



A unifying and general account of fairness measurement in recommender systems

Enrique Amigó^a, Yashar Deldjoo^b, Stefano Mizzaro^c, Alejandro Bellogín^{d,*}

^a Universidad Nacional de Educación a Distancia (UNED), Madrid, Spain

^b Polytechnic University of Bari, Bari, Italy

^c University of Udine, Udine, Italy

^d Universidad Autónoma de Madrid, Madrid, Spain

ARTICLE INFO

CCS Concepts:

Information systems: Specialized information retrieval

Keywords:

Recommender systems
Fairness
Biases
Information theory

ABSTRACT

Fairness is fundamental to all information access systems, including recommender systems. However, the landscape of fairness definition and measurement is quite scattered with many competing definitions that are partial and often incompatible. There is much work focusing on specific – and different – notions of fairness and there exist dozens of metrics of fairness in the literature, many of them redundant and most of them incompatible. In contrast, to our knowledge, there is no formal framework that covers all possible variants of fairness and allows developers to choose the most appropriate variant depending on the particular scenario. In this paper, we aim to define a general, flexible, and parameterizable framework that covers a whole range of fairness evaluation possibilities. Instead of modeling the metrics based on an abstract definition of fairness, the distinctive feature of this study compared to the current state of the art is that we start from the metrics applied in the literature to obtain a unified model by generalization. The framework is grounded on a general work hypothesis: interpreting the space of users and items as a probabilistic sample space, two fundamental measures in information theory (Kullback–Leibler Divergence and Mutual Information) can capture the majority of possible scenarios for measuring fairness on recommender system outputs. In addition, earlier research on fairness in recommender systems could be viewed as single-sided, trying to optimize some form of equity across either user groups or provider/procurer groups, without considering the user/item space in conjunction, thereby overlooking/disregarding the interplay between user and item groups. Instead, our framework includes the notion of statistical independence between user and item groups. We finally validate our approach experimentally on both synthetic and real data according to a wide range of state-of-the-art recommendation algorithms and real-world data sets, showing that with our framework we can measure fairness in a general, uniform, and meaningful way.

1. Introduction

The notion of fairness has recently attracted considerable attention. Fairness is studied in general in artificial intelligence and machine learning, typically focusing on classification problems (Mehrabian, Morstatter, Saxena, Lerman, & Galstyan, 2021; Verma & Rubin, 2018), and also in Information Retrieval (IR), with a focus on fair rankings (Celis, Straszak, & Vishnoi, 2018; Diaz, Mitra,

* Corresponding author.

E-mail addresses: enrique@lsi.uned.es (E. Amigó), yashar.deldjoo@poliba.it (Y. Deldjoo), mizzaro@uniud.it (S. Mizzaro), alejandro.bellogin@uam.es (A. Bellogín).

<https://doi.org/10.1016/j.ipm.2022.103115>

Received 15 May 2022; Received in revised form 6 October 2022; Accepted 7 October 2022

Available online 1 November 2022

0306-4573/© 2022 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Ekstrand, Biega, & Carterette, 2020; Ekstrand, Das, Burke, & Diaz, 2022). However, in the field of Recommender Systems (RSs) the notion of fairness becomes multi-faceted and arguably presents a richer scenario (Abdollahpour et al., 2020; Burke, 2017; Ekstrand et al., 2022; Li, Ge, & Zhang, 2021). To evaluate if a RS is fair, one must take into account a variety of factors, including the stakeholders (consumers, producers, side-stakeholders), the kind of benefit impacting the consumers and businesses/producers (perceived utility, item exposure), the context, morality, time among other variables (Verma & Rubin, 2018). For instance, almost every online platform we interact with, like Spotify and Amazon, functions as a marketplace connecting consumers with product producers or service providers. From the consumers' perspective, fairness mostly concerns an even distribution of effectiveness among users, avoiding the penalization of protected groups like, for example, female or black candidates in job applications.¹ Conversely, producers and item providers, who seek increased visibility, are primarily concerned with exposure fairness that should not be penalized, for example, on the basis of producers' popularity or country. Let us also remark that a fair system might provide unequal distribution of resources, as receiving a privilege can be based on merits and needs (Deldjoo et al., 2021) or fitness (Moulin, 2003). Given the complexity of such a scenario, it is not surprising that the notion of fairness in RSs lacks a unified understanding. There are many definitions, which are different if not even incompatible (Ekstrand et al., 2022; Verma & Rubin, 2018).

The fact is that measuring fairness has different facets, such as the consumer or producer perspective, modeling benefit in terms of exposure (Dong, Xie, & Li, 2021; Gómez, Boratto, & Salamó, 2022; Wundervald, 2021; Yalcin & Bilge, 2021) or utility (Abdollahpour, Mansoury, Burke, Mobasher, & Malthouse, 2021; Li, Chen, Fu, Ge, & Zhang, 2021; Wang & Joachims, 2021; Zliobaite, 2017) between groups, target benefit distribution in terms of equity or merit based, etc. To our knowledge, there is no formal framework that covers all possible variants of fairness and allows developers to choose the most appropriate variant depending on the particular scenario. The result is that in many cases the authors do not identify the most appropriate metric and in many other cases different authors apply different metrics for the same purpose. This situation is of course an obstacle to progress in the field. In perspective, providing a uniform, general, standard, and unified account of fairness in RSs would be instrumental to remove such an obstacle, and that is precisely the aim of this paper.

The main contribution of this paper is the definition of a general, flexible, and parameterizable framework that covers most possibilities in fairness measurement. There is prior research that analyzes existing fairness metrics comprehensively based on a set of dimensions (Wang, Ma, Zhang, Liu, & Ma, 2022). Theoretical frameworks of metrics that are adaptable to various circumstances have also been outlined in Deldjoo et al. (2021), Kirnap et al. (2021), Sacharidis, Mouratidis, and Klefogiannis (2019) and Wu, Mitra, Ma, Diaz, and Liu (2022). Compared to previous approaches, our proposal makes the following specific contributions:

- We define such a novel framework on the basis of a comprehensive analysis of existing metrics via categorization dimensions. Thus, we start from the metrics applied in the literature to obtain a unified model by generalization, rather than starting from a unique abstract fairness definition.
- We show that by modeling exposure, utility, and effectiveness as probability distributions over the user/item space, it is possible to capture most existing fairness metrics by means of two information theory measures, namely Kullback–Leibler Divergence and Mutual Information.
- The framework allows to model some features that are absent in previously proposed metrics, such as the independence between user/item groups or individuals regardless of any ideal target benefit distribution.
- Besides, on the basis of its coverage of existing metrics, this flexible framework is validated over a synthetic data set and recommender system outputs artificially defined to cover different fairness strengths and weaknesses. The framework behavior is also checked on real data sets.

More in detail, we adopt the following methodology. To be able to comprehensively analyze the existing metrics, we first establish five dimensions along which the metrics can be classified (Section 2). As a second step, we then perform a classification of the existing fairness metrics, still focusing on RSs but also relating to more classical fairness definitions in classification (Section 3). Although the goal of this paper is not to serve as a systematic or semi-systematic literature review, such a thorough analysis of the literature to analyze a vast majority of the existing metrics allows us to show that the above proposed five dimensions can be used to describe in an organized and coherent way the fairness metrics landscape. As a third step, we identify the most flexible metric models in the literature to propose a general and formal framework that is based on information theory and allows us to measure fairness in a unified way. We obtain a framework that, based on a series of questions, allows us to identify the most appropriate metric for each specific scenario (Section 4). As a fourth and final step, the theoretical work is complemented with an experimental analysis on both synthetic and real data where we check the behavior of all the metrics derived from the proposed framework, and evaluate the fairness of state-of-the-art recommendation algorithms, including classical and neural algorithms, tested on real-world data sets (Section 5).

2. Fairness dimensions

By analyzing the vast literature on fairness, one can note a variety of different approaches. Some research works focus on specific notions of fairness. Some others attempt to include different fairness notions. Others appear to have given up hope for a singular definition of fairness, conceding that «the English word “fairness” will need a multitude of definitions» (Krakovsky, 2022, p. 11). To

¹ In this work, we will frequently use the phrases user or customer fairness, item or producer or supplier fairness, and protected or sensitive features interchangeably.

make a historical comparison, the situation is not different from that about the notion of relevance in IR in the 1970s. At that time, the seminal paper by Saracevic (1975) was a breakthrough that helped to clarify (i) the existence of different kinds of relevance and (ii) the possibility of classifying all of them under a common framework, by identifying some classification dimensions. This study is an effort to accomplish the same for the concept of fairness, which we feel is possible or at least worth exploring.

After a careful exploratory process for the compilation of fairness metrics, we propose to define five orthogonal dimensions (D1–D5) as independent categorization criteria on which to categorize the existing metrics on recommender systems' fairness evaluation. That is, each dimension has a number of variants associated with it. Ideally, a fairness assessment framework should be flexible enough to cover all combinations of variants of the different dimensions. The dimensions defined in this paper have been compiled from different works (Ekstrand et al., 2022; Li et al., 2022; Wang et al., 2022). The five dimensions are described below; for each of them we provide a name, a set of possible values, and a brief description.

- **D1—Benefit (Exposure, Effectiveness).** The first dimension concerns the type of benefit that needs to be distributed in a fair way. It can essentially take two alternative forms. The first one is to what extent items are exposed to users. Note that some previous research has named the exposure binary case as “visibility” and the ranking case as “exposure” (Boratto, Fenu, & Marras, 2021; Gómez, Zhang, Boratto, Salamó, & Marras, 2021); herein with “exposure” we mean both of them. The second one, i.e., the effectiveness benefit criterion, is to what extent this exposure is useful to the user. Wang et al. (2022) named this dimension as *treatment* (exposure) versus *impact* (effectiveness) *optimization object*.
- **D2—Stakeholder (Users, Items/Providers).** A core characteristic of RSs is the duality of user- and vendor-centered utility (Deldjoo et al., 2021), also known as user/consumers and provider/producers fairness in the literature, or for short C-Fairness and P-Fairness (Burke, Sonboli, & Ordonez-Gauger, 2018).² They aim at a fair treatment of the users and of the producers, respectively. Both can be considered simultaneously, which is called *two-sided fairness* (Abdollahpouri & Mansoury, 2020; Ekstrand et al., 2022; Li, Chen, Xu, Ge, & Zhang, 2021; Wu, Cao, Xu, & Tan, 2021) and some authors extend this idea to multi-stakeholders, including the own system interest (da Silva, Manzato, & Durão, 2021). This dimension has been referred as *subject* (Li et al., 2022; Wang et al., 2022) and also as *Consumer vs. provider fairness* (Ekstrand et al., 2022).
- **D3—Partition Granularity (Two Groups, Many Groups, Individuals).** Fairness usually entails comparing, on average, the benefit received by the members of different groups. The granularity of the partition into groups can vary along a spectrum: at one extreme, only two groups are defined, i.e., privileged and unprivileged groups as defined by their protected attributes; in more general cases, there are several groups over which maintain equity; and at the other extreme, fairness is studied between individuals, e.g., equal recommendation effectiveness for each user. In general, fairness covers everything from the division of user or item spaces into two groups, to many groups, to consider individuals (which subsumes the previous two cases). We will see that not all metrics capture all possibilities. Wang et al. (2022) named this dimension as *target*.
- **D4—Exposure Scheme (Rating, Set, Ranking).** The existing fairness metrics differ depending on how the items are displayed to the user. Capturing different information access user interfaces is crucial for the generality of fairness measurement. In some cases, the RS exposes the items to the users according to an estimated rating (e.g., 1 to 5 scale). In other cases, the user interface consists of a set of recommended items without any priority order. In most cases, items are organized in a ranking, or into a ranking of categories, or even a ranking of rankings (Gunawardana & Shani, 2015). The fairness measurement framework should be able to weight the exposure of items in all these situations. We will see that most current metrics are oriented to specific exposure schemes, while others encapsulate this dimension in an exposure weight parameter. This dimension has been referred to as *provider representation measure* (Ekstrand et al., 2022; Kirnap et al., 2021) or *attention* (Ghosh, Dutt, & Wilson, 2021).
- **D5—Fairness Criterion (Parity, Size Proportionality, Utility Proportionality, Independence).** This last dimension concerns the overall criterion used to state that distributing in a certain way the benefits across individual users/items or groups of users/items is fair. According to Kirnap et al. (2021), there are three main ways to define such a target distribution, i.e., on the basis of parity, proportionality to the corpus presence, and proportionality to utility (Kirnap et al., 2021). Parity implies that all groups receive the same exposure mass, proportionality implies that the exposure is proportional to the group size, and utility implies that the exposure is proportional to the relevance mass of the group. Deciding the target benefit distribution that enables a reasonable allocation of resources can be task-specific (Deldjoo et al., 2021) and extremely complex for the system designer, since it can follow norms of short-term user satisfaction, long-term business growth, morality, among others.

Besides these fairness criteria, fairness can be measured on the basis of statistical *Independence* of user's or item's protected attributes. For example, in a job recommendation setting, the exposure to executive vs. low qualified jobs could be required to be independent of protected characteristics of users such as their age and gender, which is equal to say that the resource allocation should be unbiased by protected characteristics. As such, we can observe a connection between statistical independence and treatment disparity (Deldjoo et al., 2021; Ekstrand et al., 2022; Singh & Joachims, 2018), which embodies the idea that the system should make judgments (exposure) regardless of the individual's protected attributes.

Note that we do not claim that this set of dimensions is complete. However, the analysis in the next section shows that this dimension set is enough to capture the limitations of existing metrics in terms of scenario coverage. Therefore, we can use such a set for evaluating the generality of fairness measurement approaches. The reader is referred to Section 6 for a discussion on limitations and outlooks.

² Note that in the literature, consumer and user fairness are frequently used as synonyms. The same applies for item fairness, provider fairness, and producer fairness. Also, in situations where the roles of the user and the items are reversed, such as recommendation or people for a certain job (Abel, Deldjoo, Elahi, & Kohlsdorf, 2017), users and items need to be swapped.

3. Fairness metrics

Wang et al. presented an interesting survey (Wang et al., 2022) where they categorized and defined many metrics for fairness evaluation in recommender systems. Rather than presenting a comprehensive catalog of metrics and their definitions, we aim at analyzing to what extent the metrics proposed in the literature are general or if they are actually limited to different particular scenarios. More specifically, we are interested in identifying those metric schemes that allow to capture diverse fairness scenarios. Unfortunately, the number of existing metrics is very large and there would not be space in the article to include their definitions. Since the focus of this article is to evaluate the scenario coverage of the metrics, we are interested in including as many as possible and at least describe their properties in terms of coverage over possible fairness scenarios.

As a result, Table 1 summarizes how each metric or measurement approach (on the rows) captures each variant under the five dimensions presented above (columns). For each column, the meaning of the symbols is as follows: Exp and Eff (D1) mean exposure and effectiveness oriented; Us and It (D2) represent user and item (provider) serving as RSs' main stakeholders; 2g (D3) represents that the metric is defined for two groups (protected and non protected), ng represents that the metric can be defined over many groups and I represent that the metric is defined only for individuals (if all granularity levels can be captured, including individuals, we use the ✓ symbol); Rat, Set, and Rank (D4) mean that the exposure is set-, ranking-, or rating-based; P, S, U, and I (D5) represent that the fairness criterion is based on Parity, Size proportionality, Utility proportionality, or Independence. The tick symbol (✓) represents that the metric can be customized into all the variants of the corresponding dimension.

We describe each group of metrics in each of the following subsections, as indicated in the table; within each subsection we generally follow the same metric order as in the table, and we group existing metrics according to D1 (benefit) and D5 (fairness criterion). We independently analyze metrics that are based on exposure but consider utility proportionality as fairness criterion. We also consider separately some general frameworks and some metrics that capture the independence fairness criterion.

Since the scope of this article is the space of metrics that quantify fairness on the basis of system outputs rather than the recommendation algorithm process, we believe our work aligns more with the *outcome fairness* rather than *process fairness* (Wang et al., 2022). The second evaluates aspects such as what data has been used, under what principles the system makes decisions, or what are the causal relationships between inputs and outputs. In contrast, outcome fairness ignores how the system works internally and focuses on the fair distribution of benefits. Finally, we feel that outcome fairness is an established topic with a large number of metrics, making this an ideal time to build a generalization; in contrast, the many perspectives on process fairness are still in the early stages of research.

3.1. Fairness metrics based on exposure

Fairness in terms of exposure in RSs is related to previous fairness metrics for classification in Artificial Intelligence, such as *Fairness Based on Predicted Outcome* (Verma & Rubin, 2018) – also called *Statistical Parity* (Dwork, Hardt, Pitassi, Reingold, & Zemel, 2012), *Equal Acceptance Rate* (Zliobaite, 2015), and *Benchmarking* (Simoiu, Corbett-Davies, & Goel, 2016). A classifier satisfies these definitions if the probability of being assigned to the positive predicted class is equal across different item groups. In the case of RSs, we can instead speak of item exposure. This set of metrics includes those that evaluate the equity of exposure across user or item groups (D1 = Exp). In general, the exposure fairness in RSs is commonly defined from the item side. One reason is that the item providers are interested in gaining visibility.

We start by identifying a set of IR metrics that focus on the relative presence of protected and non-protected item groups (D2 = It) in the top- k ranking positions. *Ranked Group Fairness Condition* (Zehlike et al., 2017) and the *Fairness Constraint* (Celis et al., 2018) specify upper and lower bounds on the number of items from each group that are allowed to appear in the top- k positions of the ranking; both are parity oriented (D5 = P).

In other metrics, the fairness criterion is size-proportional (D5 = S). For instance, Yang and Stoyanovich (2017) proposed three metrics, namely *Normalized Discounted Difference* (rND), *Divergence* (rKL), and *Ratio* (rRD), that compare the distribution of the protected group above a certain ranking position with the group presence in the corpus. Geyik et al. (2019) proposed the metric *Skew@k* which is similar, and compares protected with non-protected groups. A common feature of all these metrics is that the fairness score is averaged across ranking position thresholds. On the contrary, other rank oriented metrics such as *Rank Parity* (Kuhlman et al., 2019) quantify exposure in terms of the cases in which items from one group are ranked above another group. The *Disparate Exposure* proposed by Boratto et al. (2021) computes the difference between the minority group representation in the item catalog and the average exposure taking into account the ranking positions of items.

A common limitation of all the previous metrics is that they are defined for two groups (D3 = 2g). The *Attention Bias Ratio* (Ghosh et al., 2021) addresses this limitation by quantifying the disparity between the groups with the lowest and highest mean exposure, considering the ranking position bias. Another metric which is able to capture multiple groups is the *Product Ranking Fairness* which computes the Kullback–Leibler Divergence (KLD) between the amount of top ranked items and the item group size (Wan et al., 2020). A similar metric is NDKL which aggregates KLD values across ranking position thresholds (Geyik et al., 2019).

In the context of RSs we found a set of metrics that are able to capture multiple groups (D3 = ✓), but limited to set-based exposure (D4 = Set). Some metrics study the exposure variance across groups. Some examples are the entropy-based metric *Inequality in Producer Exposures* (Patro et al., 2020), the *Uniform Fairness Variance* (Wu et al., 2021), the *Equity of Attention for Group Fairness* (Gharahighchi et al., 2021), the *Gini Index* (Ge et al., 2021), or the *Jain's fairness index* (Zhu et al., 2020). The *Fraction of Satisfied Producers* (Patro et al., 2020), the *Average Provider Coverage rate* (APCR) (Liu et al., 2019), and the *Group Fairness*

Table 1

Dimensions captured by fairness measures. Meaning of symbols is as follows. Exp and Eff (D1): exposure and effectiveness oriented. Us and It (D2): user and item (provider) as main stakeholders. 2g, ng, and I (D3): two groups (protected and non protected), many groups, and individuals (if all granularity levels can be captured, including individuals, we use the ✓ symbol). Rat, Set, and Rank (D4): exposure is set-, ranking-, or rating-based. P, S, U, and I (D5): fairness criterion is based on Parity, Size proportionality, Utility proportionality, or Independence. ✓: the metric can be customized into all the variants of the corresponding dimension.

| | D1: Benefit (Exp/Eff) | D2: Stakeholder (Us/It) | D3: Partition (2g/ng/I) | D4: Exposure (Rat/Set/Rank) | D5: Criterion (P/S/U/I) |
|---|--------------------------|----------------------------|----------------------------|--------------------------------|----------------------------|
| Fairness measures based on exposure (Section 3.1) | | | | | |
| Ranked Group Fairness Condition (Zehlike et al., 2017) | Exp | It | 2g | Rank | P |
| Fairness Constraint (Celis et al., 2018) | Exp | It | 2g | Rank | P |
| rND, rKL, rRD (Yang & Stoyanovich, 2017) | Exp | It | 2g | Rank | S |
| Skew@k (Geyik, Ambler, & Kenthapadi, 2019) | Exp | It | 2g | Rank | S |
| Rank Parity (Kuhlman, Valkenburg, & Rundensteiner, 2019) | Exp | It | 2g | Rank | S |
| Disparate Exposure (Boratto et al., 2021) | Exp | It | 2g | Rank | S |
| Attention Bias Ratio (Ghosh et al., 2021) | Exp | It | ✓ | Rank | S |
| Product Ranking Fairness (Wan, Ni, Misra, & McAuley, 2020) | Exp | It | ✓ | Rank | S |
| NKLD (Geyik et al., 2019) | Exp | It | ✓ | Rank | S |
| Inequality in Producer Exposures (Patro, Biswas, Ganguly, Gummadi, & Chakraborty, 2020) | Exp | It | ✓ | Set | P |
| Uniform Fairness Variance (Wu et al., 2021) | Exp | It | ✓ | Set | P |
| Equity of Attention for Group Fairness (Gharahighehi, Vens, & Pliakos, 2021) | Exp | It | ✓ | Set | P |
| Gini Index (Ge et al., 2021) | Exp | It | I | Set | P |
| Jain's fairness index (Zhu, Sun, Li, & Wang, 2020) | Exp | It | I | Set | P |
| Fraction of Satisfied Producers (Patro et al., 2020) | Exp | It | ✓ | Set | P |
| Average Provider Coverage Rate (Liu, Guo, Sonboli, Burke, & Zhang, 2019) | Exp | It | ✓ | Set | P |
| Group Fairness Measure (Mehrotra, McInerney, Bouchard, Lalmas, & Diaz, 2018) | Exp | It | ✓ | Set | P |
| Supplier Popularity Deviation (Gharahighehi et al., 2021) | Exp | It | ✓ | Set | S |
| MAD (Zhu, Hu, & Caverlee, 2018) | Exp | It | 2g | ✓ | S |
| Non-Parity Fairness (Yao & Huang, 2017) | Exp | It | 2g | ✓ | S |
| Demographic Parity (Singh & Joachims, 2018) | Exp | It | 2g | ✓ | S |
| Gupta et al. (2021) | Exp | It | ✓ | ✓ | S |
| Fairness measures based on effectiveness (Section 3.2) | | | | | |
| Absolute Difference (Zhu et al., 2018) | Eff | Us | 2g | ✓ | S |
| KS statistic (Zhu et al., 2018) | Eff | Us | 2g | ✓ | S |
| Effectiveness Standard Deviation (Patro et al., 2020; Wu et al., 2021) | Eff | Us | ✓ | ✓ | S |
| Rating Prediction Fairness (Wan et al., 2020) | Eff | ✓ | ✓ | Rat | S |
| Wang and Joachims (2021) | Eff | Us | ✓ | ✓ | S |
| Yao and Huang (2017) | Eff | Us | 2g | Rat | S |
| User Bias (Lin, Liu, Xv, & Li, 2021) | Eff | Us | 2g | ✓ | S |
| Pairwise Fairness (Beutel et al., 2019) | Eff | It | ng | ✓ | S |
| Item Bias (Lin et al., 2021) | Eff | It | 2g | ✓ | S |
| Disparate Relevance (Boratto et al., 2021) | Eff | It | 2g | Rank | S |
| Fairness measures based on utility-equalized exposure (Section 3.3) | | | | | |
| Supplier Popularity Deviation (Abdollahpour et al., 2020) | Exp | It | ✓ | Set | U |
| Mean Average Calibration (da Silva et al., 2021) | Exp | It | ✓ | Set | U |
| JS-Divergence (Modani, Jain, Soni, Gupta, & Agarwal, 2017) | Exp | It | ✓ | Set | U |
| Rank Equality (Kuhlman et al., 2019) | Exp | It | 2g | Rank | U |
| Steck (2018) | Exp | It | ✓ | ✓ | U |
| Equity of Amortized Attention (Biega, Gummadi, & Weikum, 2018) | Exp | It | ✓ | ✓ | U |
| Quality Weighted Fairness (Wu et al., 2021) | Exp | It | ✓ | ✓ | U |
| Disparate Treatment Ratio (Singh & Joachims, 2018) | Exp | It | 2g | ✓ | U |
| Flexible fairness measures (Section 3.4) | | | | | |
| Wu et al. (2022) | ✓ | ✓ | ✓ | Set | P/S/U |
| Kirnap et al. (2021) | Exp | It | ✓ | ✓ | P/S/U |
| Sacharidis et al. (2019) | Exp | ✓ | ✓ | ✓ | P/S/U |
| Deldjoo et al. (2021) | ✓ | ✓ | ✓ | ✓ | P/S/U |
| Fairness measures based on independence (Section 3.5) | | | | | |
| Relative Opportunity (Burke et al., 2018) | Exp | ✓ | 2g | Set | I |
| Bias Disparity (Tsintzou, Pitoura, & Tsaparas, 2019) | Exp | ✓ | 2g | Set | I |

Measure (Mehrotra et al., 2018) are similar but based on the number of providers (item groups) covered in single user lists. The *Supplier Popularity Deviation* (Gharahighehi et al., 2021) is also similar, but taking into account the item group size ($D5 = S$).

On the other hand, the absolute difference between mean ratings of two groups (MAD), which is extended with the Kolmogorov–Smirnov statistic (Zhu et al., 2018), captures graded exposure (including set and ranking, $D4 = \checkmark$), but it is defined for two groups ($D3 = 2g$). The same applies to the *Non-Parity Fairness* (Yao & Huang, 2017) and the *Demographic Parity* (Singh & Joachims, 2018). Finally, the *Demographic Disparity* (Gupta et al., 2021) is computed as the maximum difference of exposure between group pairs, where exposure includes the ranking discount function. A common property of these previous RS fairness metrics is that they are size proportional ($D5 = S$), i.e., group exposure is normalized according to the group size.

3.2. Fairness metrics based on effectiveness

Recommendation effectiveness is a natural benefit function. This links with the classification fairness notion *Predictive Parity*: “both protected and unprotected groups have equal probability of a subject with positive predictive value to truly belong to the positive class” (Verma & Rubin, 2018). It is also equivalent to *Outcome Test* (Simoiu et al., 2016), *Equal opportunity* (Chen et al., 2020; Hardt, Price, & Srebro, 2016; Kusner, Loftus, Russell, & Silva, 2017), and *False negative error rate balance* (Chouldechova, 2017): positive samples from different groups have equal probability to be classified as positive.

One major consumer-side group fairness problem is to determine whether the system provides comparable quality of service or utility to different groups of consumers. This family of metrics includes those that focus on the equity of recommendation effectiveness across user groups ($D1 = \text{Eff}$, $D2 = \text{Us}$). In general, since these metrics work with expected effectiveness of individual users, the fairness criterion is size proportionality ($D5 = S$).

In the contexts of IR and RSs, a common fairness evaluation procedure in the literature consists in comparing the expected effectiveness of user groups (Ekstrand & Mahant, 2017; Mehrotra et al., 2017), for instance, using the *Absolute Difference* (Zhu et al., 2018) or the Kolmogorov–Smirnov statistic (Zhu et al., 2018). The standard deviation of expected effectiveness across individuals or user groups is a common way to quantify fairness (Patro et al., 2020; Wu et al., 2021), allowing multiple groups ($D3 = \checkmark$). This method is agnostic regarding the effectiveness metric, as it captures both set and ranking based exposition ($D4 = \checkmark$). Similarly, the *Rating Prediction Fairness* proposed by Wan et al. (2020) applies the ANOVA test over the null hypothesis of independence between prediction errors and market segments. As well as accepting multiple groups ($D3 = \checkmark$), this method allows to define market segments over both user and items ($D2 = \checkmark$) but it is only rating oriented ($D4 = \text{Rat}$).

Yao and Huang (2017) defined a set of alternative metrics, namely, *Value Unfairness*, *Absolute Unfairness*, *Underestimation Unfairness*, *Overestimation Unfairness*, and *Non Parity*. They all compare, for each item, the expected score for disadvantaged and advantaged users. They are limited to two user groups ($D3 = 2g$) and top ranking heaviness is not captured since they define the recommendation problem as an item rating prediction problem ($D4 = \text{Rat}$). Wang and Joachims (2021) defined a user fairness metric that quantifies the effectiveness equity across multiple user groups through a social-welfare function. It captures multiple groups ($D3 = \checkmark$) and graded exposure ($D4 = \checkmark$). The *User Bias* proposed by Lin et al. (2021) can be also applied to any effectiveness metric ($D4 = \checkmark$), but it is defined for only two groups ($D3 = 2g$).

From the provider perspective, one major group fairness problem is to determine whether the system provides comparable quality of service or utility to different providers, i.e., useful items from different providers have equal opportunity to be exposed. However, metrics for this aspect are not very common in RSs. An exception is the *Pairwise Fairness* which computes the probability that a useful (clicked in the user feedback) item is ranked above another useless item within a certain item group (Beutel et al., 2019). It allows to compare multiple item groups ($D3 = ng$) but not individual items, and the target distribution is proportional to size since it is defined as a probability ($D5 = S$). Another exception is the *Item Bias* (Lin et al., 2021), which computes the difference between effectiveness metrics over two item sets ($D3 = 2g$). Finally, the *Disparate Relevance* proposed by Boratto et al. (2021) is somewhat particular; it computes the difference between the minority group representation in the item catalog and the estimated relevance of their exposed items.

3.3. Fairness metrics based on utility-equalized exposure

In some cases a uniform exposure distribution is not fair. It is natural to think that the exposure of suppliers should be proportional to the amount of useful items they provide. This is related to classification fairness metrics such as *Equalized odds* (Hardt et al., 2016), *conditional procedure accuracy equality* (Berk, Heidari, Jabbari, Kearns, & Roth, 2017), and *disparate mistreatment* (Zafar, Valera, Gomez-Rodriguez, & Gummadi, 2017): protected and unprotected groups have equal true positive rate, i.e., the probability of true instances (useful items in RSs) to be classified as true (exposed items in RSs). The benefit function is exposure ($D1 = \text{Exp}$) but the ideal distribution is related to item utility ($D5 = \text{U}$).

Some utility equalized exposure metrics are oriented to set exposure ($D4 = \text{Set}$). For instance, *Supplier Popularity Deviation* (Abdollahpour & Mansoury, 2020) and *Mean Average Calibration* (da Silva et al., 2021) sum the absolute differences between the ratio of recommendations and ratings that come from items of supplier. There exist some ranking oriented utility equalized metrics, like *Rank Equality* (Kuhlman et al., 2019) that computes the number of times an item of a group is falsely given a higher rank than an item of another group. It can be applied to two item groups ($D3 = 2g$). Modani et al. used the Jensen–Shannon Divergence to compare the exposure and the utility provided by item groups (Modani et al., 2017).

Other approaches are agnostic as to the type of exposure function (set, ranking, etc.). Steck defined a metric in terms of KLD between exposure weight and utility (according to the user’s previous preferences) of item groups (Steck, 2018). In the context of IR,

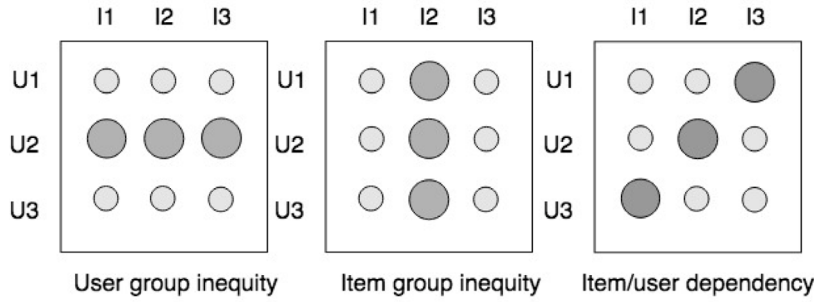


Fig. 1. The notion of independence based fairness. Each circle represents the amount of benefit (e.g., exposure) for each item (column) and user group (row).

Equity of Amortized Attention (Biega et al., 2018) is based on the L1-norm distance between accumulated exposure and relevance of single item groups. The *Disparate Treatment Ratio* compares ratios of exposure with utility of group pairs (Singh & Joachims, 2018): it is applied to two item groups ($D3 = 2g$). The *Quality Weighted Fairness* computes the variance of exposure/utility ratios across item groups, capturing the rank exposure bias and multiple item groups (Wu et al., 2021).

3.4. Flexible fairness metrics

Some authors have proposed flexible fairness measurement models that can be instantiated to particular scenarios. The common feature of these approaches is that the user/item bidimensional space is divided according to user or item groups. Then, the benefit distribution across groups is compared against an ideal (fair) distribution. These models are flexible enough to consider any fairness criterion ($D5 = P/S/U$) and one or many user/item groups ($D3 = \checkmark$).

In this line, Wu et al. (2022) proposed to compute the average differences. The effectiveness benefit function is implicitly captured by considering item relevance as target exposure. The limitations of this framework is that the management of ranking exposure is not specified ($D4 = \text{Set}$). In the context of IR, Kirnap et al. (2021) proposed a general theoretical framework consisting of: (i) the exposure distribution, which decreases with rank according to decay functions in IR evaluation metrics (Diaz et al., 2020); (ii) the target distribution, which can be *parity*, *proportionality to the corpus presence*, or *proportionality to the relevance* ($D5 = P/S/U$); (iii) the similarity between the exposure distribution of groups and the target distribution across different rank thresholds in terms of KLD or other distribution similarity metrics. In the context of RSs, the framework proposed by Sacharidis et al. (2019) is very similar. It computes the KLD between the real and the desirable exposure distribution across user or item groups. However, the way of managing ranking exposure is not specified and the model is limited to the exposure benefit criterion ($D1 = \text{Exp}$). In fact, the KLD has been used by different authors to compare the actual and the desirable benefit distribution across groups (Ge et al., 2022; Li et al., 2022; Stoica, Riederer, & Chaintreau, 2018). Deldjoo et al. (2021) proposed another similar framework. The utility distribution across user/item groups is compared with the ideal distribution via the *Generalized Cross Entropy* which is more robust against outliers. The fairness metric can be instantiated into an exposure fairness metric by considering all items equally useful. In addition, the model by Deldjoo et al. captures the effectiveness benefit function ($D1 = \checkmark$) and ranking exposure weighting ($D3 = \checkmark$).

3.5. Fairness metrics based on independence

All the metrics discussed in the previous subsections refer to equity across user or item groups. However, there also exists the notion of independence between user and item groups which is related with treatment disparity. As an example, let us assume that both men and women receive the same amount of effective items (e.g., jobs), and that all item groups are equally distributed in terms of exposure quantity and quality to users. Even in this situation, it may be the case that the user genre conditions the type of recommended jobs. Fig. 1 illustrates the notion of independence based fairness. Each panel represents the benefit distribution (e.g., exposure) across three user groups and three item groups. The larger the circles, the more the corresponding item group is exposed to the corresponding user group. In the left case there is no parity across user groups (the user group U2 receives more exposure). In the second distribution, there is no parity across item groups (the item group I2 is exposed to a larger extent). In the third case, there is parity across user and item groups (two small and one big circle for each column or row). However, user groups U1, U2, and U3 are biased to item groups I3, I2, and I1 respectively. In other words, since the benefit across user and item groups is not independent, there is treatment disparity.

In a more abstract way, we can say that the metrics described previously quantify the equity across user or item groups, whereas the metrics we are considering now deal with the dependence between attributes of users and items. Consequently, there is no target distribution in independence oriented fairness, but an independence requisite ($D5 = I$).

As far as we know, the study of fairness as independence between attributes has been addressed by very few authors. One exception is the *Relative Opportunity* metric proposed by Burke et al. (2018). In this metric, fairness is quantified as the ratio between the relative frequency of gender-protected items in one user group and in the other user group. Thus, if the group to which the user belongs and the gender of the item are statistically independent, then the mean returns the neutral value 1. This metric of

independence-based fairness is limited in that it is not generalizable to more than two groups of users and items ($D3 = 2g$) and it does not capture rank position bias ($D4 = \text{Set}$). Another exception is the *Bias Disparity* proposed by Tsintzou et al. (2019). The authors study the ratio between the frequency of the item category in a group of users versus the overall frequency of the item category. If the category of items and the group of users is statistically independent, the metric returns 1. Furthermore, this bias is compared against the original bias in the users' preferences, thus analyzing the extent to which the system introduces biases against the original data. This metric has the same limitations as the previous one.

3.6. Discussion: to what extent metrics capture diverse scenarios

Looking at Table 1, two observations can be made in relation to the analysis of fairness metrics. The first one is that there are predominant dimension variants in the metrics proposed and studied by the community. For example, effectiveness as a benefit function tends to be user-oriented rather than item group-oriented, and the item-oriented fairness metrics based on effectiveness are group size proportional. That is, none of them take as fairness criterion parity ($D5 = P$) or the utility mass provided by item groups ($D5 = U$). In other words, there exist combinations of dimension values that are not captured by specific metrics in the literature.

These alternative scenarios can be captured by flexible fairness metrics. In this respect, our second observation is that many of the generalists approaches such as Steck (2018), *Equity of Amortized Attention* (Biega et al., 2018), Deldjoo et al. (2021), and Kirnap et al. (2021) apply KLD between benefit function distributions. In particular, the Deldjoo et al. (2021) and Sacharidis et al. (2019) models capture most dimension variants.

However, the outstanding issue is still the independence fairness criterion ($D5 = I$). More specifically, *Relative Opportunity* (Burke et al., 2018) and *Bias Disparity* (Tsintzou et al., 2019) metrics only allow comparing two groups and do not capture graded exposure. Like the information theoretic metric KLD generalizes equity fairness aspects, according to our intuition, the information theoretic metric Mutual Information (MI) is the most appropriate for independence analysis.

In sum, after analyzing the great variety of existing metrics, the hypothesis on which the framework proposed in this paper is based is that, *interpreting the space of users and items as a probabilistic sample space, the two fundamental measures in information theory (KLD and MI) can capture most possible scenarios of fairness measurement on recommendation system outputs.*

4. Theoretical framework

As discussed in the previous sections, we consider fairness as *the equity or independence of user or item groups regarding a certain benefit distribution*. According to our analysis of fairness metrics, the KLD between the benefit distribution across groups and the ideal distribution (Deldjoo et al., 2021; Kirnap et al., 2021) captures most dimensions of existing fairness metrics. In our general framework, we consider this schema for equity fairness. At theoretical level, the main particularities of the proposed framework with respect to previous approaches is that, not only exposure, but also the item exposure effectiveness is modeled as a probability distribution over single user/item pairs. The second contribution is that we also define a metric based on MI to capture independence between user/item groups and individuals regardless any target benefit distribution (independence-based fairness).

4.1. Framework definition

Fig. 2 illustrates the fairness framework and its notation. Let \mathcal{U} and \mathcal{I} be the sets of users and items, respectively. To denote the elements of these sets we use $u \in \mathcal{U}$ and $i \in \mathcal{I}$. Let $\mathcal{A}_{\mathcal{U}}$ and $\mathcal{A}_{\mathcal{I}}$ be the sets of user and item attributes (e.g., $\mathcal{A}_{\mathcal{U}} = \{\text{male, female}\}$). $\psi(u, i)$ and $\phi(u, i)$ represent the *utility* and the *exposure* of the item i for the user u respectively. Both are functions from the user/item space $\mathcal{U} \times \mathcal{I}$ to $(0, 1)$. The exposure function is interpreted as item accessibility, or the probability of the user u to access i . Then, the *Exposure Effectiveness* $\text{Eff}(u, i)$ represents to what extent an item exposure to a user is effective and is modeled as: $\text{Eff}(u, i) = \phi(u, i) \cdot \psi(u, i)$. One way to interpret the above definitions is as follows: $\psi(u, i)$ is a user-defined function, while $\phi(u, i)$ is a system-driven function. While $\psi(u, i)$ answers the question of “how much user u judges item i useful”, $\phi(u, i)$ represents “how much the system provides opportunity to a user-item pair meet”. If one of these two functions, ϕ or ψ , decreases then $\text{Eff}(u, i)$ decreases as well. This effectiveness formalization is similar to the one defined in the model by Deldjoo et al., but instead of operating on an individual user (aggregate) level, i.e., $\text{Eff}(u)$, it measures effectiveness for each user/item pair $\text{Eff}(u, i)$.

We can then normalize the functions ψ , ϕ , and Eff , obtaining three probability distributions P_{θ} with $\theta \in \{\psi, \phi, \text{Eff}\}$ over the user/item space. That is, $P_{\theta}(u, i) = \frac{\theta(u, i)}{\sum_{u \in \mathcal{U}, i \in \mathcal{I}} \theta(u, i)}$. Given a distribution $P_{\theta}(u, i)$, we can infer the probability associated to user attributes (e.g., $P_{\theta}(\text{male}, \mathcal{I})$), item attributes (e.g., $P_{\theta}(\mathcal{U}, \text{action films})$), or combinations of user and item attributes (e.g., $P_{\theta}(\text{male}, \text{action films})$).

Our generalized fairness measurement model is based on two parameterizable metrics. Table 2 illustrates the possibilities. First, the inequity of user or item groups is quantified via KLD between the real (P_{θ}) and fair (Q) benefit distribution. Following the scheme proposed by other authors (Deldjoo et al., 2021; Kirnap et al., 2021; da Silva et al., 2021; Yang & Stoyanovich, 2017):

$$\text{Inequity}(\theta, Q, \mathcal{A}_X) = D_{KL}(P_{\theta} \parallel Q; \mathcal{A}_X) = \sum_{x \in \mathcal{A}_X} P_{\theta}(x) \log \frac{P_{\theta}(x)}{Q(x)}. \quad (1)$$

The groups partition \mathcal{A}_X can be user or item oriented. The benefit function θ can be based on utility (ψ), exposure (ϕ), or effectiveness (Eff). The target distribution Q can be parity (equal benefit, $Q(x) = 1/|\mathcal{A}_X|$), proportional to the user group size ($Q(x) = |\{(u, i) \in \mathcal{U} \times \mathcal{I}\}|/|\mathcal{U} \times \mathcal{I}|$), or proportional to utility ($Q(x) = P_{\psi}(x)$). In other words, being x a certain user or item attribute, in the case

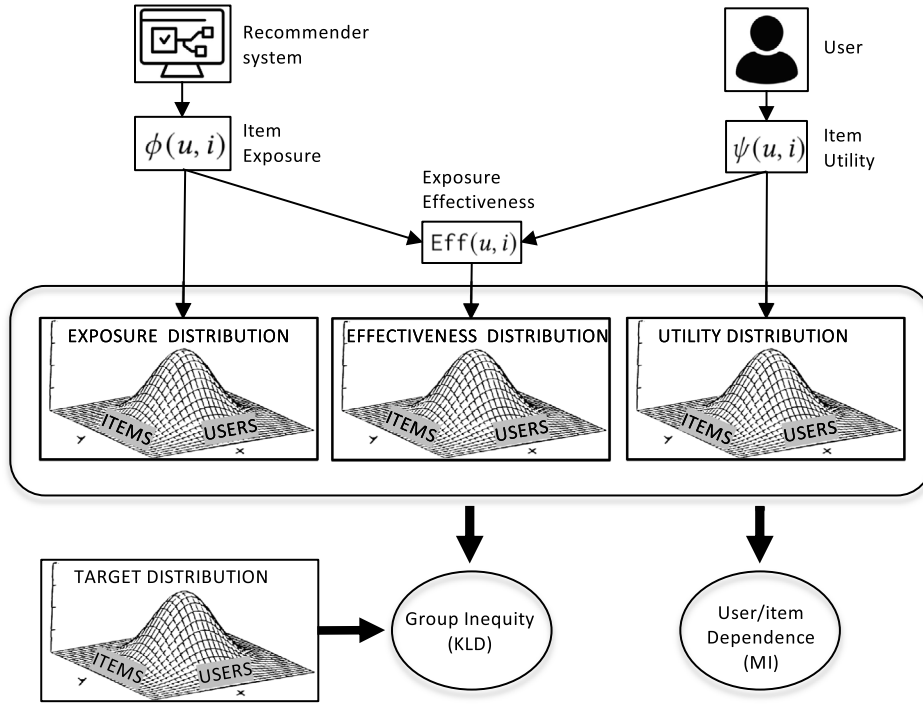


Fig. 2. The main components of the proposed fairness framework; exposure ($\phi(u, i)$), utility ($\psi(u, i)$), effectiveness ($\text{Eff}(u, i)$), their distributions over users and items, the comparison with the target distribution by means of KLD, and the measure of independence by means of MI.

of parity the distribution Q is uniform, in the case of user group size the distribution Q corresponds to the group size, and in the case of proportionality to utility the distribution Q corresponds to the group utility mass.

Second, the treatment disparity is captured via group dependence, which is measured with Mutual Information (MI):

$$\text{Dependence}(\theta, \mathcal{A}_X, \mathcal{A}_Y) = I_\theta(\mathcal{A}_X; \mathcal{A}_Y) = \sum_{x \in \mathcal{A}_X, y \in \mathcal{A}_Y} P_\theta(x, y) \log \frac{P_\theta(x, y)}{P_\theta(x) \cdot P_\theta(y)}. \quad (2)$$

\mathcal{A}_X and \mathcal{A}_Y represent the user and item partitions. When considering single items ($\mathcal{A}_Y = I$) and user groups ($\mathcal{A}_X = \mathcal{A}_{U'}$) we are measuring to what extent the user group does not influence the exposed items. In other words, the user group does not provide information about what items are recommended to the users in that group. In the same way, when considering single users ($\mathcal{A}_X = U$) and item groups ($\mathcal{A}_Y = \mathcal{A}_I$) we are measuring to what extent the item group does not influence to which users the items are exposed. When considering both user and item groups, we are checking that user and item groups do not influence each other.

4.2. Framework generalization power

Our framework provides 18 fairness metric instances. It includes both effectiveness and exposure oriented fairness depending on the benefit function $\theta \in \{\phi, \text{Eff}\}$ (Dimension D1). In addition, our effectiveness function Eff generalizes ranking metrics under the scheme proposed by Carterette (2011) in IR or Singh and Joachims (2018) in RSs, where ϕ represents the ranking decay function such as $1/\log(\text{rank}(u, i))$ in DCG or $p^{\text{rank}(u, i)}$ in RBP. In the set exposure context, Eff generalizes classification metrics such as Accuracy ($\phi(u, i) \in \{0, 1\}$ and $\psi(u, i) \in \{0, 1\}$). It can also generalize Precision or Recall by normalizing the utility with respect to the amount of relevant items in the collection or group, or the amount of exposed items. It also captures both user and producer stakeholders depending on whether we split the user/item space according to user ($\mathcal{A}_X = \mathcal{A}_{U'}$) or item attributes ($\mathcal{A}_X = \mathcal{A}_I$), complying with Dimension D2. The possibilities of Dimension D3 are also captured since we can consider two or more attribute values or even individuals ($\mathcal{A}_{U'} = U$ or $\mathcal{A}_I = I$). The variants of Dimension D4 are also captured: the exposure function ϕ and the effectiveness function Eff can be adapted to ranking or set exposure schemes; the rating scenario requires to state the exposure value for each rating or the translation to a ranking scenario (sorting items by rating). Finally, the variants in Dimension D5 are captured by the flexibility of the target distribution Q (Parity, Size Proportionality, Utility Proportionality) and the application of dependence instead of inequity.

In addition, the framework captures many possibilities that are not covered in the literature. For instance, most of user group oriented inequity metrics in the literature are oriented to effectiveness and based on group-size proportionality: $D_{KL}(P_{\text{Eff}} \parallel P; \mathcal{A}_{U'})$ (D1 = Eff, D2 = Us, and D5 = S). However, in recommendation scenarios with a variable amount of exposed items per user, metrics

Table 2

Fairness metric variants according to the general fairness measurement framework. The most popular fairness notions in the literature are highlighted in bold.

| Inequity $D_{KL}(P_\theta \parallel Q; \mathcal{A}_X)$ (Eq. (1)) | | | | | |
|--|---------------------|--------------------------------|--------------------|-----------------|-----------------------------------|
| Stakeholder | Benefit function | Fairness criterion | \mathcal{A}_X | θ | $Q(x)$ |
| User groups | Exposure Based | Parity | $\mathcal{A}_{U'}$ | ϕ | $1/ \mathcal{A}_{U'} $ |
| | | Size proportionality | | | $ \{(u,i) \in x\} / U' \times I $ |
| | | Utility proportionality | | | $P_\psi(x)$ |
| | Effectiveness Based | Parity | | Eff | $1/ \mathcal{A}_{U'} $ |
| | | Size Proportionality | | | $ \{(u,i) \in x\} / U' \times I $ |
| | | Utility Proportionality | | | $P_\psi(x)$ |
| Item groups | Exposure Based | Parity | \mathcal{A}_I | ϕ | $1/ \mathcal{A}_{U'} $ |
| | | Size proportionality | | | $ \{(u,i) \in x\} / U' \times I $ |
| | | Utility proportionality | | | $P_\psi(x)$ |
| | Effectiveness based | Parity | | Eff | $1/ \mathcal{A}_{U'} $ |
| | | Size proportionality | | | $ \{(u,i) \in x\} / U' \times I $ |
| | | Utility proportionality | | | $P_\psi(x)$ |
| Dependence $I_\theta(\mathcal{A}_X; \mathcal{A}_Y)$ (Eq. (2)) | | | | | |
| User partition | Item partition | Benefit distribution | \mathcal{A}_X | \mathcal{A}_Y | θ |
| User groups | Items | Exposure based | $\mathcal{A}_{U'}$ | I | ϕ |
| | | Effectiveness based | | | Eff |
| User | Item groups | Exposure based | U' | \mathcal{A}_I | ϕ |
| | | Effectiveness based | | | Eff |
| User group | Item groups | Exposure based | $\mathcal{A}_{U'}$ | \mathcal{A}_I | ϕ |
| | | Effectiveness based | | | Eff |

based on exposure (D1 = Exp) could be useful. One could be also interested in distributing the effective exposure mass uniformly across user groups regardless of their size (D5 = P), or proportional to their needs (D5 = U).

Conversely, within the item group inequity metrics, most of exposure based metrics are based on the uniform target distribution ($D_{KL}(P_\phi \parallel 1/|\mathcal{A}_I|; \mathcal{A}_I)$), or the utility-equalized ($D_{KL}(P_\phi \parallel P_\psi; \mathcal{A}_I)$), with the exception of Yang and Stoyanovich's approach which applies the group size proportionality ($D_{KL}(P_\phi \parallel P; \mathcal{A}_I)$) (Yang & Stoyanovich, 2017). We find in the literature two item group equity metrics oriented to effectiveness (Beutel et al., 2019; Chen et al., 2020). Both use the item group size as target distribution. However, one could be interested in giving the same amount of effective exposures to all items groups regardless the amount of items they provide (parity) or in providing effective exposures to item groups with respect to their item utility (utility-equalized).

Regarding dependence based fairness (treatment fairness), the only two metrics that we found in the literature are exposure oriented and combine user and items groups (Burke et al., 2018; Tsintzou et al., 2019) (see Section 3.5). However, the treatment fairness in terms of effective exposures can be the focus in certain scenarios.

In sum, the proposed generalized model captures most fairness notions measured in the literature while opening the door for new fairness variants. In addition, it overcomes many limitations presented by the other metrics. For instance, the item oriented exposure fairness metrics in the literature capture ranking exposures but only for two groups, or capture many groups but only on item set exposures (Section 3.1). Existing independence-based metrics capture only two groups and item set exposure.

It should be noted that there are many aspects of fairness measurement that remain unresolved. Patro et al. (2022) identify some of them as provider utility beyond position based exposure, temporal effects, cross-platform effects, or the use of positioning strategies. A positive aspect of the proposed theoretical framework is that the vast majority of these aspects can be encapsulated within the exposure or utility functions, so that the framework does not lose generality.

5. Experiments

To validate the soundness and generality of our proposal, we perform³ experiments on both synthetic and real data sets with state-of-the-art system outputs. Our research questions are:

- *RQ1. Do the metrics capture those aspects of fairness for which they are designed?* To answer this question, we use synthetic data. We artificially generate a distribution of users, items, and preferences, as well as seven system outputs with known biases to check that each metric captures specific aspects.
- *RQ2. Is there a trade-off between fairness and effectiveness?* Answering this question requires real data set and real systems. There are many studies in the literature that observe a certain trade-off between effectiveness and fairness. But is this true for every fairness criteria? We exploit the generality and completeness of our framework to check this.

³ Source code for running these experiments can be found in the following GitHub repository: [abelogin/FairnessFramework4RecSys](https://github.com/abelogin/FairnessFramework4RecSys).

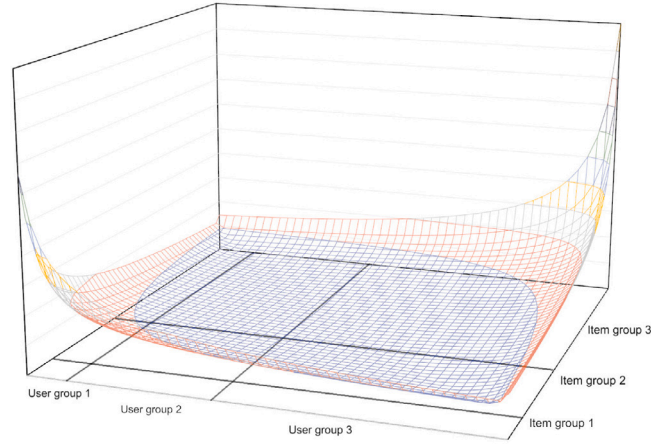


Fig. 3. Utility distribution across user and item groups in the synthetic data set.

Table 3

Name, description, and hypothesized behavior of synthetic baseline systems. $\text{Pri}(u, i)$ represents the priority of item i for user u which determines the item ranking position.

| Baseline | $\text{Pri}(u, i)$ | Description | Hypothesized behavior |
|------------------------------------|--|---|--|
| Oracle | $P_\psi(u, i)$ | Items are sorted according to user utility. | High effectiveness. Item group exposure inequities. User/item dependencies. |
| Random | $\text{Rand}()$ | Items are sorted randomly. | Low effectiveness. Group size and utility proportional equity. User/item independence. |
| Popularity | $P_\psi(i)$ | Items are exposed according to their popularity in ψ . | High effectiveness. Item group exposure inequity. Exposure independence. |
| Randomized Oracle | $P_\psi(i)$ | Adding a random factor to Oracle. | Less effective than Oracle, but more item group exposure proportionality and independence. |
| Item Group Size Norm. Oracle | $\frac{P_\psi(i)}{ g(i) }$ | Items from small groups are prioritized. Item priority is its utility divided by the group size. | Less effective than Oracle, but more item group parity (uniform distribution). |
| Item Group Exposure Cal. Oracle | $P_\psi(u, i) \cdot \frac{P_\psi(g(i))}{ g(i) }$ | The item priority is its utility multiplied by the item group utility density. | Lower effectiveness than the Oracle but higher utility proportional item group equity. |
| Item Group Exposure Cal. Random | $P_\psi(u, i) \cdot \frac{P_\psi(g(i))}{ g(i) }$ | The item priority is a random value multiplied by the item group utility density. | Higher effectiveness than Random and higher utility proportional item group equity. |
| Group Debaised Oracle | $\frac{P_\psi(u, i)}{P_\psi(g(u), g(i))}$ | It reduces the user/item group bias by dividing the Oracle utility by the user and item group utility mass. | Lower effectiveness than the Oracle but higher exposure and effectiveness independence between user and item groups. |
| Item Group Single User Deb. Oracle | $\frac{P_\psi(u, i)}{P_\psi(u, g(i))}$ | It reduces the item group vs single user bias by dividing the Oracle utility by the item group utility mass. | Lower effectiveness than the Oracle but higher exposure and effectiveness independence between single users and item groups. |
| User Group Single Item Deb. Oracle | $\frac{P_\psi(u, i)}{P_\psi(g(u), I)}$ | It reduces the user group bias across single items by dividing the Oracle utility by the user group utility mass. | Lower effectiveness than the Oracle but higher exposure and effectiveness independence between user groups and single items. |

- **RQ3. Is there a trade-off between fairness metrics?** There is work showing that there are incompatibilities between fairness metrics. That is, different fairness criteria cannot be satisfied simultaneously. Regardless of the fact that not all metrics can be maximized simultaneously, we will study on real data and systems whether improving one fairness criterion necessarily implies worsening others.
- **RQ4. Are the fairness metrics consistent across data sets?** Our framework provides 18 fairness metric instances. We hypothesize that some fairness criteria are more sensitive to particular data sets than others. We run them on two different real data sets over the same systems to check this.

5.1. Synthetic recommendation outputs

In the following experiment, we apply the fairness metrics derived from our framework to synthetic data and RS outputs. The aim is to answer RQ1. As a starting point we generate an oracle system output, in which the items are sorted by utility for each user, and

Table 4

Fairness metrics for synthetic outputs. DCG based exposure. We use the same notation as before: P, S, and U are the fairness criteria Parity, Size Proportionality, and Utility Proportionality, respectively; ϕ and Eff are Exposure and Effectiveness. Colored background denotes values that match hypotheses presented in Table 3, whereas italics are used to highlight best values for each metric (column), where except for Eff, for which higher is preferable, the lowest value indicates the best performance.

| User/Item groups: | Eff | Inequity (KLD) | | | | | | | | | Dependency (MI) | | | | | |
|--|-------|----------------|-------|-------|-------------|-------|--------|-------|--------|-------|-----------------|-------|--------------|-------|-------------|-------|
| | | User groups | | | Item groups | | | | | | User groups | | Item groups | | User groups | |
| | | A_U | | | A_I | | | | | | A_U-A_I | | A_I-U | | $I-A_U$ | |
| | | P | S | U | P | S | U | P | S | U | | | Independence | | | |
| Fairness Criterion: | | P | S | U | P | S | U | P | S | U | ϕ | Eff | ϕ | Eff | ϕ | Eff |
| Benefit Function: | DCG | Eff | Eff | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff |
| Oracle | 0.292 | 0.172 | 0.048 | 0.000 | 0.198 | 0.152 | 0.026 | 0.217 | 0.002 | 0.058 | 0.014 | 0.152 | 0.019 | 0.197 | 0.020 | 0.186 |
| Random | 0.221 | 0.171 | 0.051 | 0.000 | 0.283 | 0.166 | 0.000 | 0.053 | 0.036 | 0.000 | 0.000 | 0.050 | 0.004 | 0.072 | 0.004 | 0.071 |
| Popularity | 0.275 | 0.194 | 0.036 | 0.000 | 0.202 | 0.173 | 0.026 | 0.178 | 0.002 | 0.042 | 0.000 | 0.062 | 0.000 | 0.082 | 0.000 | 0.076 |
| Randomized Oracle | 0.273 | 0.178 | 0.044 | 0.000 | 0.215 | 0.151 | 0.016 | 0.169 | 0.006 | 0.036 | 0.008 | 0.119 | 0.012 | 0.163 | 0.012 | 0.152 |
| Item Group Size Normalized Oracle | 0.273 | 0.137 | 0.060 | 0.002 | 0.088 | 0.066 | 0.085 | 0.370 | 0.016 | 0.140 | 0.001 | 0.072 | 0.001 | 0.108 | 0.003 | 0.097 |
| Item Group Exposure Calibrated Oracle | 0.288 | 0.166 | 0.050 | 0.000 | 0.153 | 0.124 | 0.055 | 0.289 | 0.000 | 0.093 | 0.007 | 0.116 | 0.011 | 0.172 | 0.011 | 0.148 |
| Item Group Exposure Calibrated Random | 0.235 | 0.143 | 0.057 | 0.001 | 0.162 | 0.094 | 0.047 | 0.244 | 0.000 | 0.069 | 0.000 | 0.048 | 0.001 | 0.069 | 0.001 | 0.061 |
| Group Debaised Oracle | 0.279 | 0.151 | 0.055 | 0.000 | 0.095 | 0.074 | 0.074 | 0.330 | 0.014 | 0.117 | 0.001 | 0.087 | 0.004 | 0.139 | 0.005 | 0.119 |
| Item Group Single User Debaised Oracle | 0.265 | 0.121 | 0.066 | 0.005 | 0.078 | 0.052 | 0.090 | 0.399 | 0.021 | 0.160 | 0.000 | 0.055 | 0.000 | 0.078 | 0.001 | 0.073 |
| User Group Single Item Debaised Oracle | 0.281 | 0.188 | 0.039 | 0.000 | 0.186 | 0.156 | 0.034 | 0.221 | 0.000 | 0.060 | 0.000 | 0.064 | 0.009 | 0.134 | 0.000 | 0.076 |

a random system output, in which the items are sorted randomly for each user. The methodology consists of modifying artificially the oracle output or the random baseline to improve particular fairness features. Then the metric results should be consistent.

5.1.1. Data and settings.

Our synthetic data consists of 100 users and 100 items, both divided into three groups (1–10, 11–30, 31–100). The utility function is: $\psi(u, i) = \text{Max} (1/\sqrt{i-u}, 1/\sqrt{(101-i)-(101-u)})$. Fig. 3 illustrates the user/item utility distribution across groups. The resulting distribution is such that items 1 and 100 are more popular than the rest; groups are unbalanced (10, 20, and 70 items or users); user group A is biased toward item group I and user group C is biased toward item group III.

Table 3 displays the name, description, and hypothesized behavior of synthetic baseline systems. Each synthetic system output consists of a ranking of 100 items per user, ordered according to a certain priority function $\text{Pri}(u, i)$. The Oracle system output sorts items by utility ($\text{Pri}(u, i) = P_\psi(u, i)$), while the Random baseline sort items randomly. The rest of systems modify these baselines by multiplying them with a certain *fairness factor*.

5.1.2. Results

We consider the DCG decay exposure function in both the effectiveness and fairness measurements. That is, being $\text{Rank}(u, i)$ the ranking position of the item i in the user u interface according to $\text{Pri}(u, i)$, then $\phi(u, i) = 1/(\log(\text{Rank}(u, i))+1)$. Table 4 shows the fairness measurement results in all metric variants presented in Table 2. We do not consider the Exposure benefit function in user groups, given that all users receive the same amount of information in our experiment.

Numbers with colored background indicate those values that corroborate the hypotheses described in Table 3 for each of the synthetic system outputs. Oracle maximizes effectiveness at the cost of item group exposure inequities ($\text{KLD}-A_I-P-\phi = 0.198$ and $\text{KLD}-A_I-S-\phi = 0.198$) and stresses the user/item dependencies (MI fairness metrics). On the contrary, Random achieves the lowest effectiveness ($\text{DCG}=0.221$), but provides group size and utility proportional equity ($\text{KLD}-A_I-S-\phi=\text{KLD}-A_I-U-\text{Eff} = 0$) and user/item independence. Popularity provides high effectiveness and exposure user/item independence ($\text{MI}-A_U-A_I-\phi=\text{MI}-A_I-U-\phi=\text{MI}-I-A_U-\phi = 0$), since all users receive the same recommendation. The cost is a higher item group exposure inequity ($\text{KLD}-A_I-P-\phi$, $\text{KLD}-A_I-P-\text{Eff}$, $\text{KLD}-A_I-S-\phi$ and $\text{KLD}-A_I-S-\text{Eff}$). The unfairness effect of Oracle can be smoothed by a randomization fairness factor (Randomized Oracle), but at the cost of a decreased efficiency ($\text{DCG}=0.275$). We can also favor the uniform distribution of exposure across groups ($\text{KLD}-A_I-P-\phi$ and $\text{KLD}-A_I-P-\text{Eff}$) by dividing the utility of the items by the size of their group (Item Group Size Normalized Oracle). If we add the item group utility density as fairness factor (Item Group Exposure Calibrated Oracle and Item Group Exposure Calibrated Random) then we can improve the item group utility proportional equity ($\text{KLD}-A_I-U-\phi = 0$). Finally, we can improve the exposure and effectiveness independence between user and item groups or individuals (MI based metrics) by adding as fairness factor the utility mass of user and/or item groups (Group Debaised Oracle,

Item Group Single User Debiased Oracle, and User Group Single Item Debiased Oracle). In addition, we repeated the experiment but exposing 10 items per user (flat exposure). That is $\phi(u, i)$ is 1 if $\text{Rank}(\text{Pri}(u, i)) \leq 10$ and 0 otherwise. We obtained similar results.

In conclusion, according to our experiments with synthetic data, the answer to RQ1 is positive: the metrics instantiated from the general framework capture different system output features and they are consistent with the hypothesized behavior of synthetic system outputs.

5.2. Behavior of state-of-the-art RSs

In this second experiment, we analyze the behavior of the proposed framework on a real recommendation data set and system outputs in order to answer RQ2, RQ3, and RQ4.

5.2.1. Data sets

This study evaluates the performance of Collaborative Filtering (CF) approaches within the presented fairness evaluation framework using two popular data sets including explicit or implicit preferences:

- **Netflix** (explicit). The original version of this data set is one of the largest available benchmark data sets used widely for CF algorithms today (Bennett & Lanning, 2007). It has ratings collected over the course of seven years. We used a “small” variant of this data set with 9992 users, 4945 items, 607,803 ratings.
- **CiteULike-a** (implicit). The CiteULike data set⁴ is about academic citations. CiteULike is an online platform that enables registered users to establish personal libraries by archiving relevant articles. The data set consists of the papers in the users’ libraries (which are handled as “likes”), the tags provided by the users, as well as the title and abstract of the papers. CiteULike-a (Wang & Blei, 2011) data set contains 4122 users, 16,908 items, and 155,588 interactions.

5.2.2. Systems

We investigated a variety of latent factors CF models, which have been employed in previous and ongoing works of RS research to achieve excellent performance in rating and ranking tasks (Cremonesi, Koren, & Turrin, 2010; Deldjoo, Bellogín, & Di Noia, 2021; Koren, Bell, & Volinsky, 2009; Naghiaei, Rahmani, & Deldjoo, 2022).

- **MF** (Koren et al., 2009): A classical Matrix Factorization (MF) approach; in this case, the user and item factor are learned through Stochastic Gradient Descent, despite the availability of other techniques (Hu, Koren, & Volinsky, 2008). The predicted rating in MF is computed as $\hat{r}_{ui} = \mathbf{q}_i^T \mathbf{p}_u$, where $\mathbf{p}_u \in \mathbb{R}^H$ and $\mathbf{q}_i \in \mathbb{R}^H$ are the learned H -sized latent vectors for the user u and item i , respectively.
- **PMF** (Salakhutdinov & Mnih, 2007): A Maximum A Posteriori approach is used to factorize the matrix in light of a probabilistic linear model containing Gaussian noise.
- **BPR-MF** (Koren et al., 2009; Rendle, Freudenthaler, Gantner, & Schmidt-Thieme, 2009): BPR is the state-of-the-art method for personalized ranking, especially on data sets containing implicit feedback. MF is used as the predictor in BPR-MF. It is important to note that this algorithm tends to recommend popular items more often than other methods (Anelli, Bellogín, Noia, Jannach, & Pomo, 2022).
- **WMF** (Hu et al., 2008; Pan et al., 2008): Classic weighted MF model for implicit feedback data. It assumes the independence of the latent features of two items and gives lower weights to negative samples. The equivalent ALS-based approach (Hu et al., 2008) can reduce inference complexity.
- **NeuMF** (He et al., 2017): Using multi-layer perceptron and MF, this approach learns users and item features, and then uses non-linear activation functions to train a mapping between these features.
- **VAECF** (Liang, Krishnan, Hoffman, & Jebara, 2018): The method relies on variational autoencoders, which present a multinomial likelihood generative model and employ Bayesian inference for parameter estimation.

We consider the same baseline approaches as in the previous experiment (Oracle, Random, and Popularity). Note that Random has some dependency between user groups and items in effectiveness due to the original bias of the data. It also has dependencies between groups of items and individual users due to the random effect. That is, not all individual users have a uniform distribution of item groups and vice versa. The effectiveness and fairness metrics are exactly the same as in the previous experiment (Section 5.1).

5.2.3. Evaluation setup.

For each considered recommendation model, we ran them at their default hyper-parameter values according to their implementation in the Cornac recommender framework (Salah, Truong, & Lauw, 2020). The results of the recommendation were generated based on a hold-out setting (80%–20% training-test split).

⁴ <http://www.citeulike.org/>

Table 5

Fairness metrics for CiteULike data set. DCG based exposure. Same notation as in Table 4.

| | | Inequity (KLD) | | | | | | | | | Dependency (MI) | | | | | |
|---------------------|-------|----------------|-------|-------|-------------|-------|--------|-------|--------|-------|-----------------|--------------|--------------|-------|--------------|-------|
| User/Item groups: | Eff | User groups | | | Item groups | | | | | | User groups | | Item groups | | User groups | |
| | | A_U | | | A_I | | | | | | Item groups | | Single users | | Single items | |
| | | | | | | | | | | | A_I-A_U | | $I-A_U$ | | | |
| Fairness Criterion: | | P | S | U | P | S | | U | | | | Independence | | | | |
| Benefit Function: | DCG | Eff | Eff | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff |
| Oracle | 3.362 | 0.120 | 0.127 | 0.039 | 0.275 | 0.275 | 0.141 | 0.141 | 0.000 | 0.000 | 0.002 | 0.002 | 0.221 | 0.221 | 0.406 | 0.406 |
| Random | 0.002 | 0.210 | 0.219 | 0.005 | 0.023 | 0.357 | 0.000 | 0.204 | 0.171 | 0.005 | 0.000 | 0.131 | 0.091 | 0.642 | 0.305 | 0.789 |
| Popularity | 0.038 | 0.227 | 0.237 | 0.003 | 1.000 | 1.000 | 0.760 | 0.760 | 0.316 | 0.316 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.032 |
| MF | 0.002 | 0.118 | 0.125 | 0.040 | 0.001 | 0.409 | 0.036 | 0.244 | 0.376 | 0.014 | 0.000 | 0.000 | 0.091 | 0.590 | 0.042 | 0.881 |
| PMF | 0.030 | 0.064 | 0.069 | 0.089 | 0.975 | 0.977 | 0.736 | 0.739 | 0.296 | 0.298 | 0.000 | 0.003 | 0.019 | 0.022 | 0.031 | 0.423 |
| BPR-MF | 0.038 | 0.227 | 0.237 | 0.003 | 1.000 | 1.000 | 0.760 | 0.760 | 0.316 | 0.316 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.032 |
| WMF | 0.121 | 0.186 | 0.195 | 0.011 | 0.501 | 0.707 | 0.319 | 0.494 | 0.039 | 0.127 | 0.009 | 0.014 | 0.152 | 0.237 | 0.246 | 0.646 |
| NeuMF | 0.166 | 0.108 | 0.115 | 0.047 | 0.786 | 0.884 | 0.564 | 0.653 | 0.171 | 0.232 | 0.002 | 0.002 | 0.125 | 0.108 | 0.115 | 0.370 |
| VAECF | 0.167 | 0.102 | 0.109 | 0.051 | 0.864 | 0.916 | 0.634 | 0.682 | 0.219 | 0.254 | 0.000 | 0.006 | 0.087 | 0.078 | 0.067 | 0.340 |

Table 6

Fairness metrics for Netflix data set. DCG based exposure. Same notation as in Table 4.

| User/Item groups: | | Inequity (KLD) | | | | | | | | | Dependency (MI) | | | | | |
|---------------------|--------|----------------|-------|-------|-------------|-----------|---------|-------|---------|-------|----------------------------|-------|-----------------------------|-------|-----------------------------|-------|
| | | User groups | | | Item groups | | | | | | User groups Item groups | | Item groups Single users | | User groups Single items | |
| | | | | | | | | | | | | | | | | |
| | | A_U | A_I | | | A_I-A_U | A_I-U | | $I-A_U$ | | | | | | | |
| Fairness Criterion: | | P | S | U | P | S | U | | | | Independence | | | | | |
| Benefit Function: | DCG | Eff | Eff | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff | ϕ | Eff |
| Oracle | 17.212 | 0.082 | 0.306 | 0.062 | 0.328 | 0.345 | 1.816 | 1.853 | 0.000 | 0.001 | 0.002 | 0.002 | 0.174 | 0.181 | 0.101 | 0.101 |
| Random | 0.038 | 0.307 | 0.666 | 0.001 | 0.450 | 0.313 | 0.000 | 1.782 | 1.593 | 0.000 | 0.000 | 0.019 | 0.094 | 0.642 | 0.051 | 0.584 |
| Popularity | 1.296 | 0.026 | 0.185 | 0.145 | 1.000 | 1.000 | 2.977 | 2.977 | 0.300 | 0.300 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.017 |
| MF | 0.366 | 0.344 | 0.718 | 0.006 | 0.108 | 0.757 | 0.162 | 2.623 | 0.867 | 0.142 | 0.020 | 0.000 | 0.161 | 0.221 | 0.067 | 0.066 |
| PMF | 0.361 | 0.172 | 0.464 | 0.011 | 0.012 | 0.564 | 0.414 | 2.291 | 0.508 | 0.054 | 0.011 | 0.001 | 0.222 | 0.402 | 0.046 | 0.135 |
| BPR-MF | 1.320 | 0.043 | 0.226 | 0.110 | 1.000 | 1.000 | 2.977 | 2.977 | 0.300 | 0.300 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.020 |
| WMF | 1.043 | 0.034 | 0.206 | 0.126 | 0.497 | 0.818 | 2.167 | 2.718 | 0.032 | 0.176 | 0.000 | 0.002 | 0.214 | 0.161 | 0.150 | 0.274 |
| NeuMF | 1.454 | 0.000 | 0.085 | 0.275 | 0.868 | 0.933 | 2.794 | 2.888 | 0.207 | 0.250 | 0.000 | 0.000 | 0.094 | 0.062 | 0.092 | 0.156 |
| VAECF | 1.605 | 0.012 | 0.143 | 0.189 | 0.914 | 0.971 | 2.861 | 2.941 | 0.237 | 0.278 | 0.000 | 0.000 | 0.071 | 0.027 | 0.115 | 0.171 |

5.2.4. Results

The results are shown in Tables 5 (for CiteULike) and 6 (for Netflix), commented in the following with particular emphasis on the values highlighted in color. The answers to our three research questions RQ2–RQ4 can be synthesized as follows.

- **RQ2. Is there a trade-off between fairness and effectiveness?** The answer for this question is that it depends on the fairness metric. For instance, the highest DCG (3.362 and 17.212 respectively in each data set) is achieved by Oracle with a perfect fairness (zero KLD) in terms of item group utility-proportional exposure equity (KLD- $A_I-U-\phi$). In both data sets, the neural based systems (VAECF and NeuMF) achieve higher DCG values (0.166, 0.167, 1.454, and 1.605) than MF-based systems and also are more fair in terms of KLD- $A_I-U-\phi$ (0.047, 0.051, 0.275, and 0.189). On the contrary, it seems that there exists a trade-off between effectiveness and size-proportional item group exposure inequity (KLD- $A_I-S-\phi$). As an intriguing observation, we could connect this practical and general result to the well-known **accuracy-diversity** or **accuracy-novelty** trade-off phenomenon in the community (Vargas & Castells, 2011), and now we could observe a similar trade-off of effectiveness-item fairness. The Random system minimizes this inequity in both data sets (zero KLD) at the cost of DCG (0.002 and 0.038). On the other hand, Oracle maximizes effectiveness by increasing the inequity (0.141 and 1.816). The neural based systems (NeuMF and VAECF) are more effective than the others, but highly unfair in terms of item group size-proportional exposure equity (KLD- $A_I-S-\phi$) in both data sets, (0.564, 0.634, 2.794, and 2.861).
- **RQ3. Is there a trade-off between fairness metrics?** In view of the results, we cannot state that there is a trade-off between fairness metrics. However, we see that different metrics express different characteristics of the systems. For example, regarding the dependence-based fairness metrics (MI), all systems keep the independence between user and item groups (MI- $A_I-A_U-\phi$ and MI- $A_I-A_U-\text{Eff}$ are zero or almost zero for all systems). However, WMF seems to state a certain dependence between single users and item groups and between single items and user groups (MI- $A_I-U-\phi$, MI- $A_I-U-\text{Eff}$, MI- $I-A_U-\phi$, and MI- $I-A_U-\text{Eff}$); this suggests a higher personalizing degree in the recommendation. On the other hand, although BPR-MF presents item group inequities (KLD- $A_I-S-\phi$ and KLD- $A_I-S-\text{Eff}$ are 0.760 and 2.977), it keeps the independence between user/item individuals and groups in both data sets (MI- $A_I-U-\phi$, MI- $A_I-U-\text{Eff}$, MI- $I-A_U-\phi$, and MI- $I-A_U-\text{Eff}$ are all close to zero).
- **RQ4. Are the fairness metrics consistent across data sets?** Not every metric is consistent across data sets. For instance, NeuMF is more fair than VAECF in terms of size-proportional user group effectiveness (KLD- $A_U-S-\text{Eff}$) in the CiteULike data set, but

not in the Netflix data set. In addition, for this user oriented metric, the Popularity baseline is unfair in CiteULike but not in Netflix. This suggests that user group effectiveness fairness is sensitive to the evaluation benchmark. We hypothesize that the nature of systems is more determinant in item group fairness than in user group fairness which is highly related with the distribution of user preferences.

In summary, these results confirm that the different instantiations of fairness metrics in real data sets give us different information about system output bias and, in some cases, this information is sensitive to the particularities of the data set.

6. Conclusions and future work

Contributions. In this paper we have defined a formal, broad, and unified framework for measuring fairness in RSs, and validated it experimentally. The proposed framework captures the five dimensions that characterize existing fairness metrics in the literature. The practical implications of this model are essentially: (i) a tool to identify the most appropriate metric in a given scenario, (ii) the unification of fairness evaluation criteria for the comparison of results in different research works, and (iii) the identification of formal aspects of fairness that have not yet been explored, such as the statistical independence of the benefit between user and item attributes. We hope that these contributions will allow a better understanding of fairness measurement and, in perspective, to overcome the limitations imposed by the current fragmented landscape of fairness definitions and metrics.

In general, we expect both researchers and practitioners to benefit from these contributions, especially those concerned about measuring and assessing fairness from novel dimensions. This is because our framework, as defined and demonstrated throughout the paper, is two-sided (it allows capturing the notion of fairness on users and items at the same time without the need of having an *ideal* preconceived notion of fairness), flexible (because it is possible to boil down to many existing notions of fairness), and reliable (as it is focused on *independence* rather than *equity*).

Limitations and Outlook. While we make no claim that the proposed five dimensions for fairness measurement are exhaustive (as we anticipated at the end of Section 2), we believe they can serve as a useful start point for practitioners, students, and scholars. Nonetheless, we briefly outline several other dimensions that could be taken into account. For example, different scales can be defined for the utility of items (binary, rating, preferences, continuous, etc.), and the benefit distribution can be defined in terms of groups or the user past behavior itself (the notion of *calibration* Steck, 2018). However, from our point of view both aspects can be encapsulated within the notion of item utility, to which the fairness model should be agnostic. Another dimension that could be taken into account is the possibility of considering degrees of membership of items or users in groups, with a non-binary group membership function. We have not included this dimension as it is very rare in the literature, although we do take it into account in the definition of our theoretical framework.

It should be noted that some notions of fairness are not captured by our five dimensions. For example, *Fairness Through Unawareness* (Chen et al., 2020; Kusner et al., 2017; Verma & Rubin, 2018) represents to what extent certain attributes are not explicitly used in the training process. However, in this paper we focus on the evaluation of recommender output, regardless of how the system has been trained. In addition, our five dimensions focus on *group fairness* rather than *individual fairness* (Pitoura, Stefanidis, & Koutrika, 2022). However, due to definition of group fairness that incorporate as input protected features, more attention has been paid to group fairness; also, individual fairness requires a certain (arbitrary) similarity function between users or items (Chen et al., 2020; Ekstrand et al., 2018).

Although we believe that our experimental results are representative, in the future we aim to perform a more complete experimental activity, with more RSs and on more data sets. In addition, we remark that our approach seems general to be applied to any kind of information system including IR systems; we plan to do so in future work.

Still about future work, this framework offers us a uniform tool to comprehensively study the theoretical and empirical trade-off between different fairness criteria. Although there is work in the literature in this respect, the lack of a general framework for measuring fairness has not yet allowed a comprehensive analysis of the problem. Being even more ambitious, we intend to exploit this theoretical tool to identify a single measure that, even at the cost of effectiveness, ensures maximum fairness levels in all metrics. One candidate measure could be the multi-variate entropy, but this conjecture requires further study.

CRedit authorship contribution statement

Enrique Amigó: Conceptualization, Methodology, Software, Validation, Formal analysis, Visualization, Writing – original draft, Writing – review & editing. **Yashar Deldjoo:** Conceptualization, Methodology, Software, Validation, Formal analysis, Visualization, Writing – original draft, Writing – review & editing. **Stefano Mizzaro:** Conceptualization, Validation, Visualization, Writing – original draft, Writing – review & editing. **Alejandro Bellogín:** Conceptualization, Validation, Visualization, Writing – original draft, Writing – review & editing.

Data availability

Data will be made available on request.

Acknowledgments

This work was supported in part by the Ministerio de Ciencia, Innovación y Universidades, Spain (references PID2019-108965GB-I00 and PID2021-124361OB-C32). The authors thank the reviewers for their thoughtful comments and suggestions.

References

- Abdollahpour, Himan, Adomavicius, Gediminas, Burke, Robin, Guy, Ido, Jannach, Dietmar, Kamishima, Toshihiro, et al. (2020). Multistakeholder recommendation: Survey and research directions. *User Modeling and User-Adapted Interaction*, 30(1), 127–158. <http://dx.doi.org/10.1007/s11257-019-09256-1>.
- Abdollahpour, Himan, & Mansoury, Masoud (2020). Multi-sided exposure bias in recommendation. In *International workshop on industrial recommendation systems in conjunction with ACM KDD 2020*. arXiv:2006.15772.
- Abdollahpour, Himan, Mansoury, Masoud, Burke, Robin, Mobasher, Bamshad, & Malthouse, Edward C. (2021). User-centered evaluation of popularity bias in recommender systems. In Judith Masthoff, Eelco Herder, Nava Tintarev, & Marko Tkalcić (Eds.), *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization* (pp. 119–129). ACM, <http://dx.doi.org/10.1145/3450613.3456821>.
- Abel, Fabian, Deldjoo, Yashar, Elahi, Mehdi, & Kohlsdorf, Daniel (2017). RecSys challenge 2017: Offline and online evaluation. In Paolo Cremonesi, Francesco Ricci, Shlomo Berkovsky, & Alexander Tuzhilin (Eds.), *Proceedings of the eleventh ACM conference on recommender systems* (pp. 372–373). ACM, <http://dx.doi.org/10.1145/3109859.3109954>.
- Anelli, Vito Walter, Bellogín, Alejandro, Noia, Tommaso Di, Jannach, Dietmar, & Pomo, Claudio (2022). Top-N recommendation algorithms: A quest for the state-of-the-art. In Alejandro Bellogín, Ludovico Boratto, Olga C. Santos, Liliana Ardisson, & Bart Knijnenburg (Eds.), *UMAP '22: 30th ACM conference on user modeling, adaptation and personalization* (pp. 121–131). ACM, <http://dx.doi.org/10.1145/3503252.3531292>.
- Bennett, James, & Lanning, Stan (2007). The Netflix prize. In *Proceedings of KDD cup and workshop: Vol. 2007*, (p. 35). Citeseer.
- Berk, Richard A., Heidari, Hoda, Jabbari, Shahin, Kearns, Michael, & Roth, Aaron (2017). Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 50, 3–44.
- Beutel, Alex, Chen, Jilin, Doshi, Tulsee, Qian, Hai, Wei, Li, Wu, Yi, et al. (2019). Fairness in recommendation ranking through pairwise comparisons. In Ankur Teredesai, Vipin Kumar, Ying Li, Rómer Rosales, Evimaria Terzi, George Karypis (Eds.), *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 2212–2220). ACM, <http://dx.doi.org/10.1145/3292500.3330745>.
- Biega, Asia J., Gummadi, Krishna P., & Weikum, Gerhard (2018). Equity of attention: Amortizing individual fairness in rankings. In Kevyn Collins-Thompson, Qiaozhu Mei, Brian D. Davison, Yiqun Liu, & Emine Yilmaz (Eds.), *The 41st international ACM SIGIR conference on research & development in information retrieval* (pp. 405–414). ACM, <http://dx.doi.org/10.1145/3209978.3210063>.
- Boratto, Ludovico, Fenu, Gianni, & Marras, Mirko (2021). Interplay between upsampling and regularization for provider fairness in recommender systems. *User Modeling and User-Adapted Interaction*, 31(3), 421–455. <http://dx.doi.org/10.1007/s11257-021-09294-8>.
- Burke, Robin (2017). Multisided fairness for recommendation. In *2017 workshop on fairness, accountability, and transparency in machine learning*.
- Burke, Robin, Sonboli, Nasim, & Ordonez-Gauger, Aldo (2018). Balanced neighborhoods for multi-sided fairness in recommendation. In Sorelle A. Friedler, & Christo Wilson (Eds.), *Proceedings of machine learning research: Vol. 81, Conference on fairness, accountability and transparency* (pp. 202–214). PMLR, <http://proceedings.mlr.press/v81/burke18a.html>.
- Carterette, Ben (2011). System effectiveness, user models, and user utility: A conceptual framework for investigation. In Wei-Ying Ma, Jian-Yun Nie, Ricardo Baeza-Yates, Tat-Seng Chua, & W. Bruce Croft (Eds.), *Proceeding of the 34th international ACM SIGIR conference on research and development in information retrieval* (pp. 903–912). ACM, <http://dx.doi.org/10.1145/2009916.2010037>.
- Celis, L., Elisa, Straszak, Damian, & Vishnoi, Nisheeth K. (2018). Ranking with fairness constraints. In Ioannis Chatzigiannakis, Christos Kaklamani, Dániel Marx, & Donald Sannella (Eds.), *LIPIcs: Vol. 107, 45th international colloquium on automata, languages, and programming* (pp. 28:1–28:15). Schloss Dagstuhl - Leibniz-Zentrum für Informatik, <http://dx.doi.org/10.4230/LIPIcs.ICALP.2018.28>.
- Chen, Jiawei, Dong, Hande, Wang, Xiang, Feng, Fuli, Wang, Meng, & He, Xiangnan (2020). Bias and debias in recommender system: A survey and future directions. CoRR abs/2010.03240. arXiv:2010.03240.
- Chouldechova, Alexandra (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big Data*, 5(2), 153–163. <http://dx.doi.org/10.1089/big.2016.0047>.
- Cremonesi, Paolo, Koren, Yehuda, & Turrin, Roberto (2010). Performance of recommender algorithms on top-n recommendation tasks. In Xavier Amatriain, Marc Torrents, Paul Resnick, & Markus Zanker (Eds.), *Proceedings of the 2010 ACM conference on recommender systems* (pp. 39–46). ACM, <http://dx.doi.org/10.1145/1864708.1864721>.
- Deldjoo, Yashar, Anelli, Vito Walter, Zamani, Hamed, Bellogín, Alejandro, & Noia, Tommaso Di (2021). A flexible framework for evaluating user and item fairness in recommender systems. *User Modeling and User-Adapted Interaction*, 31(3), 457–511. <http://dx.doi.org/10.1007/s11257-020-09285-1>.
- Deldjoo, Yashar, Bellogín, Alejandro, & Noia, Tommaso Di (2021). Explaining recommender systems fairness and accuracy through the lens of data characteristics. *Information Processing & Management*, 58(5), Article 102662. <http://dx.doi.org/10.1016/j.ipm.2021.102662>.
- Diaz, Fernando, Mitra, Bhaskar, Ekstrand, Michael D., Biega, Asia J., & Carterette, Ben (2020). Evaluating stochastic rankings with expected exposure. In Mathieu d'Aquin, Stefan Dietze, Claudia Hauff, Edward Curry, & Philippe Cudré-Mauroux (Eds.), *CIKM '20: The 29th ACM international conference on information and knowledge management* (pp. 275–284). ACM, <http://dx.doi.org/10.1145/3340531.3411962>.
- Dong, Qiang, Xie, Shuang-Shuang, & Li, Wen-Jun (2021). User-item matching for recommendation fairness. *IEEE Access*, 9, 130389–130398. <http://dx.doi.org/10.1109/ACCESS.2021.3113975>.
- Dwork, Cynthia, Hardt, Moritz, Pitassi, Toniann, Reingold, Omer, & Zemel, Richard S. (2012). Fairness through awareness. In Shafi Goldwasser (Ed.), *Innovations in theoretical computer science 2012* (pp. 214–226). ACM, <http://dx.doi.org/10.1145/2090236.2090255>.
- Ekstrand, Michael D., Das, Anubrata, Burke, Robin, & Diaz, Fernando (2022). Fairness in information access systems. *Foundations and Trends in Information Retrieval*, 16(1–2), 1–177. <http://dx.doi.org/10.1561/15000000079>.
- Ekstrand, Michael D., & Mahant, Vaibhav (2017). Sturgeon and the cool kids: Problems with random decoys for Top-N recommender evaluation. In Vasile Rus, & Zdravko Markov (Eds.), *Proceedings of the thirtieth international Florida artificial intelligence research society conference* (pp. 639–644). AAAI Press, <https://aaai.org/ocs/index.php/FLAIRS/FLAIRS17/paper/view/15534>.
- Ekstrand, Michael D., Tian, Mucun, Azpiazu, Ion Madrazo, Ekstrand, Jennifer D., Anuyah, Oghenemaro, McNeill, David, et al. (2018). All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. In Sorelle A. Friedler, Christo Wilson (Eds.), *Proceedings of machine learning research: Vol. 81, Conference on fairness, accountability and transparency* (pp. 172–186). PMLR, <http://proceedings.mlr.press/v81/ekstrand18b.html>.
- Ge, Yingqiang, Liu, Shuchang, Gao, Ruoyuan, Xian, Yikun, Li, Yunqi, Zhao, Xiangyu, et al. (2021). Towards long-term fairness in recommendation. In Liane Lewin-Eytan, David Carmel, Elad Yom-Tov, Eugene Agichtein, Evgeniy Gabrilovich (Eds.), *WSDM '21, the fourteenth ACM international conference on web search and data mining* (pp. 445–453). ACM, <http://dx.doi.org/10.1145/3437963.3441824>.
- Ge, Yingqiang, Zhao, Xiaoting, Yu, Lucia, Paul, Saurabh, Hu, Diane J., Hsieh, Chu-Cheng, et al. (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*.
- Geyik, Sahin Cem, Ambler, Stuart, & Kethapadi, Krishnam (2019). Fairness-aware ranking in search & recommendation systems with application to LinkedIn talent search. In Ankur Teredesai, Vipin Kumar, Ying Li, Rómer Rosales, Evimaria Terzi, & George Karypis (Eds.), *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 2221–2231). ACM, <http://dx.doi.org/10.1145/3292500.3330691>.
- Gharahighchi, Alireza, Vens, Celine, & Plakos, Konstantinos (2021). Fair multi-stakeholder news recommender system with hypergraph ranking. *Information Processing & Management*, 58(5), Article 102663. <http://dx.doi.org/10.1016/j.ipm.2021.102663>.

- Ghosh, Avijit, Dutt, Ritam, & Wilson, Christo (2021). When fair ranking meets uncertain inference. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, & Tetsuya Sakai (Eds.), *SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval* (pp. 1033–1043). ACM, <http://dx.doi.org/10.1145/3404835.3462850>.
- Gómez, Elizabeth, Boratto, Ludovico, & Salamó, Maria (2022). Provider fairness across continents in collaborative recommender systems. *Information Processing & Management*, 59(1), Article 102719. <http://dx.doi.org/10.1016/j.ipm.2021.102719>.
- Gómez, Elizabeth, Zhang, Carlos Shui, Boratto, Ludovico, Salamó, Maria, & Marras, Mirko (2021). The winner takes it all: Geographic imbalance and provider (un)fairness in educational recommender systems. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, & Tetsuya Sakai (Eds.), *SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval* (pp. 1808–1812). ACM, <http://dx.doi.org/10.1145/3404835.3463235>.
- Gunawardana, Asela, & Shani, Guy (2015). Evaluating recommender systems. In Francesco Ricci, Lior Rokach, & Bracha Shapira (Eds.), *Recommender systems handbook* (pp. 265–308). Springer, http://dx.doi.org/10.1007/978-1-4899-7637-6_8.
- Gupta, Ananya, Johnson, Eric, Payan, Justin, Roy, Aditya Kumar, Kobren, Ari, Panda, Swetasudha, et al. (2021). Online post-processing in rankings for fair utility maximization. In Liane Lewin-Eytan, David Carmel, Elad Yom-Tov, Eugene Agichtein, Evgeniy Gabrilovich (Eds.), *WSDM '21, the fourteenth ACM international conference on web search and data mining* (pp. 454–462). ACM, <http://dx.doi.org/10.1145/3437963.3441724>.
- Hardt, Moritz, Price, Eric, & Srebro, Nati (2016). Equality of opportunity in supervised learning. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, & Roman Garnett (Eds.), *Advances in neural information processing systems 29: Annual conference on neural information processing systems 2016* (pp. 3315–3323). <https://proceedings.neurips.cc/paper/2016/hash/9d2682367c3935defcb1f9e247a97c0d-Abstract.html>.
- He, Xiangnan, Liao, Lizi, Zhang, Hanwang, Nie, Liqiang, Hu, Xia, & Chua, Tat-Seng (2017). Neural collaborative filtering. In Rick Barrett, Rick Cummings, Eugene Agichtein, & Evgeniy Gabrilovich (Eds.), *Proceedings of the 26th international conference on World Wide Web* (pp. 173–182). ACM, <http://dx.doi.org/10.1145/3038912.3052569>.
- Hu, Yifan, Koren, Yehuda, & Volinsky, Chris (2008). Collaborative filtering for implicit feedback datasets. In *Proceedings of the 8th IEEE international conference on data mining* (pp. 263–272). IEEE Computer Society, <http://dx.doi.org/10.1109/ICDM.2008.22>.
- Kirnap, Ömer, Diaz, Fernando, Biega, Asia, Ekstrand, Michael D., Carterette, Ben, & Yilmaz, Emine (2021). Estimation of fair ranking metrics with incomplete judgments. In Jure Leskovec, Marko Grobelnik, Marc Najork, Jie Tang, & Leila Zia (Eds.), *WWW '21: The web conference 2021* (pp. 1065–1075). ACM/TW3C2, <http://dx.doi.org/10.1145/3442381.3450080>.
- Koren, Yehuda, Bell, Robert M., & Volinsky, Chris (2009). Matrix factorization techniques for recommender systems. *Computer*, 42(8), 30–37. <http://dx.doi.org/10.1109/MC.2009.263>.
- Krakovsky, Marina (2022). Formalizing fairness. *Communications of the ACM*, 65(8), 11–13. <http://dx.doi.org/10.1145/3542815>.
- Kuhlman, Caitlin, Valkenburg, MaryAnn Van, & Rundensteiner, Elke A. (2019). FARE: Diagnostics for fair ranking using pairwise error metrics. In Ling Liu, Ryan W. White, Amin Mantrach, Fabrizio Silvestri, Julian J. McAuley, Ricardo Baeza-Yates, & Leila Zia (Eds.), *The World Wide Web conference* (pp. 2936–2942). ACM, <http://dx.doi.org/10.1145/3308558.3313443>.
- Kusner, Matt J., Loftus, Joshua R., Russell, Chris, & Silva, Ricardo (2017). Counterfactual fairness. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, & Roman Garnett (Eds.), *Advances in neural information processing systems 30: Annual conference on neural information processing systems 2017* (pp. 4066–4076). <https://proceedings.neurips.cc/paper/2017/hash/a486cd07e4ac3d270571622f4f316ec5-Abstract.html>.
- Li, Yunqi, Chen, Hanxiong, Fu, Zuohui, Ge, Yingqiang, & Zhang, Yongfeng (2021). User-oriented fairness in recommendation. In Jure Leskovec, Marko Grobelnik, Marc Najork, Jie Tang, & Leila Zia (Eds.), *WWW '21: The web conference 2021* (pp. 624–632). ACM/TW3C2, <http://dx.doi.org/10.1145/3442381.3449866>.
- Li, Yunqi, Chen, Hanxiong, Xu, Shuyuan, Ge, Yingqiang, Tan, Juntao, Liu, Shuchang, et al. (2022). Fairness in recommendation: A survey. <http://dx.doi.org/10.48550/arXiv.2205.13619>, CoRR abs/2205.13619. [arXiv:2205.13619](https://arxiv.org/abs/2205.13619).
- Li, Yunqi, Chen, Hanxiong, Xu, Shuyuan, Ge, Yingqiang, & Zhang, Yongfeng (2021). Towards personalized fairness based on causal notion. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, & Tetsuya Sakai (Eds.), *SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval* (pp. 1054–1063). ACM, <http://dx.doi.org/10.1145/3404835.3462966>.
- Li, Yunqi, Ge, Yingqiang, & Zhang, Yongfeng (2021). Tutorial on fairness of machine learning in recommender systems. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, & Tetsuya Sakai (Eds.), *SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval* (pp. 2654–2657). ACM, <http://dx.doi.org/10.1145/3404835.3462814>.
- Li, Yangkun, Hedia, Mohamed-Laid, Ma, Weizhi, Lu, Hongyu, Zhang, Min, Liu, Yiqun, et al. (2022). Contextualized fairness for recommender systems in premium scenarios. *Big Data Research*, 27, Article 100300. <http://dx.doi.org/10.1016/j.bdr.2021.100300>, <https://www.sciencedirect.com/science/article/pii/S2214579621001179>.
- Liang, Dawen, Krishnan, Rahul G., Hoffman, Matthew D., & Jebara, Tony (2018). Variational autoencoders for collaborative filtering. In Pierre-Antoine Champin, Fabien Gandon, Mounia Lalmas, & Panagiotis G. Ipeirotis (Eds.), *Proceedings of the 2018 World Wide Web conference on World Wide Web* (pp. 689–698). ACM, <http://dx.doi.org/10.1145/3178876.3186150>.
- Lin, Chen, Liu, Xinyi, Xv, Guipeng, & Li, Hui (2021). Mitigating sentiment bias for recommender systems. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, & Tetsuya Sakai (Eds.), *SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval* (pp. 31–40). ACM, <http://dx.doi.org/10.1145/3404835.3462943>.
- Liu, Weiben, Guo, Jun, Sonboli, Nasim, Burke, Robin, & Zhang, Shengyu (2019). Personalized fairness-aware re-ranking for microlending. In Toine Bogers, Alan Said, Peter Brusilovsky, & Domonkos Tikk (Eds.), *Proceedings of the 13th ACM conference on recommender systems* (pp. 467–471). ACM, <http://dx.doi.org/10.1145/3298689.3347016>.
- Mehrabi, Ninareh, Morstatter, Fred, Saxena, Nripsuta, Lerman, Kristina, & Galstyan, Aram (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6), 115:1–115:35. <http://dx.doi.org/10.1145/3457607>.
- Mehrotra, Rishabh, Anderson, Ashton, Diaz, Fernando, Sharma, Amit, Wallach, Hanna M., & Yilmaz, Emine (2017). Auditing search engines for differential satisfaction across demographics. In Rick Barrett, Rick Cummings, Eugene Agichtein, & Evgeniy Gabrilovich (Eds.), *Proceedings of the 26th international conference on World Wide Web Companion* (pp. 626–633). ACM, <http://dx.doi.org/10.1145/3041021.3054197>.
- Mehrotra, Rishabh, McInerney, James, Bouchard, Hugues, Lalmas, Mounia, & Diaz, Fernando (2018). Towards a fair marketplace: Counterfactual evaluation of the trade-off between relevance, fairness & satisfaction in recommendation systems. In Alfredo Cuzzocrea, James Allan, Norman W. Paton, Divesh Srivastava, Rakesh Agrawal, Andrei Z. Broder, Mohammed J. Zaki, K. Selçuk Candan, Alexandros Labrinidis, Assaf Schuster, & Haixun Wang (Eds.), *Proceedings of the 27th ACM international conference on information and knowledge management* (pp. 2243–2251). ACM, <http://dx.doi.org/10.1145/3269206.3272027>.
- Modani, Natwar, Jain, Deepali, Soni, Ujjawal, Gupta, Gaurav Kumar, & Agarwal, Palak (2017). Fairness aware recommendations on behavior. In Jinho Kim, Kyuseok Shim, Longbing Cao, Jae-Gil Lee, Xuemin Lin, & Yang-Sae Moon (Eds.), *Lecture notes in computer science: Vol. 10235, Advances in knowledge discovery and data mining - 21st Pacific-Asia conference* (pp. 144–155). http://dx.doi.org/10.1007/978-3-319-57529-2_12.
- Moulin, Hervé (2003). *Fair division and collective welfare*. MIT Press.
- Naghiaei, Mohammadmehdi, Rahmani, Hossein A., & Deldjoo, Yashar (2022). CPFair: Personalized consumer and producer fairness re-ranking for recommender systems. In *SIGIR '22: The 45th international ACM SIGIR conference on research and development in information retrieval* (pp. 770–779). ACM, <http://dx.doi.org/10.1145/3477495.3531959>.

- Pan, Rong, Zhou, Yunhong, Cao, Bin, Liu, Nathan Nan, Lukose, Rajan M., Scholz, Martin, et al. (2008). One-class collaborative filtering. In *Proceedings of the 8th IEEE international conference on data mining* (pp. 502–511). IEEE Computer Society.
- Patro, Gourab K., Biswas, Arpita, Ganguly, Niloy, Gummadi, Krishna P., & Chakraborty, Abhijnan (2020). FairRec: Two-sided fairness for personalized recommendations in two-sided platforms. In Yennun Huang, Irwin King, Tie-Yan Liu, & Maarten van Steen (Eds.), *WWW '20: The web conference 2020* (pp. 1194–1204). ACM/IW3C2, <http://dx.doi.org/10.1145/3366423.3380196>.
- Patro, Gourab K., Porcaro, Lorenzo, Mitchell, Laura, Zhang, Qiuyue, Zehlike, Meike, & Garg, Nikhil (2022). Fair ranking: A critical review, challenges, and future directions. In *FACt'22: 2022 ACM conference on fairness, accountability, and transparency* (pp. 1929–1942). ACM, <http://dx.doi.org/10.1145/3531146.3533238>.
- Pitoura, Evangelia, Stefanidis, Kostas, & Koutrika, Georgia (2022). Fairness in rankings and recommendations: An overview. *Vldb Journal*, 31(3), 431–458. <http://dx.doi.org/10.1007/s00778-021-00697-y>.
- Rendle, Steffen, Freudenthaler, Christoph, Gantner, Zeno, & Schmidt-Thieme, Lars (2009). BPR: Bayesian personalized ranking from implicit feedback. In Jeff A. Birmes, & Andrew Y. Ng (Eds.), *UAI 2009, proceedings of the twenty-fifth conference on uncertainty in artificial intelligence* (pp. 452–461). AUAI Press.
- Sacharidis, Dimitris, Mouratidis, Kyriakos, & Klefogiannis, Dimitrios (2019). A common approach for consumer and provider fairness in recommendations. In Marko Tkalcić, & Sole Pera (Eds.), *CEUR workshop proceedings: Vol. 2431, Proceedings of ACM RecSys 2019 late-breaking results co-located with the 13th ACM conference on recommender systems* (pp. 1–5). CEUR-WS.org, <http://ceur-ws.org/Vol-2431/paper1.pdf>.
- Salah, Aghiles, Truong, Quoc-Tuan, & Lauw, Hady W. (2020). Cornac: A comparative framework for multimodal recommender systems. *Journal of Machine Learning Research*, 21(95), 1–5, <http://jmlr.org/papers/v21/19-805.html>.
- Salakhutdinov, Ruslan, & Mnih, Andriy (2007). Probabilistic matrix factorization. In John C. Platt, Daphne Koller, Yoram Singer, & Sam T. Roweis (Eds.), *Advances in neural information processing systems 20, proceedings of the twenty-first annual conference on neural information processing systems* (pp. 1257–1264). Curran Associates, Inc..
- Saracevic, Tefko (1975). RELEVANCE: A review of and a framework for the thinking on the notion in information science. *Journal of the American Society for Information Science*, 26(6), 321–343. <http://dx.doi.org/10.1002/asi.4630260604>.
- da Silva, Diego Corrêa, Manzato, Marcelo Garcia, & Durão, Frederico Araújo (2021). Exploiting personalized calibration and metrics for fairness recommendation. *Expert Systems with Applications*, 181, Article 115112. <http://dx.doi.org/10.1016/j.eswa.2021.115112>.
- Simoiu, Camelia, Corbett-Davies, Sam, & Goel, Sharad (2016). The problem of infra-marginality in outcome tests for discrimination. *Econometrics: Econometric & Statistical Methods - General EJournal*.
- Singh, Ashudeep, & Joachims, Thorsten (2018). Fairness of exposure in rankings. In Yike Guo, & Faisal Farooq (Eds.), *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 2219–2228). ACM, <http://dx.doi.org/10.1145/3219819.3220088>.
- Steck, Harald (2018). Calibrated recommendations. In Sole Pera, Michael D. Ekstrand, Xavier Amatriain, & John O'Donovan (Eds.), *Proceedings of the 12th ACM conference on recommender systems* (pp. 154–162). ACM, <http://dx.doi.org/10.1145/3240323.3240372>.
- Stoica, Ana-Andreea, Riederer, Christopher, & Chaintreau, Augustin (2018). Algorithmic glass ceiling in social networks: The effects of social recommendations on network diversity. In *WWW '18: Proceedings of the 2018 World Wide Web conference* (pp. 923–932). <http://dx.doi.org/10.1145/3178876.3186140>.
- Tsintzou, Virginia, Pitoura, Evangelia, & Tsaparas, Panayiotis (2019). Bias disparity in recommendation systems. In Robin Burke, Himan Abdollahpour, Edward C. Malthouse, K. P. Thai, & Yongfeng Zhang (Eds.), *CEUR workshop proceedings: Vol. 2440, Proceedings of the workshop on recommendation in multi-stakeholder environments co-located with the 13th ACM conference on recommender systems*. CEUR-WS.org, <http://ceur-ws.org/Vol-2440/short4.pdf>.
- Vargas, Saul, & Castells, Pablo (2011). Rank and relevance in novelty and diversity metrics for recommender systems. In Bamshad Mobasher, Robin D. Burke, Dietmar Jannach, & Gediminas Adomavicius (Eds.), *Proceedings of the 2011 ACM conference on recommender systems* (pp. 109–116). ACM, <https://dl.acm.org/citation.cfm?id=2043955>.
- Verma, Sahil, & Rubin, Julia (2018). Fairness definitions explained. In Yuriy Brun, Brittany Johnson, & Alexandra Meliou (Eds.), *Proceedings of the international workshop on software fairness* (pp. 1–7). ACM, <http://dx.doi.org/10.1145/3194770.3194776>.
- Wan, Mengting, Ni, Jianmo, Misra, Rishabh, & McAuley, Julian J. (2020). Addressing marketing bias in product recommendations. In James Caverlee, Xia (Ben) Hu, Mounia Lalmas, & Wei Wang (Eds.), *WSDM '20: The thirteenth ACM international conference on web search and data mining* (pp. 618–626). ACM, <http://dx.doi.org/10.1145/3336191.3371855>.
- Wang, Chong, & Blei, David M. (2011). Collaborative topic modeling for recommending scientific articles. In Chid Apté, Joydeep Ghosh, & Padhraic Smyth (Eds.), *Proceedings of the 17th ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 448–456). ACM, <http://dx.doi.org/10.1145/2020408.2020480>.
- Wang, Lequn, & Joachims, Thorsten (2021). User fairness, item fairness, and diversity for rankings in two-sided markets. In Faegheh Hasibi, Yi Fang, & Akiko Aizawa (Eds.), *ICTIR '21: The 2021 ACM SIGIR international conference on the theory of information retrieval* (pp. 23–41). ACM, <http://dx.doi.org/10.1145/3471158.3472260>.
- Wang, Yifan, Ma, Weizhi, Zhang, Min, Liu, Yiqun, & Ma, Shaoping (2022). A survey on the fairness of recommender systems. *ACM Transactions on Information Systems*, arXiv:2206.03761.
- Wu, Yao, Cao, Jian, Xu, Guandong, & Tan, Yudong (2021). TFROM: A two-sided fairness-aware recommendation model for both customers and providers. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, & Tetsuya Sakai (Eds.), *SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval* (pp. 1013–1022). ACM, <http://dx.doi.org/10.1145/3404835.3462882>.
- Wu, Haolun, Mitra, Bhaskar, Ma, Chen, Diaz, Fernando, & Liu, Xue (2022). Joint multisided exposure fairness for recommendation. In Enrique Amigó, Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, & Gabriella Kazai (Eds.), *SIGIR '22: The 45th international ACM SIGIR conference on research and development in information retrieval* (pp. 703–714). ACM, <http://dx.doi.org/10.1145/3477495.3532007>.
- Wundervald, Bruna D. (2021). Cluster-based quotas for fairness improvements in music recommendation systems. *International Journal of Multimedia Information Retrieval*, 10(1), 25–32. <http://dx.doi.org/10.1007/s13735-020-00203-0>.
- Yalcin, Emre, & Bilge, Alper (2021). Investigating and counteracting popularity bias in group recommendations. *Information Processing & Management*, 58(5), Article 102608. <http://dx.doi.org/10.1016/j.ipm.2021.102608>.
- Yang, Ke, & Stoyanovich, Julia (2017). Measuring fairness in ranked outputs. In *Proceedings of the 29th international conference on scientific and statistical database management* (pp. 22:1–22:6). ACM, <http://dx.doi.org/10.1145/3085504.3085526>.
- Yao, Sirui, & Huang, Bert (2017). Beyond parity: Fairness objectives for collaborative filtering. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, & Roman Garnett (Eds.), *Advances in neural information processing systems 30: Annual conference on neural information processing systems 2017* (pp. 2921–2930). <https://proceedings.neurips.cc/paper/2017/hash/e6384711491713d29bc63f5eeb5ba4f-Abstract.html>.
- Zafar, Muhammad Bilal, Valera, Isabel, Gomez-Rodriguez, Manuel, & Gummadi, Krishna P. (2017). Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In Rick Barrett, Rick Cummings, Eugene Agichtein, & Evgeniy Gabrilovich (Eds.), *Proceedings of the 26th international conference on World Wide Web* (pp. 1171–1180). ACM, <http://dx.doi.org/10.1145/3038912.3052660>.
- Zehlike, Meike, Bonchi, Francesco, Castillo, Carlos, Hajian, Sara, Megahed, Mohamed, & Baeza-Yates, Ricardo (2017). FA*IR: A fair top-k ranking algorithm. In Ee-Peng Lim, Marianne Winslett, Mark Sanderson, Ada Wai-Chee Fu, Jimeng Sun, J. Shane Culpepper, Eric Lo, Joyce C. Ho, Debora Donato, Rakesh Agrawal, Yu Zheng, Carlos Castillo, Aixin Sun, Vincent S. Tseng, & Chenliang Li (Eds.), *Proceedings of the 2017 ACM conference on information and knowledge management* (pp. 1569–1578). ACM, <http://dx.doi.org/10.1145/3132847.3132938>.

- Zhu, Ziwei, Hu, Xia, & Caverlee, James (2018). Fairness-aware tensor-based recommendation. In Alfredo Cuzzocrea, James Allan, Norman W. Paton, Divesh Srivastava, Rakesh Agrawal, Andrei Z. Broder, Mohammed J. Zaki, K. Selçuk Candan, Alexandros Labrinidis, Assaf Schuster, & Haixun Wang (Eds.), *Proceedings of the 27th ACM international conference on information and knowledge management* (pp. 1153–1162). ACM, <http://dx.doi.org/10.1145/3269206.3271795>.
- Zhu, Qiliang, Sun, Qibo, Li, Zengxiang, & Wang, Shangguang (2020). FARM: A fairness-aware recommendation method for high visibility and low visibility mobile APPs. *IEEE Access*, 8, 122747–122756. <http://dx.doi.org/10.1109/ACCESS.2020.3007617>.
- Zliobaite, Indre (2015). On the relation between accuracy and fairness in binary classification. In *2nd workshop on fairness, accountability, and transparency in machine learning*.
- Zliobaite, Indre (2017). Measuring discrimination in algorithmic decision making. *Data Mining and Knowledge Discovery*, 31(4), 1060–1089. <http://dx.doi.org/10.1007/s10618-017-0506-1>.