

FA*IR: A Fair Top-k Ranking Algorithm

Meike Zehlike
TU Berlin
Berlin, Germany
meike.zehlike@tu-berlin.de

Francesco Bonchi
ISI Foundation
Turin, Italy
francesco.bonchi@isi.it

Carlos Castillo
Universitat Pompeu Fabra
Barcelona, Catalunya, Spain
chato@acm.org

Sara Hajian
NTENT
Barcelona, Catalunya, Spain
shajian@ntent.com

Mohamed Megahed
TU Berlin
Berlin, Germany
mohamed.megahed@campus.
tu-berlin.de

Ricardo Baeza-Yates
Universitat Pompeu Fabra
Barcelona, Catalunya, Spain
rbaeza@acm.org

ABSTRACT

In this work, we define and solve the Fair Top- k Ranking problem, in which we want to determine a subset of k candidates from a large pool of $n \gg k$ candidates, maximizing utility (i.e., select the “best” candidates) subject to group fairness criteria.

Our ranked group fairness definition extends group fairness using the standard notion of protected groups and is based on ensuring that the proportion of protected candidates in every prefix of the top- k ranking remains statistically above or indistinguishable from a given minimum. Utility is operationalized in two ways: (i) every candidate included in the top- k should be more qualified than every candidate not included; and (ii) for every pair of candidates in the top- k , the more qualified candidate should be ranked above.

An efficient algorithm is presented for producing the Fair Top- k Ranking, and tested experimentally on existing datasets as well as new datasets released with this paper, showing that our approach yields small distortions with respect to rankings that maximize utility without considering fairness criteria. To the best of our knowledge, this is the first algorithm grounded in statistical tests that can mitigate biases in the representation of an under-represented group along a ranked list.

KEYWORDS

Algorithmic fairness, Bias in Computer Systems, Ranking, Top- k selection.

1 INTRODUCTION

People search engines are increasingly common for job recruiting and even for finding companionship or friendship. A top- k ranking algorithm is typically used to find the most suitable way of ordering items (persons, in this case), considering that if the number of people matching a query is large, most users will not scan the entire list. Conventionally, these lists are ranked in descending order of some measure of the relative quality of items.

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

CIKM'17, November 6–10, 2017, Singapore.

© 2017 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4918-5/17/11...\$15.00

DOI: <https://doi.org/10.1145/3132847.3132938>

The main concern motivating this paper is that a biased machine learning model that produces ranked lists can further systematically reduce the visibility of an already disadvantaged group [10, 31] (corresponding to a legally protected category such as people with disabilities, racial or ethnic minorities, or an under-represented gender in a specific industry).

According to [14] a computer system is *biased* “if it systematically and unfairly discriminate[s] against certain individuals or groups of individuals in favor of others. A system discriminates unfairly if it denies an opportunity or a good or if it assigns an undesirable outcome to an individual or a group of individuals on grounds that are unreasonable or inappropriate.” Yet “unfair discrimination alone does not give rise to bias unless it occurs systematically” and “systematic discrimination does not establish bias unless it is joined with an unfair outcome.” On a ranking, the desired good for an individual is to appear in the result and to be ranked amongst the top- k positions. The outcome is unfair if members of a protected group are systematically ranked lower than those of a privileged group. The ranking algorithm discriminates unfairly if this ranking decision is based fully or partially on the protected feature. This discrimination is systematic when it is embodied in the algorithm’s ranking model. As shown in earlier research, a machine learning model trained on datasets incorporating *preexisting bias* will embody this bias and therefore produce biased results, potentially increasing any disadvantage further, reinforcing existing bias [28].

Based on this observation, in this paper we study the problem of producing a fair ranking given one legally-protected attribute,¹ i.e., a ranking in which the representation of the minority group does not fall below a minimum proportion p at any point in the ranking, while the utility of the ranking is maintained as high as possible.

We propose a post-processing method to remove the systematic bias by means of a *ranked group fairness criterion*, that we introduce in this paper. We assume a ranking algorithm has given an undesirable outcome to a group of individuals, but the algorithm itself cannot determine if the grounds were appropriate or not. Hence we expect the user of our method to know that the outcome is based on unreasonable or inappropriate grounds and provide p as input which can originate in a legal mandate or in voluntary commitments. For instance, the US Equal Employment Opportunity Commission sets a goal of 12% of workers with disabilities in

¹We make the simplifying assumption that there is a dominant legally-protected attribute of interest in each case. The extension to deal with multiple protected attributes is left for future work.

federal agencies in the US,² while in Spain, a minimum of 40% of political candidates in voting districts exceeding a certain size must be women [35]. In other cases, such quotas might be adopted voluntarily, for instance through a diversity charter.³ In general these measures do not mandate perfect parity, as distributions of qualifications across groups can be unbalanced for legitimate, explainable reasons [29, 39].

The ranked group fairness criterion compares the number of protected elements in every prefix of the ranking with the expected number of protected elements if they were picked at random using Bernoulli trials (independent “coin tosses”) with success probability p . Given that we use a statistical test for this comparison, we include a significance parameter α corresponding to the probability of a Type I error, which means rejecting a fair ranking in this test.

Example. Consider the three rankings in Table 1 corresponding to searches for an “economist,” “market research analyst,” and “copywriter” in XING⁴, an online platform for jobs that is used by recruiters and headhunters, mostly in German-speaking countries, to find suitable candidates in diverse fields (this data collection is reported in detail on §5.2). While analyzing the extent to which candidates of both genders are represented as we go down these lists, we can observe that such proportion keep changing and is not uniform (see, for instance, the top-10 vs. the top-40). As a consequence, recruiters examining these lists will see different proportions depending on the point at which they decide to stop. Corresponding with [14] this outcome systematically disadvantages individuals of one gender by preferring the other at the top- k positions. As we do not know the learning model behind the ranking, we assume that the result is at least partly based on the protected attribute *gender*.

Let $k = 10$. Our notion of *ranked group fairness* imposes a fair representation with proportion p and significance α at each top- i position with $i \in [1, 10]$ (formal definitions are given in §3). Consider for instance $\alpha = 0.1$ and suppose that the required proportion is $p = 0.4$. This translates (see Table 2) to having at least one individual from the protected minority class in the first 5 positions: therefore the ranking for “copywriter” would be rejected as unfair. However, it also requires to have at least 2 individuals from the protected group in the first 9 positions: therefore also the ranking for “economist” is rejected as unfair, while the ranking for “market research analyst” is fair for $p = 0.4$. However, if we would require $p = 0.5$ then this translates in having at least 3 individuals from the minority group in the top-10, and thus even the ranking for “market research analyst” would be considered unfair. We note that for simplicity, in this example we have not adjusted the significance α to account for multiple statistical tests; this is not trivial, and is one of the key contributions of this paper.

Our contributions. In this paper, we define and analyze the FAIR TOP- k RANKING PROBLEM, in which we want to determine a subset of k candidates from a large pool of $n \gg k$ candidates, in a way that maximizes utility (selects the “best” candidates), subject to group fairness criteria. The running example we use in this paper is that of selecting automatically, from a large pool of potential candidates, a smaller group that will be interviewed for a position.

Table 1: Example of non-uniformity of the top-10 vs. the top-40 results for different queries in XING (Jan 2017).

	Position										top 10	top 10	top 40	top 40
	1	2	3	4	5	6	7	8	9	10	male	female	male	female
Econ.	f	m	m	m	m	m	m	m	m	m	90%	10%	73%	27%
Analyst	f	m	f	f	f	f	f	m	f	f	20%	80%	43%	57%
Copywr.	m	m	m	m	m	m	f	m	m	m	90%	10%	73%	27%

Our notion of utility assumes that we want the interview the most qualified candidates, while their qualification is equal to a relevance score calculated by a ranking algorithm. This score is assumed to be based on relevant metrics for evaluating candidates for a position, which depending on the specific skills required for the job could be their grades (e.g., Grade Point Average), their results in a standardized knowledge/skills test specific for a job, their typing speed in words per minute for typists, or their number of hours of flight in the case of pilots. We note that this measurement will embody *preexisting bias* (e.g. if black pilots are given less opportunities to flight they accumulate less flight hours), as well as *technical bias*, as learning algorithms are known to be susceptible to direct and indirect discrimination [15, 16].

The utility objective is operationalized in two ways. First, by a criterion we call *selection utility*, which prefers rankings in which every candidate included in the top- k is more qualified than every candidate not included, or in which the difference in their qualifications is small. Second, by a criterion we call *ordering utility*, which prefers rankings in which for every pair of candidates included in the top- k , either the more qualified candidate is ranked above, or the difference in their qualifications is small.

Our definition of *ranked group fairness* reflects the legal principle of group under-representation in obtaining a benefit [11, 26]. We use the standard notion of a protected group (e.g., “people with disabilities”); where protection emanates from a legal mandate or a voluntary commitment. We formulate a criterion applying a statistical test on the proportion of protected candidates on every prefix of the ranking, which should be indistinguishable or above a given minimum. We also show that the verification of the ranked group fairness criterion can be implemented efficiently.

Finally, we propose an efficient algorithm, named FA*IR, for producing a top- k ranking that maximizes utility while satisfying ranked group fairness, as long as there are “enough” protected candidates to achieve the desired minimum proportion. We also present extensive experiments using both existing and new datasets to evaluate the performance of our approach compared to the so-called “color-blind” ranking with respect to both the utility of ranking and the fairness degree.

Summarizing, the main contributions of this paper are:

- (1) the principled definition of *ranked group fairness*, and the associated FAIR TOP- k RANKING PROBLEM;
- (2) the FA*IR algorithm for producing a top- k ranking that maximizes utility while satisfying ranked group fairness.

Our method can be used within an anti-discrimination framework such as *positive actions* [33]. We do not claim these are the only way of achieving fairness, but we provide *an algorithm grounded in statistical tests that enables the implementation of a positive action policy in the context of ranking*.

²US EEOC: <https://www1.eeoc.gov/eeoc/newsroom/release/1-3-17.cfm>, Jan 2017.

³European Commission: <http://ec.europa.eu/justice/discrimination/diversity/charters/>

⁴<https://www.xing.com/>

The rest of this paper is organized as follows. The next section presents a brief survey of related literature, while Section 3 introduces our ranked group fairness and utility criteria, our model adjustment approach, and a formal problem statement. Section 4 describes the FA*IR algorithm. Section 5 presents experimental results. Section 6 presents our conclusions and future work.

2 RELATED WORK

Anti-discrimination has only recently been considered from an algorithmic perspective [15]. Some proposals are oriented to discovering and measuring discrimination (e.g., [1, 3, 31]); while others deal with mitigating or removing discrimination (e.g., [10, 16, 17, 21, 38]). All these methods are known as *fairness-aware algorithms*.

2.1 Group fairness and individual fairness

Two basic frameworks have been adopted in recent studies on algorithmic discrimination: (i) *individual fairness*, a requirement that individuals should be treated consistently [10]; and (ii) *group fairness*, also known as statistical parity, a requirement that the protected groups should be treated similarly to the advantaged group or the populations as a whole [30, 31].

Different fairness-aware algorithms have been proposed to achieve group and/or individual fairness, mostly for predictive tasks. Calders and Verwer [4] consider three approaches to deal with naive Bayes models by modifying the learning algorithm. Kamiran et al. [21] modify the entropy-based splitting criterion in decision tree induction to account for attributes denoting protected groups. Kamishima et al. [22] apply a regularization (i.e., a change in the objective minimization function) to probabilistic discriminative models, such as logistic regression. Zafar et al. [37] describe fairness constraints for several classification methods.

Feldman et al. [12] study *disparate impact* in data, which corresponds to an unintended form of group discrimination, in which a protected group is less likely to receive a benefit than a non-protected group [2]. Besides measuring disparate impact, the authors also propose a method for removing it: we use this method as one of our experimental baselines in §5.3. Their method “repairs” the scores of the protected group to make them have the same or similar distribution as the scores of the non-protected group, which is one particular form of positive action. Recently, other fairness-aware algorithms have been proposed for mostly supervised learning algorithms and different bias mitigation strategies [6, 9, 13, 18, 19].

2.2 Fair Ranking

Yang and Stoyanovich [36] studied the problem of fairness in rankings. They propose a statistical parity measure based on comparing the distributions of protected and non-protected candidates (for instance, using KL-divergence) on different prefixes of the list (e.g., top-10, top-20, top-30) and then averaging these differences in a discounted manner. The discount used is logarithmic, similarly to Normalized Discounted Cumulative Gain (NDCG, a popular measure used in Information Retrieval [20]). Finally, they show very preliminary results on incorporating their statistical parity measure into an optimization framework for improving fairness of ranked outputs while maintaining accuracy. We use the synthetic ranking

generation procedure of Yang and Stoyanovich [36] to calibrate our method, and optimize directly the utility of a ranking that has statistical properties (ranked group fairness) resembling the ones of a ranking generated using that procedure; in other words, unlike [36], we connect the creation of the ranking with the metric used for assessing fairness.

Kulshrestha et al. [24] determine search bias in rankings, proposing a quantification framework that measures the bias of the results of a search engine. This framework discerns to what extent this output bias is due to the input dataset that feeds into the ranking system, and how much is due to the bias introduced by the system itself. In contrast to their work, which mostly focus on auditing ranking algorithms to identify the sources of bias in the data or algorithm, our paper focuses on generating fair ranked results.

A recent work [7] proposes algorithms for constrained ranking, in which the constraint is a $k \times \ell$ matrix with k the length of the ranking and ℓ the number of classes, indicating the maximum number of elements of each class (protected or non-protected in the binary case) that can appear at any given position in the ranking. The objective is to maximize a general utility function that has a positional discount, i.e., gives more weight to placing a candidate with more qualifications in a higher position. Differently from [7], in our work we show how to construct the constraint matrix by means of a statistical test of ranked group fairness (a problem they left open), and our measure of utility is based on individuals, which allows to determine which individuals are more affected by the re-ranking with respect to the non-fairness-aware solution §3.4.

2.3 Diversity

Additionally, the idea that we want to avoid showing only items of the same class has been studied in the Information Retrieval community for many years. The motivation there is that the user query may have different intents and we want to cover several of them with results. The most common approach, since Carbonell and Goldstein [5], is to consider distances between elements, and maximize a combination of relevance (utility) with a penalty for adding to the ranking an element that is too similar to an element already appearing at a higher position. A similar idea is used for diversification in recommender systems through various methods [8, 25]. They deal with different kinds of bias such as presentation bias, where, only a few items are shown and most of the items are not shown, and also popularity bias and a negative bias towards new items. An exception is Sakai and Song [32], that provides a framework for per-intent NDCG for evaluating diversity, in which an “intent” could be mapped to a protected/non-protected group in the fairness ranking setting. Their method, however, is concerned with evaluating a ranking, similar to the NDCG-based metrics described by Yang and Stoyanovich [36] that we describe before, and not with a construction of such ranking, as we do in this paper. In contrast with most of the research on diversity of ranking results or recommender systems, our work operates on a discrete set of classes (not based on similarity to previous items).

3 THE FAIR TOP-K RANKING PROBLEM

In this section, we first present the needed notation (§3.1), then the ranked group fairness criterion (§3.2-§3.3) and criteria for utility (§3.4). Finally we provide a formal problem statement (§3.5).

3.1 Preliminaries and Notation

Notation. Let $[n] = \{1, 2, \dots, n\}$ represent a set of candidates; let q_i for $i \in [n]$ denote the “quality” of candidate i : this can be interpreted as an overall summary of the fitness of candidate i for the specific job, task, or search query, and that could be obtained by the combination of several different attributes, possibly by means of a machine learning model, and potentially including preexisting and technical bias with respect to the protected group. We will consider two kinds of candidates, protected and non-protected, and we will assume there are enough of them, i.e., at least k of each kind. Let $g_i = 1$ if candidate i is in the protected group, $g_i = 0$ otherwise. Let $\mathcal{P}_{k,n}$ represent all the subsets of $[n]$ containing exactly k elements. Let $\mathcal{T}_{k,n}$ represent the union of all permutations of sets in $\mathcal{P}_{k,n}$. For a permutation $\tau \in \mathcal{T}_{k,n}$ and an element $i \in [n]$, let

$$r(i, \tau) = \begin{cases} \text{rank of } i \text{ in } \tau & \text{if } i \in \tau, \\ |\tau| + 1 & \text{otherwise.} \end{cases}$$

We further define τ_p to be the number of protected elements in τ , i.e., $\tau_p = |\{i \in \tau : g_i = 1\}|$. Let $c \in \mathcal{T}_{n,n}$ be a permutation such that $\forall i, j \in [n], r(i, c) < r(j, c) \Rightarrow q_i \geq q_j$. We call this the *color-blind* ranking of $[n]$, because it ignores whether elements are protected or non-protected. Let $c|_k = \langle c(1), c(2), \dots, c(k) \rangle$ be a prefix of size k of this ranking.

Fair top- k ranking criteria. We would like to obtain $\tau \in \mathcal{T}_{k,n}$ with the following characteristics, which we describe formally next:

- Criterion 1. Ranked group fairness: τ should fairly represent the protected group;
- Criterion 2. Selection utility: τ should contain the most qualified candidates; and
- Criterion 3. Ordering utility: τ should be ordered by decreasing qualifications.

We will provide a formal problem statement in §3.5, but first, we need to provide a formal definition of each of the criteria.

3.2 Group Fairness for Rankings

We operationalize criterion 1 of §3.1 by means of a *ranked group fairness criterion*, which takes as input a protected group and a minimum target proportion of protected elements in the ranking, p . Intuitively, this criterion declares the ranking as unfair if the observed proportion is far below the target one.

Specifically, the ranked group fairness criterion compares the number of protected elements in every prefix of the ranking, with the expected number of protected elements if they were picked at random using Bernoulli trials. The criterion is based on a statistical test, and we include a significance parameter (α) corresponding to the probability of rejecting a fair ranking (i.e., a Type I error).

Definition 3.1 (Fair representation condition). Let $F(x; n, p)$ be the cumulative distribution function for a binomial distribution of parameters n and p . A set $\tau \subseteq \mathcal{T}_{k,n}$, having τ_p protected candidates fairly represents the protected group with minimal proportion p and significance α , if $F(\tau_p; k, p) > \alpha$.

This is equivalent to using a statistical test where the null hypothesis H_0 is that the protected elements are represented with a sufficient proportion p_t ($p_t \geq p$), and the alternative hypothesis H_a is that the proportion of protected elements is insufficient ($p_t < p$).

Table 2: Example values of $m_{\alpha,p}(k)$, the minimum number of candidates in the protected group that must appear in the top k positions to pass the ranked group fairness criteria with $\alpha = 0.1$.

p \ k	1	2	3	4	5	6	7	8	9	10	11	12
0.1	0	0	0	0	0	0	0	0	0	0	0	0
0.2	0	0	0	0	0	0	0	0	0	0	1	1
0.3	0	0	0	0	0	0	1	1	1	1	1	2
0.4	0	0	0	0	1	1	1	1	2	2	2	3
0.5	0	0	0	1	1	1	2	2	3	3	3	4
0.6	0	0	1	1	2	2	3	3	4	4	5	5
0.7	0	1	1	2	2	3	3	4	5	5	6	6

In this test, the p-value is $F(\tau_p; k, p)$ and we reject the null hypothesis, and thus declare the ranking as unfair, if the p-value is less than or equal to the threshold α .

The ranked group fairness criterion enforces the fair representation constraint over all prefixes of the ranking:

Definition 3.2 (Ranked group fairness condition). A ranking $\tau \in \mathcal{T}_{k,n}$ satisfies the ranked group fairness condition with parameters p and α , if for every prefix $\tau|_i = \langle \tau(1), \tau(2), \dots, \tau(i) \rangle$ with $1 \leq i \leq k$, the set $\tau|_i$ satisfies the fair representation condition with proportion p and significance $\alpha_c = c(\alpha, k, p)$. Function $c(\alpha, k, p)$ is a corrected significance to account for multiple testing (described in §3.3).

We remark that a larger α means a larger probability of declaring a fair ranking as unfair. In our experiments (§5), we use a relatively conservative setting of $\alpha = 0.1$. The ranked group fairness condition can be used to create a *ranked group fairness measure*. For a ranking τ and probability p , the ranked group fairness measure is the maximum $\alpha \in [0, 1]$ for which τ satisfies the ranked group fairness condition. Larger values indicate a stricter adherence to the required number of protected elements at each position.

Verifying ranked group fairness. Note that ranked group fairness can be verified efficiently in time $O(k)$, by having a pre-computed table of the percent point function with parameters k and p , i.e., the inverse of $F(x; k, p)$. Table 2 shows an example of such a table, computed for $\alpha = 0.1$. For instance, for $p = 0.5$ we see that at least 1 candidate from the protected group is needed in the top 4 positions, and 2 protected candidates in the top 7 positions.

3.3 Model Adjustment

Our ranked group fairness definition requires an adjusted significance $\alpha_c = c(\alpha, k, p)$. This is required because it tests multiple hypotheses (k of them). If we use $\alpha_c = \alpha$, we might produce false negatives, rejecting fair rankings, at a rate larger than α .

The adjustment we propose is calibrated using the generative model of Yang and Stoyanovich [36], which creates a ranking that we will consider fair by: (i) starting with an empty list, and (ii) incrementally adding the best available protected candidate with probability p , or the best available non-protected candidate with probability $1 - p$.

Figure 1 shows the probability that a fair ranking generated with $p = 0.5$ is rejected by our ranked group fairness test with $p = 0.5$, $\alpha_c = 0.1$. The curve is computed analytically by the method

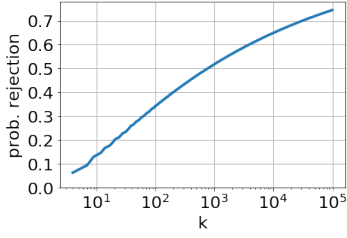


Figure 1: Example showing the need for multiple tests correction. The probability that a ranking generated by the method of Yang and Stoyanovich [36] with $p = 0.5$ fails the ranked group fairness test with $p = 0.5$ using $\alpha_c = 0.1$, is in general larger than $\alpha = 0.1$. Note the scale of k is logarithmic.

Table 3: Example of $m(\cdot)$, $m^{-1}(\cdot)$, and $b(\cdot)$ for $p = 0.5$, $\alpha = 0.1$.

k	1	2	3	4	5	6	7	8	9	10	11	12
$m(k)$	0	0	0	1	1	1	2	2	3	3	3	4
Inverse	$m^{-1}(1) = 4$				$m^{-1}(2) = 7$			$m^{-1}(3) = 9$		$m^{-1}(4) = 12$		
Blocks	$b(1) = 4$				$b(2) = 3$			$b(3) = 2$		$b(4) = 3$		

we describe in the following paragraphs, and it experimentally matches the result of simulations we performed. We can see that the probability of a Type-I error (declaring this fair ranking as unfair) is in general higher than $\alpha = 0.1$. If the k tests were independent, we could use $\alpha_c = 1 - (1 - \alpha)^{1/k}$ (i.e., Šidák's correction), but given the positive dependence, the false negative rate is smaller than the bound given by Šidák's correction.

The probability that a ranking generated using the process of Yang and Stoyanovich [36] with parameter p passes the ranked group fairness criteria where each test is done with parameters (p, α_c) can be computed as follows: Let $m(k) = m_{\alpha, p}(k)$ be as before the number of protected elements required up to position k . Let $m^{-1}(i) = k$ s.t. $m(k) = i$ be the position at which i protected elements are required. Let $b(i) = m^{-1}(i) - m^{-1}(i-1)$ (with $m^{-1}(0) = 0$) be the size of a "block," that is, the gap between one increase and the next in $m(\cdot)$. An example is shown on Table 3.

Let $I_\ell = \{(i_1, i_2, \dots, i_\ell) : \forall \ell' \in [\ell], 0 \leq i_{\ell'} \leq b(\ell') \wedge \sum_{j=1}^{\ell'} i_j \geq \ell'\}$ represent all possible ways in which a fair ranking generated by the method of Yang and Stoyanovich [36] can pass the ranked group fairness test up to block ℓ , with i_j corresponding to the number of protected elements in block $1 \leq j \leq k$. The probability of considering this ranking of k elements ($m^{-1}(k)$ blocks) unfair, is:

$$1 - \sum_{v \in I_{m^{-1}(k)}} \prod_{j=1}^{m^{-1}(k)} f(v_j; b(j), p) \quad (1)$$

where $f(x; b(j), p) = \Pr(X = x)$ is the probability density function of a binomially-distributed variable $X \sim \text{Bin}(b(j), p)$.

The above expression is intractable because of the large number of combinations in $I_{m^{-1}(k)}$; however, there is an efficient iterative process to compute this quantity, shown in Algorithm 1. This algorithm maintains a vector S that at iteration ℓ holds in position $S[i]$ the probability of having obtained i protected elements in the first ℓ blocks, conditioned on obtaining at least j protected elements up to each block $1 \leq j \leq \ell$. This has running time $O(k^2)$, but

Algorithm 1: Algorithm ADJUSTSIGNIFICANCE used to compute model adjustment. Note that for notational convenience, vector indexes start at zero. Operator " $>>$ " shifts vector components to the right, padding on the left with zeros.

input : k , the size of the ranking to produce; p , the expected proportion of protected elements; α_c , the significance for each individual test.

output : The probability of rejecting a fair ranking.

```

1  $(m_{\text{old}}, i_{\text{old}}) \leftarrow (0, 0)$  // Auxiliary vectors
2 for  $i \leftarrow 1$  to  $k$  do
3    $m[i] \leftarrow F^{-1}(\alpha_c; i, p)$ 
4   if  $m[i] > m_{\text{old}}$  then
5      $m^{-1}[m_{\text{old}}] \leftarrow i$ 
6      $b[m_{\text{old}}] \leftarrow i - i_{\text{old}}$ 
7      $(m_{\text{old}}, i_{\text{old}}) \leftarrow (m[i], i)$ 
8   end
9 end
10  $S[0] \leftarrow 1$  // Success probabilities
11 for  $j \leftarrow 1$  to  $m(k)$  do
12    $S_{\text{new}} \leftarrow$  zero vector of dimension  $j$ 
13   for  $i \leftarrow 0$  to  $b(j)$  do
14     //  $f(i; b(j), p)$  is the prob. mass of  $\text{Bin}(b(j), p)$ 
15      $S_{\text{new}} \leftarrow S_{\text{new}} + (S >> i) \cdot f(i; b(j), p)$ 
16   end
17    $S_{\text{new}}[j-1] \leftarrow 0$ 
18    $S \leftarrow S_{\text{new}}$ 
19 end
20 return probability of rejecting a fair ranking:  $1 - \sum S[i]$ 

```

Table 4: Adjusted significance α_c obtained by using ADJUSTSIGNIFICANCE with $\alpha = 0.1$ for selected k, p . For small values of k, p there is no α_c that yields the required significance.

$p \backslash k$	40	100	1,000	1,500
0.1	–	–	0.0140	0.0122
0.2	–	–	0.0115	0.0101
0.3	–	0.0220	0.0103	0.0092
0.4	–	0.0222	0.0099	0.0088
0.5	0.0168	0.0207	0.0096	0.0084
0.6	0.0321	0.0209	0.0093	0.0085
0.7	0.0293	0.0216	0.0094	0.0084

we note it has to be executed only once, as it does not depend on the dataset, only on k . The summation of the probabilities in this vector S is the probability that a fair ranking is accepted when using α_c . This algorithm can be used to determine the value of α_c at which the acceptance probability becomes $1 - \alpha$, for instance, by performing binary search. This adds a logarithmic factor that depends on the desired precision. The values of α_c obtained using this procedure for selected k, p and $\alpha = 0.1$ appear on Table 4.

3.4 Utility

Our notion of utility reflects the desire to select candidates that are potentially better qualified, and to rank them as high as possible. In contrast with previous works [7, 36], we do not assume we know the contribution of having a given candidate at a particular

position, but instead base our utility calculation on losses due to non-monotonicity. The qualifications may have been even proven to be biased against a protected group, as is the case with the COMPAS scores [1] that we use in the experiments of §5, but our approach can bound the effect of that bias, because the utility maximization is subject to the ranked group fairness constraint.

Ranked utility. The ranked individual utility associated to a candidate i in a ranking τ , compares it against the least qualified candidate ranked above it.

Definition 3.3 (Ranked utility of an element). The ranked utility of an element $i \in [n]$ in ranking τ , is:

$$\text{utility}(i, \tau) = \begin{cases} \bar{q} - q_i & \text{if } \bar{q} \triangleq \min_{j: r(j, \tau) < r(i, \tau)} q_j < q_i \\ 0 & \text{otherwise} \end{cases}$$

By this definition, the maximum ranked individual utility that can be attained by an element is zero.

Selection utility. We operationalize criterion 2 of §3.1 by means of a *selection utility* objective, which we will use to prefer rankings in which the more qualified candidates are included, and the less qualified, excluded.

Definition 3.4 (Selection utility). The selection utility of a ranking $\tau \in \mathcal{T}_{k,n}$ is $\min_{i \in [n], i \notin \tau} \text{utility}(i, \tau)$.

Naturally, a “color-blind” top- k ranking $c|_k$ maximizes selection utility, i.e., has selection utility zero.

Ordering utility and in-group monotonicity. We operationalize criterion 3 of §3.1 by means of an *ordering utility* objective and an *in-group monotonicity constraint*, which we will use to prefer top- k lists in which the more qualified candidates are ranked above the less qualified ones.

Definition 3.5 (Ordering utility). The ordering utility of a ranking $\tau \in \mathcal{T}_{k,n}$ is $\min_{i \in \tau} \text{utility}(i, \tau)$.

The ordering utility of a ranking is only concerned with the candidate attaining the worst (minimum) ranked individual utility. Instead, the in-group monotonicity constraints refer to all elements, and specifies that both protected and non-protected candidates, independently, must be sorted by decreasing qualifications.

Definition 3.6 (In-group monotonicity). A ranking $\tau \in \mathcal{T}_{k,n}$ satisfies the in-group monotonicity condition if $\forall i, j$ s.t. $g_i = g_j$, $r(i, \tau) < r(j, \tau) \Rightarrow q_i \geq q_j$.

Again, the “color-blind” top- k ranking $c|_k$ maximizes ordering utility, i.e., has ordering utility zero; it also satisfies the in-group monotonicity constraint.

Connection to the individual fairness notion. Our notion of utility is centered on individuals, for instance by taking the minima instead of averaging. While other choices are possible, this has the advantage that we can trace loss of utility to specific individuals. These are the people who are ranked below a less qualified candidate, or excluded from the ranking, due to the ranked group fairness constraint. This is connected to the notion of individual fairness, which requires people to be treated consistently [10]. Under this interpretation, a consistent treatment should require that two people with the same qualifications be treated equally, and any deviation from this is in our framework a utility loss. This allows trade-offs to be made explicit.

3.5 Formal Problem Statement

The criteria we have described allow for different problem statements, depending on whether we use ranked group fairness as a constraint and maximize ranked utility, or vice versa.

PROBLEM (FAIR TOP-K RANKING). Given a set of candidates $[n]$ and parameters k, p , and α , produce a ranking $\tau \in \mathcal{T}_{k,n}$ that:

- (i) satisfies the in-group monotonicity constraint;
- (ii) satisfies ranked group fairness with parameters p and α ;
- (iii) achieves optimal selection utility subject to (i) and (ii); and
- (iv) maximizes ordering utility subject to (i), (ii), and (iii).

Related problems. Alternative problem definitions are possible with the general criteria described in §3.1. For instance, instead of maximizing selection and ordering utility, we may seek to keep the utility loss bounded, e.g., producing a ranking that satisfies in-group monotonicity and ranked group fairness, and that produces an ϵ -bounded loss with respect to ordering and/or selection utility. If the ordering does not matter, we have a FAIR TOP- k SELECTION PROBLEM, in which we just want to maximize selection utility. Conversely, if the entire set $[n]$ must be ordered, we have a FAIR RANKING PROBLEM, in which we just want to maximize ordering utility. If k is not specified, we have a FAIR SELECTION PROBLEM, which resembles a classification problem, and in which the objective might be to maximize a combination of ranked group fairness, selection utility, and ordering utility. This multi-objective problem would require a definition of how to combine the different criteria.

4 ALGORITHM

We present the FA*IR algorithm (§4.1) and prove it is correct (§4.2).

4.1 Algorithm Description

Algorithm FA*IR, presented in Algorithm 2, solves the FAIR TOP- k RANKING problem. As input, FA*IR takes the expected size k of the ranking to be returned, the qualifications q_i , indicator variables g_i indicating if element i is protected, the target minimum proportion p of protected candidates, and the adjusted significance level α_c .

First, the algorithm uses q_i to create two priority queues with up to k candidates each: P_0 for the non-protected candidates and P_1 for the protected candidates. Next (lines 5-7), the algorithm derives a ranked group fairness table similar to Table 2, i.e., for each position it computes the minimum number of protected candidates, given p, k and α_c . Then, FA*IR greedily constructs a ranking subject to candidate qualifications, and minimum protected elements required, resembling the method by Celis et al. [7] for the case of a single protected attribute (the main difference being that we compute the table m , while [7] assumes it is given). If the previously computed table m demands a protected candidate at the current position, the algorithm appends the best candidate from P_1 to the ranking (Lines 11-12); otherwise, it appends the best candidate from $P_0 \cup P_1$ (Lines 14-20).

FA*IR has running time $O(n + k \log k)$; which includes building the two $O(k)$ size priority queues from n items and processing them to obtain the ranking, where we assume $k < O(n/\log n)$. If we already have two ranked lists for both classes of elements, FA*IR can avoid the first step and obtain the top- k in $O(k \log k)$ time. Our method is applicable as long as there is a protected group and there

Algorithm 2: Algorithm FA*IR finds a ranking that maximizes utility subject to in-group monotonicity and ranked group fairness constraints. Checks for special cases (e.g., insufficient candidates of a class) are not included for clarity.

input : $k \in [n]$, the size of the list to return; $\forall i \in [n]$: q_i , the qualifications for candidate i , and g_i an indicator that is 1 iff candidate i is protected; $p \in]0, 1[$, the minimum proportion of protected elements; $\alpha_c \in]0, 1[$, the adjusted significance for each fair representation test.

output : τ satisfying the group fairness condition with parameters p , σ , and maximizing utility.

```

1  $P_0, P_1 \leftarrow$  empty priority queues with bounded capacity  $k$ 
2 for  $i \leftarrow 1$  to  $n$  do
3   | insert  $i$  with value  $q_i$  in priority queue  $P_{g_i}$ 
4 end
5 for  $i \leftarrow 1$  to  $k$  do
6   |  $m[i] \leftarrow F^{-1}(\alpha_c; i, p)$ 
7 end
8  $(t_p, t_n) \leftarrow (0, 0)$ 
9 while  $t_p + t_n < k$  do
10  | if  $t_p < m[t_p + t_n + 1]$  then
11    | // add a protected candidate
12    |  $t_p \leftarrow t_p + 1$ 
13    |  $\tau[t_p + t_n] \leftarrow \text{dequeue}(P_1)$ 
14  | else
15    | // add the best candidate available
16    | if  $q(\text{peek}(P_1)) \geq q(\text{peek}(P_0))$  then
17      |  $t_p \leftarrow t_p + 1$ 
18      |  $\tau[t_p + t_n] \leftarrow \text{dequeue}(P_1)$ 
19    | else
20      |  $t_n \leftarrow t_n + 1$ 
21      |  $\tau[t_p + t_n] \leftarrow \text{dequeue}(P_0)$ 
22    | end
23  | end
24 end
25 return  $\tau$ 

```

are enough candidates from that group; if there are k from each group, the algorithm is guaranteed to succeed, otherwise the “head” of the ranking will satisfy the ranked group fairness constraint, but the “tail” of the ranking may not.

4.2 Algorithm Correctness

By construction, a ranking τ generated by FA*IR satisfies in-group monotonicity, because protected and non-protected candidates are selected by decreasing qualifications. It also satisfies the ranked group fairness constraint, because for every prefix of size i the list, the number of protected candidates is at least $m[i]$. What we must prove is that τ achieves optimal selection utility, and that it maximizes ordering utility. This is done in the following lemmas.

LEMMA 4.1. *If a ranking satisfies the in-group monotonicity constraint, then the utility loss (ordering or selection utility different from zero) can only happen across protected/non-protected groups.*

PROOF. This comes directly from Definition 3.6 given that for two elements i, j , the only case in which $r(i, \tau) < r(j, \tau) \wedge q_i < q_j$ is when $g_i \neq g_j$. \square

LEMMA 4.2. *The optimal selection utility among rankings satisfying in-group monotonicity (i) and ranked group fairness (ii), is either zero, or is due to a non-protected candidate ranked below a less qualified protected candidate.*

PROOF. Let i, j be the two elements that attain the optimal selection utility, with $i \in \tau, j \in [n] \setminus \tau$. We will prove by contradiction: let us assume i is a non-protected element ($g_i = 0$) and j is a protected element ($g_j = 1$). By in-group monotonicity, we know i is the last non-protected element in τ . Let us swap i and j , moving i outside τ and j inside the ranking, and then moving down j if necessary to place it in the correct ordering among the protected elements below its position (given that i is the last non-protected element in τ). The new ranking continues to satisfy in-group monotonicity as well as ranked group fairness (as it has not decreased the number of protected elements at any position in the ranking), and has a larger selection utility. This is a contradiction because the selection utility was optimal. Hence, i is a protected element and j a non-protected element. \square

LEMMA 4.3. *Given two rankings ρ, τ satisfying in-group monotonicity (i), if they have the same number of protected elements $\rho_p = \tau_p$, then both rankings contain the same k elements (possibly in different order), and hence both rankings have the same selection utility.*

PROOF. Both rankings contain a prefix of size τ_p of the list of protected candidates ordered by decreasing qualifications, and a prefix of size $k - \tau_p$ of the list of non-protected candidates ordered by decreasing qualifications. Hence, $\forall i \in [n], i \in \tau \Leftrightarrow i \in \rho$, so the elements not included in the rankings are also the same elements, and the selection utility of both rankings is the same. \square

The previous lemma means selection utility is determined by the number of protected candidates in a ranking.

LEMMA 4.4. *Algorithm FA*IR achieves optimal selection utility among rankings satisfying in-group monotonicity (i) and ranked group fairness (ii).*

PROOF. Let τ be the ranking produced by FA*IR, and τ^* be the ranking achieving the optimal selection utility. We will prove that $\tau_p = \tau_p^*$ by contradiction. Suppose $\tau_p < \tau_p^*$. Then, we could take the least qualified protected element in τ_p^* and swap it with the most qualified non-protected element in $[n] \setminus \tau_p^*$, re-ordering as needed. This would increase selection utility and still satisfy the constraints, which is a contradiction with the fact that τ_p^* achieved the optimal selection utility. Suppose $\tau_p > \tau_p^*$. Then, at the position at which the least qualified protected element in τ is found, we could have placed a non-protected element with higher qualifications, as τ^* satisfies ranked group fairness and has less protected elements. This is a contradiction with the way in which FA*IR operates, as it only places a protected element with lower qualifications when needed to satisfy ranked group fairness. Hence, $\tau_p = \tau_p^*$ and by Lemma 4.3 it achieves the same selection utility. \square

LEMMA 4.5. *Algorithm FA*IR maximizes ordering utility among rankings satisfying in-group monotonicity (i), ranked group fairness (ii), and achieving optimal selection utility (iii).*

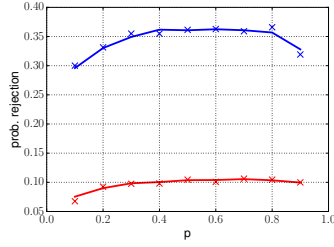


Figure 2: Probability of considering a fair ranking generated by [36] as unfair for $k = 1,000$; $\alpha_c = 0.01$ (bottom curve) and for $k = 1,500$; $\alpha_c = 0.05$ (top curve). Model represented by lines, experimental results (avg. of 10,000 runs) by crosses.

Table 5: Datasets and experimental settings.

Dataset	n	k	Quality criterion	Protected group	Protected %
D1 COMPAS [1]	18K	1K	~recidivism	Afr.-Am.	51.2%
D2 "	"	"	"	male	80.7%
D3 "	"	"	"	female	19.3%
D4 Ger. credit [27]	1K	100	credit rating	female	69.0%
D5 "	"	"	"	< 25 yr.	14.9%
D6 "	"	"	"	< 35 yr.	54.8%
D7 SAT [34]	1.6 M	1.5K	test score	female	53.1%
D8 XING [ours]	40	40	ad-hoc score	f/m/f	27/43/27%

PROOF. By lemmas 4.3 and 4.4 we know that satisfying the constraints and achieving optimal selection utility implies having a specific number of protected elements τ_p^* . Hence, we need to show that among rankings having this number of protected elements, FA*IR achieves the maximum ordering utility. By Lemma 4.1 we know that loss of ordering utility is due only to non-protected elements placed below less qualified protected elements. However, we know that in FA*IR this only happens when necessary to satisfy ranked group fairness, and having less protected elements at any given position than the ranking produced by FA*IR would violate the ranked group fairness constraint. \square

5 EXPERIMENTS

In the first part of our experiments we create synthetic datasets to demonstrate the correctness of the adjustment done by Algorithm ADJUSTSIGNIFICANCE (§5.1). In the second part, we consider several public datasets, as well as new datasets that we make public, for evaluating algorithm FA*IR (datasets in §5.2, metrics and comparison with baselines in §5.3, and results in §5.4).

5.1 Verification of Multiple Tests Adjustment

We empirically verified the adjustment formula and the ADJUSTSIGNIFICANCE method using randomly generated data. We repeatedly generated multiple rankings of different lengths k using the algorithm by Yang and Stoyanovich [36] and evaluated these rankings with our ranked group fairness test, determining the probability that this ranking, which we consider fair, was declared unfair. Example results are shown on Figure 2 for some combinations of k and α_c . As expected, the experiment results closely resemble the output of ADJUSTSIGNIFICANCE.

5.2 Datasets

Table 5 summarizes the datasets used in our experiments. Each dataset contains a set of people with demographic attributes, plus a quality attribute. For each dataset, we consider a value of k that is a small round number (e.g., 100, 1,000, or 1,500), or $k = n$ for a small dataset. For the purposes of these experiments, we considered several scenarios of protected groups. We remark that the choice of protected group is not arbitrary: it is determined completely by law or voluntary commitments; for the purpose of experimentation we test different scenarios, but in a real application there is no ambiguity about which is the protected group and what is the minimum proportion. An experiment consists of generating a ranking using FA*IR and then comparing it with baseline rankings according to the metrics introduced in the next section.

We used the two publicly-available datasets used in [36] (COMPAS [1] and German Credit [27]), plus another publicly available dataset (SAT [34]), plus a new dataset created and released with this paper (XING), as we describe next.

COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) is an assessment tool for predicting recidivism based on a questionnaire of 137 questions. It is used in several jurisdictions in the US, and has been accused of racial discrimination by producing a higher likelihood to recidivate for African Americans [1]. In our experiment, we test a scenario in which we want to create a fair ranking of the top- k people who are least likely to recidivate, who could be, for instance, considered for a pardon or reduced sentence. We observe that African Americans as well as males are given a larger recidivism score than other groups; for the purposes of this experiment we select these two categories as the protected groups.

German Credit is the Statlog German Credit Data collected by Hans Hofmann [27]. It is based on credit ratings generated by Schufa, a German private credit agency based on a set of variables for each applicant, including age, gender, marital status, among others. Schufa Score is an essential determinant for every resident in Germany when it comes to evaluating credit rating before getting a phone contract, a long-term apartment rental or almost any loan. We use the credit-worthiness as qualification, as [36], and note that women and younger applicants are given lower scores; for the purposes of these experiments, we use those groups as protected.

SAT corresponds to scores in the US Scholastic Assessment Test, a standardized test used for college admissions in the US. We generate this data using the actual distribution of SAT results from 2014, which is publicly available for 1.6 million applicants in fine-grained buckets of 10 points (out of a total of 2,400 points) [34]. The qualification attribute is set to be the achieved SAT score, and the protected group is women (female students), who scored about 25 points lower on average than men in this test.

XING (<https://www.xing.com/>) is a career-oriented website from which we automatically collected the top-40 profiles returned for 54 queries, using three for which there is a clear difference between top-10 and top-40. We used a non-personalized (not logged in) search interface and confirmed that it yields the same results from different locations. For each profile, we collected gender, list of positions held, list of education details, and the number of times each profile has been viewed in the platform, which is a measure

of popularity of the profile. With this information, we constructed an ad-hoc score: the months of work experience plus the months of education, multiplied by the number of views of the profile. This score tends to be somewhat higher for profiles in the first positions of the search results, but in general does not approximate the proprietary ordering in which profiles are shown. We include this score and its components in our anonymized data release. We use the appropriate gender for each query as the protected group.

5.3 Baselines and Metrics

For each dataset, we generate various top- k rankings with varying targets of minimum proportion of protected candidates p using FA*IR, plus two baseline rankings:

Baseline 1: Color-blind ranking. The ranking $c|_k$ that only considers the qualifications of the candidates, without considering group fairness, as described in Section 3.1.

Baseline 2: Feldman et al. [12]. This ranking method aligns the probability distribution of the protected candidates with the non-protected ones. Specifically, for a candidate i in the protected group, we replace its score $q_i \leftarrow q_j$ by choosing a candidate j in the non-protected group having $F_n(j) = F_p(i)$, with $F_p(\cdot)$ (respectively, $F_n(\cdot)$) being the quantile of a candidate among the protected (respectively, non-protected) candidates.

Utility. We report the loss in ranked utility after score normalization, in which all q_i are normalized to be within $[0, 1]$. We also report the maximum rank drop, *i.e.*, the number of positions lost by the candidate that realizes the maximum ordering utility loss.

NDCG. We report a normalized weighted summation of the quality of the elements in the ranking, $\sum_{i=1}^k w_i q(\tau_i)$, in which the weights are chosen to have a logarithmic discount in the position: $w_i = \frac{1}{\log_2(i+1)}$. This is a standard measure to evaluate search rankings [20]. This is normalized so that the maximum value is 1.0.

5.4 Results

Table 6 summarizes the results. We report on the result using p as a multiple of 0.1 close to the proportion of protected elements in each dataset. First, we observe that in general changes in utility with respect to the color-blind ranking are minor, as the utility is dominated by the top positions, which do not change dramatically. Second, FA*IR achieves higher or equal selection utility than the baseline [12] in all but one of the experimental conditions (D7). Third, FA*IR achieves higher or equal ordering utility in all conditions. This is also reflected in the rank loss of the most unfairly treated candidate included in the ranking (*i.e.*, the candidate that achieves the maximum ordering utility loss).

Interestingly, FA*IR allows to create rankings for multiple values of p , something that cannot be done directly with the baselines (Feldman et al. [12] allows what they call a “partial repair,” but through an indirect parameter determining a mixture of the original and a transformed distribution). Figure 3 shows results when varying p in dataset D4 (German credit, the protected group is people under 25 years old). This means that FA*IR allows a wide range of positive actions, for instance, offering favorable credit conditions to people with good credit rating, with a preference towards younger customers. In this case, the figure shows that we can double the

Table 6: Experimental results, highlighting in boldface the best non-color-blind result. Both FA*IR and the baseline from Feldman et al. achieve the same target proportion of protected elements in the output and the same selection unfairness, but in general FA*IR achieves it with less ordering unfairness, and with less maximum rank drop (the number of positions that the most unfairly ordered element drops).

	Method	% Prot. output	NDCG	Ordering utility loss	Rank drop	Selection utility loss
D1 (51.2%)	Color-blind	25%	1.0000	0.0000	0	0.0000
COMPAS, race=Afr.-Am.	FA*IR $p=0.5$	46%	0.9858	0.2026	319	0.1087
	Feldman et al.	51%	0.9779	0.2281	393	0.1301
D2 (80.7%)	Color-blind	73%	1.0000	0.0000	0	0.0000
COMPAS, gender=male	FA*IR $p=0.8$	77%	1.0000	0.1194	161	0.0320
	Feldman et al.	81%	0.9973	0.2090	294	0.0533
D3 (19.3%)	Color-blind	28%	1.0000	0.0000	0	0.0000
COMPAS, gender=female	FA*IR $p=0.2$	28%	0.9999	0.2239	1	0.0000
	Feldman et al.	19%	0.9972	0.3028	278	0.0533
D4 (69.0%)	Color-blind	74%	1.0000	0.0000	0	0.0000
Ger. cred, gender=female	FA*IR $p=0.7$	74%	1.0000	0.0000	0	0.0000
	Feldman et al.	69%	0.9988	0.1197	8	0.0224
D5 (14.9%)	Color-blind	9%	1.0000	0.0000	0	0.0000
Ger. cred, age < 25	FA*IR $p=0.2$	15%	0.9983	0.0436	7	0.0462
	Feldman et al.	15%	0.9952	0.1656	8	0.0462
D6 (54.8%)	Color-blind	24%	1.0000	0.0000	0	0.0000
Ger. cred, age < 35	FA*IR $p=0.6$	50%	0.9913	0.1137	30	0.0593
	Feldman et al.	55%	0.9853	0.2123	36	0.0633
D7 (53.1%)	Color-blind	49%	1.0000	0.0000	0	0.0000
SAT, gender=female	FA*IR $p=0.6$	57%	0.9996	0.0167	365	0.0083
	Feldman et al.	56%	0.9996	0.0167	241	0.0042
D8a (27.5%)	Color-blind	28%	1.0000	0.0000	0	0.0000
Economist, gender=female	FA*IR $p=0.3$	28%	1.0000	0.0000	0	0.0000
	Feldman et al.	28%	0.9935	0.6109	5	0.0000
D8b (42.5%)	Color-blind	43%	1.0000	0.0000	0	0.0000
Mkt. Analyst, gender=male	FA*IR $p=0.4$	43%	1.0000	0.0000	0	0.0000
	Feldman et al.	43%	0.9422	1.0000	5	0.0000
D8c (29.7%)	Color-blind	30%	1.0000	0.0000	0	0.0000
Copywriter, gender=female	FA*IR $p=0.3$	30%	1.0000	0.0000	0	0.0000
	Feldman et al.	30%	0.9782	0.4468	10	0.0000

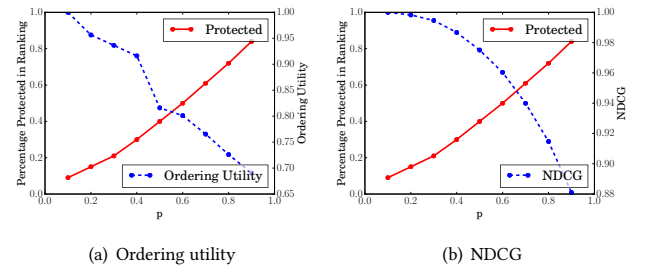


Figure 3: Depiction of possible trade-offs using FA*IR. Increase in the percentage of protected candidates in D5 (German credit, protected group age < 25) for increasing values of p , compared to ordering utility and NDCG.

proportion of young people in the top- k ranking (from the original 15% up to 30%) without introducing a large ordering utility loss and maintaining NDCG almost unchanged.

6 CONCLUSIONS

The method we have presented can generate a ranking with a guaranteed ranked group fairness, and as we have observed, does not introduce a large utility loss. Compared to the baseline of Feldman et al. [12], in general we introduce the same or less utility loss. We also do not assume that the distributions of qualifications in the protected and non-protected groups have a similar shape. More importantly, we can directly control through a parameter p the trade-off between fairness and utility.

Future work. For simplicity, we have considered a situation where people belong to a protected or a non-protected group, and leave the case of multiple protected groups or combinations of protected attributes for future work; we plan to adapt methods based on linear programming to achieve this [7]. We are also experimenting with the related problems we considered in §3.5, including directly bounding the maximum utility loss (ordering or selection), while maximizing ranked group fairness, or weighing the three criteria.

One of the main challenges is to create an in-processing ranking method instead of a post-processing one. However, we must also be cautious as results by Kleinberg et al. [23] stating that one cannot have a predictor of risk that is well calibrated and satisfies statistical parity requirements, may imply that having a fair ranking by construction is not possible. We should also consider explainable discrimination [39], or even try to show a causal relation between protected attributes and qualification scores.

Experimentally, there are several directions. For instance, we have used real datasets that exhibit some differences among protected and non-protected groups; experiments with synthetic datasets would allow to test with more extreme differences that are more rarely found in real data. Further experimental work may be done to measure robustness to noise in qualifications, and later to evaluate the impact of this in a real application.

Reproducibility. All the code and data used on this paper is available online at https://github.com/MilkaLichtblau/FA-IR_Ranking.

7 ACKNOWLEDGEMENTS

This research was supported by the German Research Foundation, Eurecat and the Catalonia Trade and Investment Agency (ACCIÓ). M.Z. and M.M. were supported by the GRF. C.C. and S.H. worked on this paper while at Eurecat. C.C., S.H., and F.B. were supported by ACCIÓ. We would also like to show our gratitude to the three anonymous reviewers for their comments and insights.

REFERENCES

- [1] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine Bias. *ProPublica* (May 2016).
- [2] Solon Barocas and Andrew D Selbst. 2014. Big data's disparate impact. *SSRN 2477899* (2014).
- [3] Francesco Bonchi, Sara Hajian, Bud Mishra, and Daniele Ramazzotti. 2017. Exposing the probabilistic causal structure of discrimination. *International Journal of Data Science and Analytics* 3, 1 (2017), 1–21.
- [4] Toon Calders and Sicco Verwer. 2010. Three naive Bayes approaches for discrimination-free classification. *Data Mining and Knowledge Discovery* 21, 2 (2010), 277–292.
- [5] Jaime Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proc. of SIGIR*. ACM Press, 335–336.
- [6] L Elisa Celis, Amit Deshpande, Tarun Kathuria, and Nisheeth K Vishnoi. 2016. How to be Fair and Diverse? *arXiv:1610.07183* (2016).
- [7] L Elisa Celis, Damian Straszak, and Nisheeth K Vishnoi. 2017. Ranking with Fairness Constraints. *arXiv:1704.06840* (2017).
- [8] Sushma Channamsetty and Michael D Ekstrand. 2017. Recommender Response to Diversity and Popularity Bias in User Profiles (short paper). In *30th International Florida Artificial Intelligence Research Society Conference*.
- [9] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. 2017. Algorithmic decision making and the cost of fairness. *arXiv:1701.08230* (2017).
- [10] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proc. of ITCS*. ACM Press, 214–226.
- [11] Evelyn Ellis and Philippa Watson. 2012. *EU anti-discrimination law*. Oxford University Press.
- [12] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In *Proc. of KDD*. ACM Press, 259–268.
- [13] Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2016. On the (im) possibility of fairness. *arXiv:1609.07236* (2016).
- [14] Batya Friedman and Helen Nissenbaum. 1996. Bias in computer systems. *ACM Transactions on Information Systems* 14, 3 (1996), 330–347.
- [15] Sara Hajian, Francesco Bonchi, and Carlos Castillo. 2016. Algorithmic Bias: From Discrimination Discovery to Fairness-aware Data Mining. In *KDD Tutorials*.
- [16] Sara Hajian and Josep Domingo-Ferrer. 2013. A methodology for direct and indirect discrimination prevention in data mining. *IEEE TKDE* 25, 7 (2013).
- [17] Sara Hajian, Josep Domingo-Ferrer, and Oriol Farràs. 2014. Generalization-based privacy preservation and discrimination prevention in data publishing and mining. *Data Mining and Knowledge Discovery* 28, 5–6 (2014), 1158–1188.
- [18] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. In *Proc. of NIPS*. Curran Associates, Inc., 3315–3323.
- [19] Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, and Aaron Roth. 2016. Fair Learning in Markovian Environments. *arXiv:1611.03071* (2016).
- [20] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems* 20, 4 (2002), 422–446.
- [21] Faisal Kamiran, Toon Calders, and Mykola Pechenizkiy. 2010. Discrimination aware decision tree learning. In *Proc. of ICDM*. IEEE CS, 869–874.
- [22] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2012. Fairness-aware classifier with prejudice remover regularizer. In *Machine Learning and Knowledge Discovery in Databases*. Springer, 35–50.
- [23] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2016. Inherent trade-offs in the fair determination of risk scores. *arXiv:1609.05807* (2016).
- [24] Juhí Kulshrestha, Muhammad B. Zafar, Motahare Eslami, Saptarshi Ghosh, Johnatan Messias, and Krishna P. Gummadi. 2017. Quantifying search bias: Investigating sources of bias for political searches in social media. In *Proc. of CSCW*.
- [25] Matevž Kunaver and Tomaž Požrl. 2017. Diversity in recommender systems—A survey. *Knowledge-Based Systems* 123 (2017), 154–162.
- [26] Natan Lerner. 2003. *Group rights and discrimination in international law*. Vol. 77. Martinus Nijhoff Publishers.
- [27] M. Lichman. 2013. UCI Machine Learning Repository. (2013).
- [28] Cathy O'Neil. 2016. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Crown Publishing Group.
- [29] Dino Pedreschi, Salvatore Ruggieri, and Franco Turini. 2009. Integrating induction and deduction for finding evidence of discrimination. In *Proc. of AI and Law*. ACM Press, 157–166.
- [30] Dino Pedreschi, Salvatore Ruggieri, and Franco Turini. 2009. Measuring Discrimination in Socially-Sensitive Decision Records. In *Proc. of SDM*. SIAM, 581–592.
- [31] Dino Pedreschi, Salvatore Ruggieri, and Franco Turini. 2008. Discrimination-aware data mining. In *Proc. of KDD*. ACM Press, 560–568.
- [32] Tetsuya Sakai and Ruihua Song. 2011. Evaluating diversified search results using per-intent graded relevance. In *Proc. of SIGIR*. ACM Press, 1043–1052.
- [33] Thomas Sowell. 2005. *Affirmative action around the world: An empirical analysis*. Yale University Press.
- [34] The College Board. 2014. SAT Percentile Ranks. (2014).
- [35] Tania Verge. 2010. Gendering representation in Spain: Opportunities and limits of gender quotas. *J. of Women, Politics & Policy* 31, 2 (2010), 166–190.
- [36] Ke Yang and Julia Stoyanovich. 2016. Measuring Fairness in Ranked Outputs. In *Proc. of FATML*.
- [37] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. 2015. Fairness constraints: A mechanism for fair classification. *arXiv:1507.05259* (2015).
- [38] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *Proc. of ICML*. 325–333.
- [39] Indre Žliobaite, Faisal Kamiran, and Toon Calders. 2011. Handling conditional discrimination. In *Proc. of ICDM*. IEEE CS, 992–1001.