. ACM (2023)
59. Yin, D., Lopes, R.G., Shlens, J., Cubuk, E.D., Gilmer, J.: A fourier perspective on model robustness in computer vision. In: Proc. of the Annual Conference on Neural Information Processing Systems, NeurIPS. pp. 13255–13265 (2019)
60. Yuan, F., Yao, L., Benatallah, B.: Exploring missing interactions: A convolutional generative adversarial network for collaborative filtering. In: Proc. of the 29th ACM International Conference on Information and Knowledge Management, CIKM. pp. 1773–1782. ACM (2020)

61. Yue, Z., Zeng, H., Kou, Z., Shang, L., Wang, D.: Defending substitution-based profile pollution attacks on sequential recommenders. In: Proc. of the Sixteenth ACM Conference on Recommender Systems, RecSys. pp. 59–70. ACM (2022)
62. Zeng, Z., Tan, H., Zhang, H., Li, J., Zhang, Y., Zhang, L.: An extensive study on pre-trained models for program understanding and generation. In: Proc. of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA. pp. 39–51. ACM (2022)
63. Zhang, K., Cao, Q., Sun, F., Wu, Y., Tao, S., Shen, H., Cheng, X.: Robust recommender system: A survey and future directions. CoRR **abs/2309.02057** (2023)
64. Zhang, S., Yin, H., Chen, T., Huang, Z., Cui, L., Zhang, X.: Graph embedding for recommendation against attribute inference attacks. In: Proc. of The Web Conference 2021, TheWebConf. pp. 3002–3014. ACM / IW3C2 (2021)
65. Zhang, S., Yin, H., Chen, T., Nguyen, Q.V.H., Huang, Z., Cui, L.: Gcn-based user representation learning for unifying robust recommendation and fraudster detection. In: Proc. of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR. pp. 689–698. ACM (2020)
66. Zhang, Y., Chen, X.: Explainable recommendation: A survey and new perspectives. *Found. Trends Inf. Retr.* **14**(1), 1–101 (mar 2020)
67. Zhao, W.X., Mu, S., Hou, Y., Lin, Z., Chen, Y., Pan, X., Li, K., Lu, Y., Wang, H., Tian, C., Min, Y., Feng, Z., Fan, X., Chen, X., Wang, P., Ji, W., Li, Y., Wang, X., Wen, J.: Recbole: Towards a unified, comprehensive and efficient framework for recommendation algorithms. In: Proc. of the 30th ACM International Conference on Information and Knowledge Management, CIKM. pp. 4653–4664. ACM (2021)
68. Zheng, J., Ma, Q., Gu, H., Zheng, Z.: Multi-view denoising graph auto-encoders on heterogeneous information networks for cold-start recommendation. In: Proc. of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD. pp. 2338–2348. ACM (2021)