

# Investigating the Impact of Gender on Rank in Resume Search Engines

**Le Chen**  
Northeastern University  
leonchen@ccs.neu.edu

**Anikó Hannák**  
Central European University  
ancsaaa@gmail.com

**Ruijun Ma**  
Rutgers University  
rma@stat.rutgers.edu

**Christo Wilson**  
Northeastern University  
cbw@ccs.neu.edu

## ABSTRACT

In this work we investigate gender-based inequalities in the context of *resume search engines*, which are tools that allow recruiters to proactively search for candidates based on keywords and filters. If these ranking algorithms take demographic features into account (directly or indirectly), they may produce rankings that disadvantage some candidates. We collect search results from Indeed, Monster, and CareerBuilder based on 35 job titles in 20 U. S. cities, resulting in data on 855K job candidates. Using statistical tests, we examine whether these search engines produce rankings that exhibit two types of *indirect discrimination: individual and group unfairness*. Furthermore, we use controlled experiments to show that these websites do not use inferred gender of candidates as explicit features in their ranking algorithms.

## ACM Classification Keywords

H.3.5 Online Information Services: Web-based services; J.4 Social and Behavioral Sciences: Sociology; K.4.2 Social Issues: Employment

## Author Keywords

information retrieval; algorithm auditing; discrimination

## INTRODUCTION

The internet is fundamentally changing the labor economy. Millions of people use services like LinkedIn, Indeed, Monster, and CareerBuilder to find employment [42, 71, 13]. These online services offer innovative mechanisms for recruiting and organizing employment, often driven by algorithmic systems that rate, sort, and recommend workers and employers.

There is potential for online labor markets to mitigate some of the mechanisms that cause discrimination in traditional labor markets. In online contexts, workers' demographics may be less clear or even anonymized, which limits the potential for cognitive biases to skew recruiting decisions. For example,

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](mailto:permissions@acm.org).

CHI 2018, April 21–26, 2018, Montreal, QC, Canada

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-5620-6/18/04...\$15.00

DOI: <https://doi.org/10.1145/3173574.3174225>

Indeed, Monster, and CareerBuilder do not ask job seekers to input their demographics, or upload a profile image.

Yet, evidence indicates that inequalities persist in many different online labor contexts. Scholars have uncovered cases of unequal opportunities presented to women in online ads [18]; biases in social feedback for gig-economy workers based on gender and race [37]; and discrimination against online customers based on socioeconomic [87]. In 2017, the Illinois Attorney General sent letters to six major hiring websites after users complained about age discrimination [63]. Although there are policies and best practices that employers may adopt to address biases in traditional hiring contexts, detecting and mitigating these issues in online, algorithmically-driven contexts remains an open challenge [4, 57].

Our goal in this work is to investigate gender-based inequalities in the ranking algorithms used by three major hiring websites: Indeed, Monster, and CareerBuilder. A common feature offered by these (and many other) hiring websites is a *resume search engine*, which allows recruiters to proactively search for candidates based on keywords and filters. Like any search engine, these tools algorithmically rank candidates, with those at the top being more likely to be seen and clicked on by recruiters [79, 17, 56]. However, if the ranking algorithm takes demographic features into account (explicitly or inadvertently through a proxy feature), it may produce rankings that systematically disadvantage some candidates. Although candidates on hiring websites rarely self-report their gender, gender can be inferred with high-accuracy from other information, such as their first name [26, 86, 64, 54, 93].

**First**, we examine *indirect discrimination*, which is defined as correlations between the output of a system and sensitive user features (e.g., gender), even if those features are not explicitly used by the system [76, 20, 95]. A ranking algorithm that exhibits indirect discrimination may cause disparate impact on the individuals being ranked. To facilitate our investigation, we ran queries on each site's resume search engine for 35 job titles across 20 American cities between May and October 2016, and recorded the search results. Our final dataset includes over 855K job candidates. Intuitively, our data corresponds to a recruiter's perspective of these sites, including the candidates' profiles and their rank in the search results.

Using this dataset, we leverage statistical tests to quantify whether the resume search engines exhibit *individual fair-*

ness (which is defined as placing candidates with similar features, excluding gender, at similar ranks) and *group fairness* (which is defined as assigning similar distributions of ranks men and women) [20, 95, 94].<sup>1</sup> We use two measures of fairness because they correspond to different assumptions about the world, and consequently have different normative consequences [27]. If we assume that candidates’ personal profiles and resumes accurately reflect their intrinsic skills, then individual fairness is the appropriate accountability standard for search engine output. However, if we assume that candidates’ data is impacted by structural inequalities in society, then the raw data is not an accurate reflection of intrinsic skills, and therefore group fairness is a more appropriate standard. We make no assumptions about how structural inequalities may impact our dataset, and thus we evaluate both types of fairness.

Our statistical tests reveal a complicated picture with respect to gender fairness on the three hiring websites:

- **Individual fairness:** By fitting mixed linear models, we find that inferred gender is a significant, negative feature on all three websites ( $p \leq 0.05$ ), indicating that feminine candidates appear at lower ranks<sup>2</sup> than masculine candidates, even when controlling for all visible candidate features. However, **the size of the gender effect is small:** on CareerBuilder, for example, at rank 30 men appear 1.4 ranks above equally qualified women on average (95% CI: [0.73, 2.13]). We demonstrate that these findings are robust by replicating them using matched subsets of candidates [39], varying subsets of the top  $k$  candidates, and candidate populations that include search filters.
- **Group fairness:** Using the Mann-Whitney  $U$  test, we find that 8.5–13.2% of job title/city pairs exhibit significant group unfairness ( $p \leq 0.05$ ). In 12 of the 35 job titles, the search results consistently favor masculine candidates; Bartender is the only job title that favors women.

**Second**, we examine *direct discrimination* on these resume search engines, which is defined as the explicit use of inferred gender as a feature when ranking candidates [76]. We performed controlled tests using resumes uploaded by us to test whether the ranking algorithms use features extracted directly from the resumes to rank candidates, including inferred gender, *almae mater*, and unemployment status.

Overall, our examination of resume search engines leads to mixed conclusions. On the positive side, we find that **the ranking algorithms used by all three hiring sites do not use candidates’ inferred gender as a feature**. Furthermore, our regressions demonstrate that the three ranking algorithms are, for the most part, individually fair with respect to gender. The small, significant gender effects that we observe are likely caused by some ranking feature that serves as a weak proxy for gender. On the negative side, however, we do observe significant and consistent group unfairness against feminine

candidates in roughly 1/3 of the job titles we examine. This may be of particular concern in technical professions like Electrical, Mechanical, Network, and Software Engineering that are known to be gender-imbalanced [72].

Whether the hiring websites should adopt ranking algorithms that strive for group fairness is a fraught question. Our analysis conclusively shows that these ranking algorithms do not “create” group unfairness with respect to gender: the algorithms are mostly “gender blind.” Rather, the algorithms are likely reflecting structural gender inequalities that are embedded in the raw data. Ultimately, we hope that our work furthers the dialog about the adoption of *algorithmic affirmative action* policies that benefit marginalized populations.

**Limitations.** There are several limitations of our work. *First*, there are no user studies that quantify how recruiters use resume search engines, including how many results they view and click on, or how they construct queries and filters. We attempt to address this in our analysis by examining gender effects under a large variety of use cases, e.g., 35 different job titles, search results lists of differing lengths, and queries with and without filters. Furthermore, it is unknown what fraction of online recruiting is active (using resume search engines) or passive (using advertisements for open positions). We leave user studies of recruiters as future work.

*Second*, the nature of our dataset restricts us to inferring binary gender labels for candidates. This is a common limitation of work that uses observational datasets to examine gender biases in online contexts [55, 37].

## RELATED WORK

Labor discrimination is a long standing, troubling aspect of society that may impact workers’ wages or opportunities for advancement. In this paper we specifically focus on *hiring discrimination*, which occurs when discrimination impacts the candidates that are selected to fill open positions. Hiring discrimination still impacts the modern job market [75], and may be based on gender, race [89, 5], sexual orientation [92], disability [29], or age [5, 25]. Unfortunately, it is one of the most difficult types of discrimination to prove in court [31].

One of the key tools used to study hiring discrimination is the *audit* or *correspondence study*. In this methodology, researchers probe the hiring practices of a target by submitting carefully crafted resumes, or by sending human participants in for interviews [78, 5, 6]. By carefully constructing the treatments to only differ by specific demographic features (e.g., gender), the researchers can measure the correlation between these variables and hiring outcomes [74].

Scholars and regulators have begun to focus on the ways that big data and algorithms can create hiring discrimination. Pauline T. Kim thoroughly catalogs how data-driven systems that evaluate job seekers may introduce new forms of discrimination against members of protected classes [57]. One potential driver of algorithmic discrimination identified by Kim and by Barocas and Selbst [4] occurs when subpopulations are not well-represented in training data. This issue is not hypothetical: in 2017, the Illinois Attorney General

<sup>1</sup>Dwork et al. originally defined these terms in the context of classification algorithms [20]. We have adapted them slightly to the context of ranking algorithms.

<sup>2</sup>In this paper, we use the terms “top” and “high” to refer to desirable ranks in search results (e.g., rank 1). This is the standard terminology used in Information Retrieval literature [17, 46].

launched an investigation against six major online job boards (including Indeed, Monster, and CareerBuilder) after receiving complaints that their design excluded older job seekers [63].

### Defining Fairness

Defining and operationalizing “fair” and “non-discriminatory” algorithms is an active area of research. Pedreshi et al. defined *direct* and *indirect discrimination*, where the former refers to algorithms that explicitly take sensitive features as input [76]. Direct discrimination is analogous to *disparate treatment* when the use of a sensitive input is legally proscribed (e.g., gender and race). An algorithm exhibits indirect discrimination if the outputs are strongly correlated with sensitive features, even if those sensitive features are not explicitly used as inputs [76]. This is analogous to *disparate impact*, and it may occur in practice when “neutral” features in a dataset act as *proxies* for sensitive features. A classic example is the use of zipcode in place of race to implement redlining [9, 34].

Dwork et al. introduced two types of fairness to address indirect discrimination: *individual fairness* states that similar individuals should be treated similarly, while *group fairness* states that demographic subsets of the population should be treated the same as the entire population [20]. Dwork et al. likened group fairness to “fair affirmative action,” as it “equalizes outcomes across protected and non-protected groups.” [20] Unfortunately, it is not always possible to achieve both types of fairness simultaneously if the base rates of one or more critical features are different across subpopulations. Thus, most existing fair classification systems optimize for individual [50, 53, 66, 52, 20] or group fairness [48, 9, 49, 11, 96, 51, 10, 34, 24]. Zemel et al. present classifiers that allow the operator to tune the tradeoff between individual and group fairness [95].

Friedler et al. present a framework for grappling with the assumptions that underly individual and group fairness in algorithmic scenarios [27]. If we assume that a dataset accurately encodes the intrinsic characteristics of a population, then it is appropriate to judge fairness at the individual-level. However, if we assume that a dataset is influenced by systematic or structural biases, then it is no longer individually-accurate, and we should instead pursue group fairness. In this work, we evaluate whether hiring websites exhibit individual and group fairness, since we do not make assumptions about whether the data on the websites is impacted by structural bias.

### Biases in Online Systems

Researchers have begun to empirically investigate whether algorithmic systems may (inadvertently) cause harm to users. The process of examining a black-box computer system has become known as an *algorithm audit*, as the methodology draws inspiration from classic audit studies [80]. Prior algorithm audits have examined search engines [35, 58], online maps [84], social networks [23], e-commerce [69, 70, 36, 15], and online advertising [33, 61]. To our knowledge, ours is the first audit to focus on online job boards.

There are many causes for bias in online systems. Some cases are direct manifestations of societal biases by users, e.g., gender biases on Wikipedia [60, 77, 91]. In other cases, biases are “learned” by an algorithm that is trained on biased data, e.g.,

racist ad targeting on Google Search [85], or sexist word associations in language models [12]. Finally, bias may arise due to a combination of user-driven, structural bias, and algorithm design [47]. In this work, we examine whether volunteered information from job seekers and/or algorithm design decisions cause hiring websites to produce biased search results.

Closely related to our work are studies that have uncovered discrimination on “gig-” and “sharing-economy” services. Studies have found examples of workers/service providers discriminating against customers on TaskRabbit [87], AirBNB [21], and Uber [28], as well as cases where customers discriminate against workers on TaskRabbit and Fiverr [37].

**The Importance of Rank.** In this study, we examine three websites that present job seekers in ranked lists in response to queries from recruiters. If these ranking algorithms systematically elevate candidates with specific demographic attributes, this may recreate real-world social inequality in an online context because the top items of ranked lists are much more likely to be clicked by users [30, 32, 79, 17]. For example, Keane et al. demonstrated that even if users are unknowingly presented with an inverted list of Google Search results, they still overwhelmingly click on the top-ranked items [56].

**Search Engines.** Our work falls within the literature that examines social harms that can occur when search engines present misleading or biased information. Researchers have examined misinformation about vaccines [2, 65], biased scientific information [73], and climate change denial [90]. Epstein et al. use controlled experiments to show that biased political information presented by a search engine can significantly alter users’ voting patterns [22].

Three studies have specifically examined demographic biases on search engines. Hannák et al. found negative correlations between race and rank on the gig-economy website TaskRabbit, even after controlling for all other worker-related features [37]. Kay et al. and Magno et al. both examined gender on Google Image Search, and found that women are often depicted stereotypically [55, 67]. Kay et al.’s work is particularly relevant, since they examined gendered images after querying for occupations. They found that users preferred images containing people that matched an occupation’s gender stereotype (e.g., a male CEO), and that over-representation of a gender in the search results shifted a user’s perception of gender balance in that occupation. If these findings apply to recruiters, then this suggests that hiring websites should strive for group fairness in search results, as this would improve the perceived gender balance in occupations, as well as encourage recruiters to engage with candidates that do not match an occupation’s gender stereotype (e.g., female CEOs).

Two studies have proposed metrics for quantifying bias in search results. Kulshrestha et al. separately quantify the amount of bias in a search engine’s corpus and output [59]. Unfortunately, we cannot use this metric because we do not have the entire corpus of candidates from the hiring websites. Yang et al. define a family of bias metrics that are related to normalized Discounted Cumulative Gain (nDCG) [94], which is a common Information Retrieval metric that we use in this

Limits	Indeed	Monster	CB
Number of candidates shown per page	50	20	20-30
Maximum number of candidates per query	5000	1000	5000
Maximum resume views per month	No limit	100	50
Cost per month	Free	\$700	\$400

Table 1: Search result format and limits for all three sites.

paper. However, the Yang et al. metrics are difficult to use in practice, since they rely on a normalization term that is computed stochastically. In this paper, we rely on standard statistical tests, since they are easier to interpret, provide confidence guarantees, and have been used successfully by prior work to examine inequalities in algorithmic systems [81].

## BACKGROUND

In this section, we introduce the three websites that are the focus of our study. We discuss how recruiters use these *resume search engines*, and specific details about their user interfaces. We also briefly discuss how candidates use these websites.

### Hiring Websites

We chose three hiring sites to examine in this study: Indeed, Monster, and CareerBuilder. We chose these sites for three reasons. *First*, they are three of the most popular employment websites in the U. S. (along with LinkedIn and Glassdoor) [1]. Each of these websites claims to serve millions of unique visitors and job queries per month [42, 71, 13]. *Second*, all three sites provide economical access to their resume search engine, as shown in Table 1. Contrast this to LinkedIn, which charges \$9000 for access to its unrestricted recruiter tools. *Third*, as we discuss next, all three sites have similar user interfaces for candidates and recruiters. This makes the sites roughly comparable in terms of usability and features, which allows us to contrast our results across the sites.

### Recruiter’s Perspective

In this study, we examine the *resume search engines* provided by Indeed, Monster, and CareerBuilder. All three sites offer similar search engines that are designed to help recruiters identify and recruit new employees. The *corpus* of each resume search engine contains resumes and personal profiles uploaded by candidates seeking employment. Recruiters *query* the corpus by entering a free-text job title, a geographic location, and (optional) *filters* to refine the results (e.g., years of experience, minimum educational attainment, *etc.*). **None of the sites allow recruiters to filter or order search results by demographics** (e.g., gender, ethnicity), but proxy variables exist in some cases (e.g., years of experience as an indicator of age).

The resume search engines use rankings algorithms to determine which candidates are relevant to a given query and their ordering, subject to any specified filters. By default, only candidates within 20–50 miles of the specified location are deemed relevant, and the search results are sorted by opaque metrics (e.g., “most relevant”), although the recruiter may re-sort the list by objective metrics like years of experience.

Table 1 lists details about the search result format and data access restrictions on the three sites. Candidate features we

U. S. State	Cities
Texas	Austin
Iowa	Des Moines
Wisconsin	Madison, Milwaukee
Louisiana	New Orleans
New York	New York City, Buffalo
Nebraska	Omaha
Utah	Salt Lake City
California	San Francisco, Stockton, San Bernardino, Los Angeles
Missouri	Springfield
Michigan	Detroit
Ohio	Toledo, Cleveland, Cincinnati
Tennessee	Memphis
Illinois	Chicago

Table 3: Cities used in our queries.

can observe in search results for each site, besides candidates’ names and current job titles, are given in Table 2.

**Scope.** In this work, we focus on the ranking of candidates, rather than the composition of search results (i.e., which candidates are deemed relevant). Although examination of the overall candidate pool would be interesting, we cannot do so because there is no way for us to enumerate all candidates in a hiring website’s corpus.

### Candidate’s Perspective

Indeed, Monster, and CareerBuilder are very similar from a candidate’s perspective. Candidates must register for a free account, and possibly fill out a personal profile and upload a resume. The amount of profile information that is mandatory varies across the sites; on Monster, users must provide their name, location, educational attainment, and previous job, while CareerBuilder allows users to leave their profile empty. However, all three sites remind users to upload more information, especially a resume, since the job recommendation functionality on the sites depends on this information. Once a candidate has created a profile, they can browse open positions and respond to solicitations from recruiters.

## DATA COLLECTION

In this section, we describe our data collection methodology, and the specific variables we extract from the data.

**Crawl.** To collect data for this study, we use an automated web browser to search for candidates on Indeed, Monster, and CareerBuilder. Intuitively, our crawler is designed to emulate how a recruiter would search for candidates on these hiring websites. We ran queries for 35 job titles in 20 U. S. cities (described below) on all three sites, and recorded the resulting lists of candidates. We also queried for a subset of 490, 700, and 700 job title/city pairs with one, two, and three search filters, on Monster, Indeed, and CareerBuilder, respectively. For binary filters (e.g., willingness to relocate), we queried with both options; for non-binary filters (e.g., minimum years of experience), we set three different values for the filter, chosen such that the values select from 0–33%, 33–66%, and 66–100% of the candidate population.<sup>3</sup> On Indeed we recorded candidates’ resumes, but not on Monster or CareerBuilder since they only allow recruiters to view  $\leq 100$  resumes

<sup>3</sup>We calculate these filter values based on the Cumulative Density Function of the corresponding features from the unfiltered datasets.

Origin	Feature	Description	Present On		
			Indeed	Monster	CB
Observed in Search Results	Job Title Relevance	Relevance of the searched job title to the candidate's current job title	✓	✓	✓
	Skills Relevance	Relevance of the searched job title to the candidate's skills	✓	✓	✗
	Education Level	Education level of the candidate	✓	✓	✓
	Job Popularity	Popularity of the candidate's current job title among all candidates returned for the searched job title	✓	✓	✓
	Last Modified	The recency of the candidate's profile and resume	✓	✓	✓
	Experience	The experience of the candidate in years	✓	✓	✓
	Relocate	Whether the candidate is willing to relocate (binary)	✗	✓	✗
	Skills Popularity	Popularity of the candidate's skills among all candidates returned for the searched job title	✗	✓	✗
	Information Relevance	Relevance of the searched job title to additional profile info	✓	✗	✗
	Bio Relevance	Relevance of the searched job title to the candidate provided description of the working experience	✓	✗	✗
	Skills Match	Whether the entire searched job title is present in the candidate's skill set (binary)	✓	✗	✗
	Information Match	Whether the entire searched job title is present in the candidate's additional profile information	✓	✗	✗
	Bio Match	Whether the entire searched job title is present in the candidate's bio	✓	✗	✗
Inferred	Gender	Probability of the candidate being masculine	✓	✓	✓

Table 2: Per-candidate features we extract from search results on Indeed, Monster, and CareerBuilder. We infer *Gender* from each candidate's first name; other features are directly observed. Not all features are present on all hiring websites.

per month. All of our data was collected between May and October 2016, and the crawling took two months on each site.

To obtain a broad sample of candidates, we ran queries for 35 job titles in 20 cities. Table 6 lists our 35 job titles; 19 were chosen because they are the most commonly searched job titles [14], while the remaining 16 do not require high-school-level education [8], which adds diversity to our queries. Table 3 shows the 20 cities we focus on, which were chosen to give us broad geographic and demographic variety. We ultimately gathered 521,783, 265,172, and 67,580 candidates on Indeed, Monster, and CareerBuilder after running our queries.

**Candidate Features.** Next, we extract information about candidates in the search results. We focus on three types of features: 1) profile data (e.g., experience, education, *etc.*); 2) inferred gender; and 3) rank in search results. Table 2 lists the features we are able to extract on each site. We normalize all features in Table 2 to be between 0 and 1 for consistency. Details of how we encode and normalize each feature can be found in the supplementary materials.

**Inferring Gender.** Since our goal is to examine resume search engines with respect to gender, we need to label each candidate's gender. However, none of the sites we focus on collect this information. Instead, we infer each candidate's gender based on their first name, which is a common method to infer gender in Western societies [26, 86, 64, 54, 93].

In this work, we rely on the U. S. baby name dataset [83] to infer candidates' gender.<sup>4</sup> We assign a *probability of being masculine* to each candidate based on the fraction of times their first name was given to a male baby in the name dataset. We represent gender as a probability since this corresponds to how recruiters perceive candidates' genders. For example, a candidate named "John" is almost certainly masculine,

<sup>4</sup>We also tried inferring candidates' gender using the Genni and SexMachine datasets [82, 88]. However, this resulted in 13% and 19% candidates having unknown gender, respectively, versus 8% for our method. Additionally, we fit mixed linear models to the Genni and SexMachine labeled data, and found that the *Probability of Being Masculine* was significant ( $p \leq 0.001$ ) and negative on all three hiring websites.

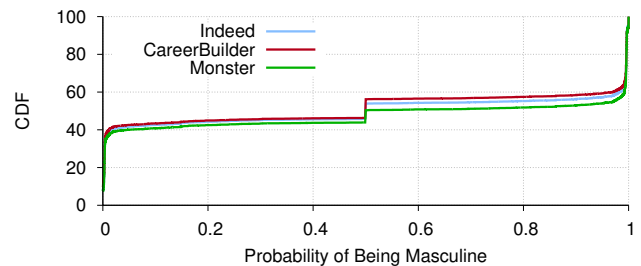


Figure 1: Inferred probability of being masculine for all candidates in our dataset.

while "Madison" is ambiguous. We assign a score of 0.5 to candidates whose names do not appear in the dataset.

Figure 1 shows the CDF of the probability of being masculine for all candidates on our dataset. Only 8% of candidates have ambiguous gender, all but 45 of which correspond to candidates whose names do not appear in the dataset. This means that we can be very confident in the vast majority of our gender labels. The plot also shows that the masculine/feminine ratio is approximately 1:1 on all three sites.

**Limitations.** There are two limitations of our dataset and labeling methodology that are worth noting. *First*, the candidate attributes that we extract from search results may not match a candidates' true attributes. Fortunately, this limitation does not impact our analysis, since the ranking algorithms we are auditing, as well as recruiters that rely on resume search engines, base their decisions on candidates' reported attributes.<sup>5</sup> Thus, throughout this paper, when we compare the capabilities of candidates, we are referring to their reported attributes, rather than their true attributes.

*Second*, we do not know candidates' true genders. As above, this limitation does not impact our analysis, since recruiters and ranking algorithms must also rely on inferred gender when making decisions (if they use this information at all). Throughout this paper when we refer to gender, we are referring to

<sup>5</sup>Recruiters do not learn a candidate's true attributes until much later in the hiring process, e.g., after interviewing the candidate and checking their references.

inferred gender based on first name. Furthermore, as noted in the introduction, we are limited to analyzing binary genders.

**Ethics.** We were careful to conduct our study in an ethical manner. This study was approved under Northeastern IRB #16-01-19. To protect the contextual privacy of candidates, we will not be releasing our crawled data. Furthermore, we limited the impact of our crawler on the hiring sites by restricting it to one query every 30 seconds, and at no point did we contact candidates. In our controlled experiments (detailed in the next section), we only uploaded two resumes at a time to the hiring sites, meaning we decreased the rank of other candidates by at most two. Although all three sites prohibit crawling in their terms of service, we believe our study is acceptable under the norms that have been established by prior algorithm audit studies [80, 37].

## ANALYSIS

In this section, we investigate the rankings of candidates produced by resume search engines with respect to inferred gender. We organize our analysis around three questions: (1) Do the resume search engines exhibit individual fairness? (2) Do they exhibit group fairness? (3) Do they explicitly rank candidates based on inferred gender?

### Individual Fairness

In this section, we investigate whether the rankings of candidates produced by resume search engines are individually fair with respect to inferred gender. Recall that to be individually fair, the ranking algorithms must rank candidates with similar features (excluding gender) at similar ranks [20, 95, 94].

To investigate individual fairness, we use regression tests. Our goal is to examine the effect of inferred gender on candidate’s rank while controlling for all other observable candidate features. If the gender feature is significant and has a non-zero coefficient in the fitted model, this indicates that the ranking algorithm in question is not individually fair, as candidates with equivalent features but different inferred genders are not assigned the same rank.

Throughout this section, we focus on the top 100 candidates returned in search results, since recruiters are unlikely to browse to candidates at lower ranks [79, 17]. However, to ensure the robustness of our results and avoid making assumptions about the behavior of recruiters, we fit additional models to many different subsets of candidates, which we describe below.

**Model Specification.** We adopt a Mixed Linear Model (MLM) for our regressions. We regress on individual candidates, specifying the model as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$ .  $\mathbf{y}$  is a vector of the responses ( $\log_2(\text{rank})$ ) of each candidate, explained below).  $\mathbf{X}$  is the design matrix, and the predictors include the features from Table 2.  $\boldsymbol{\beta}$  is the vector of fixed-effects parameters, including the global intercept.  $\boldsymbol{\mu}$  is the vector of random-effects intercepts. In other words, search results for a specific job title and city have a shared global fixed-effects intercept and an individual random-effects intercept.

We choose a model with fixed and random effects because it agrees with two fundamental assumptions about our data:

1. Each hiring website has a single ranking algorithm with features and weights do not vary by query term or location. This corresponds to fixed effects ( $\boldsymbol{\beta}$ ) regardless of job title and location. This assumption is reasonable because it is impractical for a hiring website to implement different algorithms or feature weight vectors for an unbounded number of free-text job titles and locations.
2. Candidates with identical features that appear in different search result lists may be assigned dramatically different ranks. This is true because the population of candidates, and their relative qualifications, varies across locations and professions. This corresponds to random-effects group intercepts ( $\boldsymbol{\mu}$ ) that vary by job title and location.

We use  $\log_2(\text{rank})$  as the dependent variable in our regressions for two reasons:

- It prioritizes top-ranked candidates by giving them higher weight, while decaying the importance for lower-ranked candidates.
- The  $\log()$  function is monotonically increasing, which maintains the ordering of candidates after the transformation.

Logarithms have been found to be widely applicable in the IR literature. Empirical studies [79, 17], backed up by eye-tracking surveys [30, 32, 68], have found that the probability of search result items being clicked decays logarithmically. Similarly,  $\log_2()$  is used to calculate normalized Discount Cumulative Gain (nDCG), which is a standard metric used to quantify the similarity of search result lists [43, 44]. By using logarithmic decay, nDCG affords higher weight to the important items at the top of the search result lists.

**Model Fitting.** We fit three MLMs, one for each hiring website on the top 100 candidates in each search result list. Negative coefficients signify effects with higher rank. We conduct the Variance Inflation Factors (VIF) test to remove multicollinearity before fitting the models; all the variables remain after the test. The correlation matrix is available in the supplementary materials.

To assess how well the MLMs fit for our data, we evaluate their predictive power. To do this, we treat the each model as a ranking algorithm: we input the ground-truth feature vectors of candidates from a given search result list  $\mathbf{R}$  into a fit model, which then outputs a predicted log ranking  $\hat{\mathbf{y}}$  (corresponding to a predicted ranking  $\hat{\mathbf{R}}$ ). Next, we use the nDCG metric to compute the similarity between  $\hat{\mathbf{R}}$  and the original ranking  $\mathbf{R}$ . nDCG is a standard metric used in the IR literature to compare ranked lists [43, 44]. The DCG of a ranking  $\mathbf{R} = [r_1, r_2, \dots, r_k]$  is calculated as

$$\text{DCG}(\mathbf{R}) = g(r_1) + \sum_{i=2}^k (g(r_i) / \log_2(i))$$

where  $g(r_i)$  is the “gain” or score assigned to result  $r_i$ . nDCG is defined as  $\text{DCG}(\hat{\mathbf{R}}) / \text{DCG}(\mathbf{R})$ , where  $\mathbf{R}$  is the ideal ranking of the items in  $\hat{\mathbf{R}}$ . In our case,  $\mathbf{R}$  is the original ranking produced by a hiring site (i.e., we treat the original ranking as the baseline), and  $\hat{\mathbf{R}}$  is a ranking predicted by the model. An

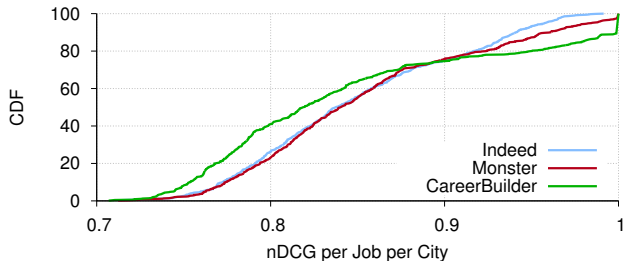


Figure 2: nDCG comparison of the predicted rankings  $\hat{R}$  produced by our MLMs versus the original rankings  $R$ .

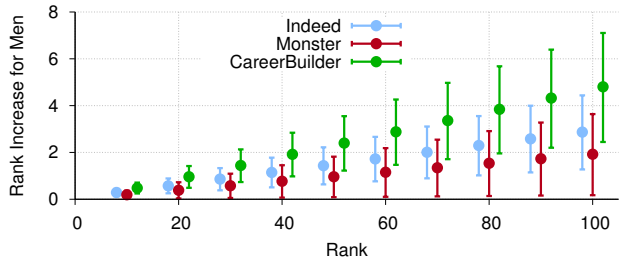


Figure 3: Effect sizes with 95% confidence intervals for the gender coefficients. Points have been jittered horizontally to prevent overlap.

nDCG value of 1 indicates that the two rankings are identical. In essence, nDCG for rank responses is analogous to  $R^2$  for continuous responses.

Figure 2 shows the evaluation of the predictive power of our MLMs. We observe that 60–77% of the nDCG scores are  $\geq 0.8$  across all three websites. To put this in perspective, state-of-the-art learning-to-rank algorithms produce nDCG scores in the 0.4–0.8 range, depending on the context and the benchmark dataset that is used [40, 19, 45]. This demonstrates that our MLMs reproduce most of the search results with high accuracy, and that the models are a good fit for our data.

**Results.** Table 4 shows the results of our MLM regressions on the top 100 candidates in search results. We observe that the majority of features have significant, negative effect on log rank on all three hiring sites, such as *Job Title Relevance* and *Job Popularity*. On Indeed, *Last Modified* has the largest coefficient by far, which matches their documentation stating that they tend to rank candidates by resume update time [41]. Interestingly, *Education Level* and *Experience* have significant, positive effects on log rank on Indeed, which suggest that there are more candidates with lower degrees and less experience in the top ranks (possibly newly graduated students). In contrast, Monster and CareerBuilder both exhibit significant, negative effects of *Experience* on log rank.

Lastly, we observe that the *Probability of Being Masculine* feature is significant ( $p < 0.05$  in all cases) and negative in all three models, meaning that overall, men rank higher than women with equivalent features.

Figure 3 shows the effect sizes and 95% confidence intervals for the gender coefficient in our models. The effect sizes

Feature	Dependent Variable: $\log_2(\text{rank})$		
	Indeed	Monster	CB
Fixed Effect Intercept	4.803***	5.938***	4.129***
Job Title Relevance	-0.518***	-1.3***	-1.258***
Skills Relevance (1)	-0.051	-0.31***	
Skills Relevance (2)		-0.109***	
Skills Relevance (3)		-0.108***	
Education Level	0.042**	-0.061*	-0.027
Job Popularity	-0.115***	-0.004	-0.147***
Last Modified	-2.053***	-0.197***	-0.149***
Experience	0.116***	-1.303***	-0.185***
Relocate		-0.021	
Skills Popularity (1)		-0.048**	
Skills Popularity (2)		-0.062**	
Skills Popularity (3)		-0.017	
Bio Relevance	0.041		
Information Relevance	-0.255***		
Skills Match	-0.072		
Information Match	-0.093*		
Bio Match	-0.262***		
Random Effect (s.d.)	0.01	0.106	0.018
Prob. of Being Masculine	-0.042***	-0.028*	-0.071***
Observations	67410	50813	28289

Table 4: Estimated coefficients and standard deviation of mixed linear regressions on the top 100 candidates in search results from each hiring website, grouped by city and job title. Significance level is unavailable for *Random Effect*. Note: \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$ .

Top k Candidates	Indeed Gender Coef.	Monster Gender Coef.	CareerBuilder Gender Coef.
Top 10	-0.023	-0.031	0.027
Top 20	-0.009	-0.056*	-0.023
Top 50	-0.036*	-0.047**	-0.053*
Top 100	-0.042***	-0.028*	-0.071***
Top 200	-0.029**	-0.026**	-0.071***
Top 500	-0.022***	-0.040***	-0.055***
Top 1000	-0.019***	-0.043***	-0.039**

Table 5: Estimated *Probability of Being Masculine* coefficient from mixed linear regressions as the length of the search result list  $k$  is varied. Note: \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$ .

increase as rank decreases due to our  $\log_2()$  transformation. For very high ranks, the increase in rank for men is negligible: for example, on Monster at rank 10 the mean increase in rank is 0.2 (95% CI: [0.02, 0.36]). However, by rank 50 on Monster, ties between men and women are consistently broken in favor of men (mean increase 0.96, 95% CI: [0.08, 1.82]), and by rank 80 the increase is large enough to push men across pagination boundaries (mean increase 1.54, 95% CI: [0.14, 2.91]). Indeed and CareerBuilder both exhibit larger effects than Monster.

**Robustness.** In addition to the MLMs we fit to the top 100 candidates in our dataset, we fit hundreds of other MLMs to sub- and supersets of our candidate population to gauge the robustness of our models. Specifically, we fit models to: (1) the top  $k$  candidates in search results as  $k$  is varied from 10 to 1000 (chosen because it is the maximum number of candidates returned by Monster); (2) the top 100 and 1000 candidates in a *matched* subset of the population (propensity score matching is a technique for reducing selection bias in observational datasets [38]); (3) the top 100 candidates in search results that include combinations of one, two, and three filters (e.g., minimum experience).

Job Title	Indeed		Monster		CB	
	%	TDRC	%	TDRC	%	TDRC
Accountant	20	-0.05	15	0.02	7.5	-0.01
Auditor	5	0.06	0	0.05	0	0.26
Bartender	15	0.00	5	0.00	11.7	-0.02
Business Dev. Manager	0	0.07	0.6	0.04	0	0.06
Call Center Director	25	0.05	-	-	-	-
Cashier	<b>10</b>	<b>0.01</b>	<b>5</b>	<b>0.02</b>	<b>37.1</b>	<b>0.09</b>
Casino Manager	5.5	0.03	0	0.06	-	-
Concierge	0.6	0.04	4.1	-0.01	0	0.03
Corrections Officer	<b>20</b>	<b>0.15</b>	<b>0</b>	<b>0.19</b>	-	-
Customer Service	<b>5</b>	<b>0.03</b>	<b>20</b>	<b>0.02</b>	<b>8.3</b>	<b>0.17</b>
Electrical Engineer	<b>22.8</b>	<b>0.53</b>	<b>20</b>	<b>0.13</b>	-	-
Elevator Technician	0	0.01	-	-	-	-
Financial Analyst	6.1	0.16	6.1	-0.02	36.7	0.00
Human Res. Specialist	<b>10</b>	<b>0.03</b>	<b>0</b>	<b>0.15</b>	<b>70</b>	<b>0.11</b>
Janitor	<b>5</b>	<b>0.35</b>	<b>5</b>	<b>0.17</b>	<b>0</b>	<b>0.02</b>
Laborer	<b>20</b>	<b>0.36</b>	<b>40</b>	<b>0.36</b>	<b>0</b>	<b>0.19</b>
Mail Carrier	<b>10</b>	<b>0.16</b>	<b>9.3</b>	<b>0.22</b>	-	-
Manufacturing Engineer	25	0.03	0	-0.10	-	-
Marketing Manager	5.5	-0.01	0	0.08	20	0.00
Mechanical Engineer	<b>16.1</b>	<b>0.01</b>	<b>0</b>	<b>0.29</b>	-	-
Network Engineer	<b>47.6</b>	<b>0.48</b>	<b>0</b>	<b>0.27</b>	<b>0</b>	<b>0.00</b>
Occupational Therapist	7.5	0.04	0	-0.01	-	-
Payroll Specialist	5	0.04	0	0.00	0	0.00
Personal Trainer	5.5	0.01	0	0.21	0	0.04
Pharmacist	35	0.08	0	0.09	0	-0.03
Physical Therapist	15	0.02	11.7	-0.04	-	-
Real Estate Agent	5	0.06	0.3	0.15	0	0.00
Registered Nurse	5	0.03	0.9	0.01	0	0.05
Retail Sales	5	-0.09	55	0.05	25.8	0.00
Speech Pathologist	28.3	0.01	-	-	-	-
Software Engineer	<b>25</b>	<b>0.59</b>	<b>30</b>	<b>0.42</b>	<b>11.7</b>	<b>0.04</b>
Tax Manager	5	0.07	0	0.02	0	0.00
Taxi Driver	28.3	-0.04	0	0.02	-	-
Technical Recruiter	20	-0.15	15	0.12	0	0.00
Truck Driver	<b>15</b>	<b>0.50</b>	<b>9.3</b>	<b>0.49</b>	<b>45</b>	<b>0.00</b>
<b>Overall Percentage/Total</b>	<b>13.2</b>	<b>3.68</b>	<b>8.5</b>	<b>3.48</b>	<b>13.2</b>	<b>1.45</b>

Table 6: Evaluation of group fairness. The “%” columns show the percentage of cities with  $p < 0.05$  in the Mann-Whitney  $U$  test on the rank of the top 1000 women and men for a given job title and hiring website. We subtract 5% from each percentage (with a floor of 0) to account for multiple hypothesis testing. The “TDRC” columns show the total difference in area under the recall curves. Negative (positive) TDRC indicates that feminine (masculine) candidates are ranked higher overall. “-” marks instances where there are no cities with sufficiently large populations to test. We also remove the 8% of candidates with ambiguous genders. Job titles with significant unfairness in a constant direction are **bolded**.

Overall, these models exhibit the same significance, sign, and effect sizes for the *Probability of Being Masculine* feature as the top 100, unmatched, non-filtered models that we examine in Table 4. For example, the gender coefficients from our top  $k$  models are shown in Table 5; we observe that once the sample sizes are sufficiently large ( $k \geq 50$ ), the gender coefficients become uniformly significant and negative. Taken together, these models demonstrate that our results are robust to under a wide variety of conditions. Detailed results for our top 1000 unmatched and matched models are available in the supplementary materials.

### Group Fairness

Next, we investigate whether the resume search engines exhibit group fairness with respect to inferred gender. To be group fair, the ranking algorithms should assign a similar distribution of ranks to masculine and feminine candidates [20, 95, 94].

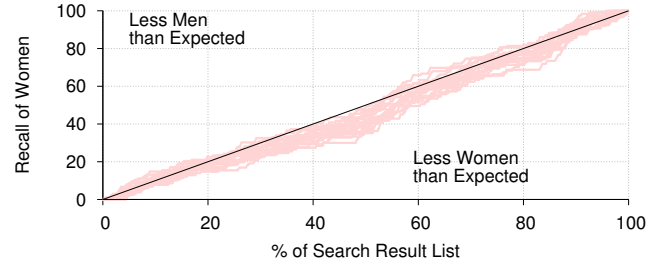


Figure 4: Recall of women when searching for a “Software Engineer” on Indeed. The 20 lines correspond to different cities.

**Metrics.** To examine group fairness, we use two metrics. First, we use the Mann-Whitney  $U$  (M-W  $U$ ) test to compare the distribution of ranks for men and women in a given list of search results [16]. The M-W  $U$  test is a nonparametric test of the null hypothesis that there is no tendency for ranks of one class to be significantly different from the other. We omit search result lists where the number of feminine or masculine candidates is less than 20, as the result of M-W  $U$  are not reliable when samples are this small. Out of  $35 \times 20 = 700$  samples on each hiring website, 648, 421, and 181 are suitable for analysis on Indeed, Monster, and CareerBuilder, respectively. We also remove the 8% of candidates with uncertain gender (*Probability of Being Masculine*  $\geq 0.2$  and  $\leq 0.8$ ).

Table 6 shows the results of the M-W  $U$  tests on our search result samples. Each cell in the “%” columns shows the percentage of samples across cities where the M-W  $U$  test was significant at  $p \leq 0.05$  for a given job title and hiring website. We adjust each percentage downward by 5% to correct for multiple hypothesis testing, i.e., we assume that all tests are independent, and that there is a uniform 5% false positive rate. This correction approach is more appropriate for our scenario than Bonferroni correction, since we are interested in the overall amount of positive tests, not the specific cities that exhibit significant differences.

The disadvantage of M-W  $U$  is that it does not tell us the direction or magnitude of group unfairness. To answer these questions, we calculate the area under *recall curves*. To generate recall, we iterate from the first candidate (i.e., rank 1) to the last candidate in a given search result list  $R$ . At rank  $i$ , we calculate a tuple  $(x_i, y_i) = (i/|R|, |R_{f,i}|/|R_f|)$ , where  $|R|$  is the total number of candidates in the list,  $|R_f|$  is the total number of *feminine* candidates in  $R$ , and  $|R_{f,i}|$  is the number of feminine candidates observed between ranks 1 and  $i$ .

As an illustrative example, Figure 4 shows the recall for candidates on Indeed when we search for “Software Engineer.” Each red line corresponds to the search results in a different city. If women and men are evenly distributed in the list, the resulting line is along the diagonal. In this case, we see that most of the lines are below the diagonal, meaning that women are under-represented in the search results relative to their overall percentage of the candidate population.

Finally, we calculate the area under the recall curves and subtract the area under the diagonal; the Difference between



Feature	Indeed	Monster	CB
Inferred Gender	×	×	×
School Ranking	×	✓	×
Employment	×	✓	✓
Number of Keywords	×	×	×
Resume Length	×	×	×
Job Churn	×	×	✓
Contact Information	×	×	×
Company Name	×	×	×

Table 7: Features tested in our controlled resume experiments. Green check marks denote cases where the feature is taken into account by the ranking algorithm.

the Recall Curves (DRC) exists in the range  $[-0.5, 0.5]$ , with negative (positive) values indicating that feminine (masculine) candidates are favored in the rankings. Table 6 shows the Total DRC (TDRC) for each job title summed across cities.

**Results.** Table 6 shows that overall 8.5–13.2% of job title/city pairs exhibit significant group unfairness, and that the aggregate directionality favors masculine candidates. In some of the significant job titles (e.g., Accountant, Physical Therapist, Retail Sales, and Technical Recruiter), the direction of unfairness is inconsistent, suggesting that the unfairness may be due to natural variations in the candidate populations. However, in 12 of the significant job titles (highlighted in **bold**), the direction of unfairness is consistently in favor of men. The magnitude of favor towards men in some of these occupations (e.g., Network Engineer, Software Engineer, and Truck Driver) is very large, often  $\sim 0.5$ . The consistent directionality of unfairness in these job titles suggests that the underlying cause is structural. The only job title that is significant and approaches uniform unfairness in favor of women is Bartender, but the magnitude of unfairness is relatively small.

### Direct Discrimination And Hidden Features

Up to now, our analysis has focused on candidate features that are directly observable in the search results. However, there is an element of each candidates’ profile that we cannot observe (on Monster and CareerBuilder), but that may be taken into account by the ranking algorithm: resumes. For example, Monster and CareerBuilder ask candidates to enter their education level into their profile, but actual almae matres are likely stated in each candidate’s resume. The ranking algorithm could parse this additional information from the PDF-format resume and use it when ranking candidates. Parsing resumes makes sense from a design standpoint, in that it allows the websites to collect detailed information about candidates without having to ask them to enter it manually, which can be tedious.

To test if resume content influences ranking, we conducted controlled experiments using resumes created and uploaded by us. We create two user accounts,  $A$  and  $B$ , in that temporal order, with identical profile information and resumes. We then verify that the two accounts appear directly adjacent in search results in the order  $B, A$ . Next, we delete the old resumes, upload two new resumes (starting with  $A$ ) that differ by exactly one feature, then query for our users again.<sup>6</sup> In each

<sup>6</sup>The time delay between uploading a resume and its inclusion in search results varies between 4–6 hours, so we query periodically

treatment, we assign  $A$  the “stronger” value for the feature, e.g.,  $A$  attended an Ivy League school while  $B$  attended community college. If user  $A$  appears *before* user  $B$ , it means the treatment variable in the resumes has flipped the rank ordering, thus revealing that the algorithm takes that particular resume feature into account. We repeat this process on all three hiring websites, with the different treatment features shown in Table 7. Furthermore, we repeated each treatment three times to ensure that the observed result was consistent.

Table 7 shows the results of our controlled experiments. Crucially, our inferred gender treatment did not influence the order of our users, which confirms that none of the hiring websites engage in direct discrimination with respect to inferred gender. Only three features influenced the ranking algorithms: Monster’s algorithm ranks users by the strength of their alma mater and whether they are currently employed, while CareerBuilder’s algorithm takes employment status and frequency of job changes into account. All of these findings were consistent across repeated trials, which is expected if we assume that these websites use deterministic ranking algorithms.

**Limitations.** Hiring websites may extract features from resumes that are not covered by our treatments in Table 7. We manually examined dozens of real resumes from the hiring websites to inform our selection of features, and were surprised by how consistent the types of informational content were across resumes. We hypothesize that freely available templates and “optimizers”<sup>7</sup> may encourage resume homogeneity across online hiring services. Thus, we believe that our treatments cover the most salient features of typical resumes.

### CONCLUDING DISCUSSION

In this study, we examine gender inequality on the resume search engines provided by Indeed, Monster, and CareerBuilder. We crawled search results for 35 job titles across 20 U. S. cities; these contain data on 855K candidates. Using statistical tests, we examine two types of algorithmic fairness with respect to inferred gender:

- **Individual fairness:** We find statistically significant ( $p \leq 0.05$ ), negative correlations between rank and inferred gender in our dataset. This means that even when controlling for all other visible candidate features, there is a slight penalty against feminine candidates. These results are robust under a variety of conditions. However, the effect size is small: only by rank 30–50 (depending on the website) is the gender effect large enough that masculine candidates receive a substantive increase in rank.
- **Group fairness:** We observe that 8.5–13.2% of job title/city pairs show statistically significant group unfairness. In 12 of 35 job titles, the unfairness benefits men.

Using controlled experiments, we find that none of the hiring sites are directly discriminatory with respect to inferred gender. This concurs with the design of these websites, which do not ask candidates to input their gender. However, we see that

until our users appear. Thus, it is unlikely that our users receive many clicks from recruiters before we measure their ranks. This is important, because clicks may be a feature used to rank candidates.

<sup>7</sup>E.g. <https://www.jobscan.co/>

other hidden features (unemployment and alma mater) are taken into account.

### Why Is There Unfairness?

One unsatisfying aspect of our study is that we are not able to say definitively why there is unfairness with respect to inferred gender on these resume search engines. This is a common criticism of algorithm audits that rely on observational data [37].

We hypothesize that there are two potential causes for the slight individual unfairness we observe. First, the ranking algorithms may rely on a hidden feature that is extracted from resumes that is (weakly) correlated with gender. Our controlled experiments rule out direct discrimination as a cause, and our regressions control for indirect discrimination that might be caused by visible candidate features. Unfortunately, we cannot isolate the hidden feature(s) that may be causing individual unfairness because we do not have access to all candidate resumes on Monster and CareerBuilder.

A second possibility is that small amounts of individual unfairness occur because the algorithms adjust the rank of candidates based on the volume of clicks they receive from recruiters (a so-called learning-to-rank approach [40, 19, 45]). If recruiters are biased, they may generate more clicks on candidates with desirable demographic traits. Testing this hypothesis is challenging, since it would require uploading many resumes with varying features and then waiting weeks in the hope of collecting sufficient clicks to trigger changes in rank.

With respect to group unfairness, the likely cause is structural inequality. It is unlikely to be a coincidence that the job titles where we observe the largest magnitudes of group unfairness include technical professions (e.g., Software Engineer), truck driver, and laborer, i.e., all professions that are historically gendered. Thus, it is fair to say that the ranking algorithms on these hiring sites are not increasing group unfairness on top of what already exists at large in society; rather, they reflect an unfortunate status quo that persists in many professions.

### Interpretation

Are the ranking algorithms on these hiring sites fair, and if so, who is responsible for addressing the situation? Answering these complex questions requires grappling with the desired goals of fairness, and the role of companies in society.

When John F. Kennedy introduced the modern usage of the term *affirmative action*, he asked companies to “take affirmative action to ensure that applicants are employed . . . without regard to their race, creed, color, or national origin” [3]. JFK’s definition of affirmative action was quite conservative, in that he was calling for equal treatment for equivalent candidates. This roughly corresponds to individual fairness.

If we judge these three resume search engines’ output by JFK’s definition of affirmative action, then we conclude that they are largely successful. Although we do observe slightly negative gender coefficients in our models, the effect sizes are sufficiently small that recruiters would need to browse 50 or more candidates before the gender effect substantively impacted the search results. It is unclear whether recruiters browse this deeply into results.

However, there are other interpretations of the meaning of affirmative action. Lyndon B. Johnson famously advocated for “active recruitment,” which encouraged companies to make their workforces more reflective of the overall population by actively hiring underrepresented workers [3]. As noted by Dwork et al., this roughly corresponds to group fairness [20].

If LBJ’s definition of affirmative action is our barometer, then we conclude that these three resume search engines are less successful. It is heartening to observe that the search results for the majority of job titles we investigate are group fair, yet there are professions like Software Engineer where large group unfairness to women remains. Although the ranking algorithms are not responsible for creating this group unfairness, we have to consider whether it is appropriate for these algorithms to perpetuate widely criticized structural inequalities [72].

To address structural inequalities, hiring websites will need to adopt ranking algorithms that are group fair by design. This would ensure that recruiters see the strongest men and women at a rate that reflects the underlying population distribution. Admittedly, having group fair search results does not necessarily mean that candidates from the minority class will be hired more frequently, but it does correct the fundamental problem that unseen candidates have no chance of being interviewed or hired. Furthermore, presenting search results that are strictly representative of population demographics may have the positive effect of combating entrenched stereotypes that discourage hiring in some professions [55].

### Limitations and Future Work

The primary limitation of our work is that we do not know exactly how recruiters interact with resume search engines. It is safe to assume that well-documented psychological biases (e.g., ordering effects) do impact how recruiters use these systems [30, 32, 79, 17, 56]. However, there are differences between the goals of recruiters and search engine users in general that may cause behavioral changes; for example, search engine users often want a single result, while a recruiter may want to interview multiple candidates.

We advocate for future user studies of recruiters. Epstein et al. and Kay et al. both present methodological frameworks that could be adapted to quantify recruiters’ behavior when they interact with resume search engines, as well as the causal impacts of showing unfair search results [22, 55].

Furthermore, additional work is needed to study other hiring websites, especially LinkedIn and Glassdoor [1]. Studying LinkedIn will be complicated due to the monetary cost, the complexity introduced by their social graph, and their history of litigation against third-parties that crawl their data [7, 62].

### Acknowledgments

We thank Christoph Riedl for methodological suggestions, Vette Torvik for help inferring genders, and all of the anonymous reviewers for their insightful suggestions. This research was supported in part by NSF grants IIS-1408345 and IIS-1553088. Any opinions, findings, conclusions, or recommendations expressed herein are those of the authors and do not necessarily reflect the views of the NSF.

## REFERENCES

1. Alexa 2016. Ranking of Emploment Sites. Alexa. (2016). <http://www.alexa.com/topsites/category/Business/Employment>.
2. Ahmed Allam, Peter Johannes Schulz, and Kent Nakamoto. 2014. The impact of search engine selection and sorting criteria on vaccination beliefs and attitudes: two experiments manipulating Google output. *Journal of medical Internet research* 16, 4 (2014), e100.
3. Terry H. Anderson. 2004. *The Pursuit of Fairness: A History of Affirmative Action*. Oxford University Press.
4. Solon Barocas and Andrew D. Selbst. 2016. Big Data's Disparate Impact. *104 California Law Review* 671 (2016).
5. Marc Bendick, Charles W Jackson, and Victor A Reinoso. 1994. Measuring employment discrimination through controlled experiments. *The Review of Black Political Economy* 23, 1 (1994), 25–48.
6. Marianne Bertrand and Sendhil Mullainathan. 2004. Are Emily and Greg More Employable than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination. *The American Economic Review* 94, 4 (2004), 991–1013.
7. Joshua Brustein. 2014. LinkedIn Sues Unknown Hackers in an Attempt to Find Out Who They Are. Bloomberg Businessweek. (January 2014). <https://www.bloomberg.com/news/articles/2014-01-08/linkedin-sues-unknown-hackers-in-an-attempt-to-find-out-who-they-are>.
8. Quoc Trung Bui. 2014. The Most Common Jobs For The Rich, Middle Class And Poor. NPR. (October 2014). <http://n.pr/2iSVi00>.
9. Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. 2009. Building Classifiers with Independency Constraints. In *Proc. of ICDM Workshops*.
10. Toon Calders, Asad Karim, Faisal Kamiran, Wesam Ali, and Xiangliang Zhang. 2013. Controlling Attribute Effect In Linear Regression. In *Proc. of ICDM*.
11. Toon Calders and Sicco Verwer. 2010. Three Naive Bayes Approaches For Discrimination-free Classification. *Data Mining and Knowledge Discovery* 21, 2 (2010), 277–292.
12. Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 6334 (April 2017), 183–186.
13. CareerBuilder 2016. Careerbuilder About US. Careerbuilder. (2016). <http://www.careerbuilder.com/share/aboutus/>.
14. CEB 2011. Most Common Job Titles Posted Online. CEB. (August 2011). <http://bit.ly/2iSYb9h>.
15. Le Chen, Alan Mislove, and Christo Wilson. 2016. An Empirical Analysis of Algorithmic Pricing on Amazon Marketplace. In *Proc. of WWW*.
16. Gregory W Corder and Dale I Foreman. 2014. *Nonparametric statistics: A step-by-step approach*. John Wiley & Sons.
17. Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. An Experimental Comparison of Click Position-bias Models. In *Proc. of WSDM*.
18. Amit Datta, Michael Carl Tschantz, and Anupam Datta. 2015. Automated Experiments on Ad Privacy Settings: A Tale of Opacity, Choice, and Discrimination. In *Proc. of PETS*.
19. Clebson C.A. de Sá, Marcos A. Gonçalves, Daniel X. Sousa, and Thiago Salles. 2016. Generalized BROOF-L2R: A General Framework for Learning to Rank Based on Boosting and Random Forests. In *Proc. of SIGIR*.
20. Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness Through Awareness. In *Proc. of ITCS*.
21. Benjamin G. Edelman, Michael Luca, and Dan Svirsky. 2015. Racial Discrimination in the Sharing Economy: Evidence from a Field Experiment. (2015). <http://ssrn.com/abstract=2701902>.
22. Robert Epstein and Ronald E Robertson. 2015. The search engine manipulation effect (SEME) and its possible impact on the outcomes of elections. *Proceedings of the National Academy of Sciences* 112, 33 (2015), E4512–E4521.
23. Motahhare Eslami, Amirhossein Aleyasen, Karrie Karahalios, Kevin Hamilton, and Christian Sandvig. 2015. Feedvis: A path for exploring news feed curation algorithms. In *Proc. of CSCW Companion*.
24. Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and Removing Disparate Impact. In *Proc. of KDD*.
25. Lisa M Finkelstein, Michael J Burke, and Manbury S Raju. 1995. Age discrimination in simulated employment contexts: An integrative analysis. *Journal of Applied Psychology* 80, 6 (1995), 652.
26. Kevin Fiscella and Allen M Fremont. 2006. Use of Geocoding and Surname Analysis to Estimate Race and Ethnicity (*Health Serv Res*).
27. Sorelle A. Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2016. On the (im)possibility of fairness. *CoRR* abs/1609.07236 (2016).
28. Yanbo Ge, Christopher R. Knittel, Don MacKenzie, and Stephen Zoepf. 2016. *Racial and Gender Discrimination in Transportation Network Companies*. Working Paper 22776. National Bureau of Economic Research.
29. W. Drew Gouvier, Sara Sytsma-Jordan, and Stephen Mayville. 2003. Patterns of discrimination in hiring job applicants with disabilities: The role of disability type, job complexity, and public contact. *Rehabilitation Psychology* 48, 3 (2003), 175.

30. Laura A Granka, Thorsten Joachims, and Geri Gay. 2004. Eye-tracking analysis of user behavior in WWW search. In *Proc. of SIGIR*.
31. Steven Greenhouse. 2013. Supreme Court Raises Bar to Prove Job Discrimination. New York Times. (June 2013). <http://nyti.ms/2ppuG3w>.
32. Zhiwei Guan and Edward Cutrell. 2007. An eye tracking study of the effect of target rank on web search. In *Proc. of CHI*.
33. Saikat Guha, Bin Cheng, and Paul Francis. 2010. Challenges in Measuring Online Advertising Systems. In *Proc. of IMC*.
34. Sara Hajian and Josep Domingo-Ferrer. 2013. A Methodology For Direct And Indirect Discrimination Prevention In Data Mining. *IEEE Transactions on Knowledge and Data Engineering* 25, 7 (2013), 1445–1459.
35. Anikó Hannák, Piotr Sapieżyński, Arash Molavi Kakhki, Balachander Krishnamurthy, David Lazer, Alan Mislove, and Christo Wilson. 2013. Measuring Personalization of Web Search. In *Proc. of WWW*.
36. Anikó Hannák, Gary Soeller, David Lazer, Alan Mislove, and Christo Wilson. 2014. Measuring Price Discrimination and Steering on E-commerce Web Sites. In *Proc. of IMC*.
37. Anikó Hannák, Claudia Wagner, David Garcia, Alan Mislove, Markus Strohmaier, and Christo Wilson. 2017. Bias in Online Freelance Marketplaces: Evidence from TaskRabbit and Fiverr. In *Proc. of CSCW*.
38. Daniel E. Ho, Kosuke Imai, Gary King, and Elizabeth A. Stuart. 2007. Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Political Analysis* 15, 3 (2007), 199–236.
39. Daniel E. Ho, Kosuke Imai, Gary King, and Elizabeth A. Stuart. 2011. MatchIt: Nonparametric Preprocessing for Parametric Causal Inference. *Journal of Statistical Software* 42, 8 (2011), 1–28.
40. Muhammad Ibrahim and Mark Carman. 2016. Comparing Pointwise and Listwise Objective Functions for Random-Forest-Based Learning-to-Rank. *ACM Trans. Inf. Syst.* 34, 4 (2016), 20:1–20:38.
41. Indeed 2016a. How do you rank search results? Indeed. (2016). <http://support.indeed.com/hc/en-us/articles/204488980-How-do-you-rank-search-results->.
42. Indeed 2016b. Indeed Hits Record 200 Million Unique Visitors. Indeed. (2016). <http://blog.indeed.com/2016/02/08/indeed-200-million-unique-visitors/>.
43. Kalervo Järvelin and Jaana Kekäläinen. 2000. IR Evaluation Methods for Retrieving Highly Relevant Documents. In *Proc. of SIGIR*.
44. Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated Gain-based Evaluation of IR Techniques. *ACM Trans. Inf. Syst.* (Oct. 2002).
45. Shan Jiang, Yuening Hu, Changsung Kang, Tim Daly, Jr., Dawei Yin, Yi Chang, and Chengxiang Zhai. 2016. Learning Query and Document Relevance from a Web-scale Click Graph. In *Proc. of SIGIR*.
46. Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. 2017. Unbiased Learning-to-Rank with Biased Feedback. In *Proc. of WSDM*.
47. Isaac Johnson, Connor McMahon, Johannes Schöning, and Brent Hecht. 2017. The Effect of Population and “Structural” Biases on Social Media-based Algorithms – A Case Study in Geolocation Inference Across the Urban-Rural Spectrum. In *Proc. of CHI*.
48. Faisal Kamiran and Toon Calders. 2009. Classifying without discriminating. In *Proc. of Conference on Computer, Control and Communication*.
49. Faisal Kamiran and Toon Calders. 2010. Classification with No Discrimination by Preferential Sampling. In *Proc. of Machine Learning Conference of Belgium and The Netherlands*.
50. Faisal Kamiran, Toon Calders, and Mykola Pechenizkiy. 2010. Discrimination Aware Decision Tree Learning. In *Proc. of ICDM*.
51. Faisal Kamiran, Asad Karim, and Xiangliang Zhang. 2012. Decision Theory for Discrimination-Aware Classification. In *Proc. of ICDM*.
52. Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2012. Fairness-Aware Classifier with Prejudice Remover Regularizer. In *Proc. of ECML PKDD*.
53. Toshihiro Kamishima, Shotaro Akaho, and Jun Sakuma. 2011. Fairness-aware Learning through Regularization Approach. In *Proc. of ICDM Workshops*.
54. Fariba Karimi, Claudia Wagner, Florian Lemmerich, Mohsen Jadidi, and Markus Strohmaier. 2016. Inferring Gender from Names on the Web: A Comparative Evaluation of Gender Detection Methods (*WWW ’16 Companion*).
55. Matthew Kay, Cynthia Matuszek, and Sean A. Munson. 2015. Unequal Representation and Gender Stereotypes in Image Search Results for Occupations. In *Proc. of CHI*.
56. Mark T. Keane, Maeve O’Brien, and Barry Smyth. 2008. Are People Biased in Their Use of Search Engines? *Commun. ACM* 51, 2 (Feb. 2008), 49–52.
57. Pauline T. Kim. 2017. Data-Driven Discrimination at Work. *William & Mary Law Review* 58 (2017).
58. Chloe Kliman-Silver, Anikó Hannák, David Lazer, Christo Wilson, and Alan Mislove. 2015. Location, Location, Location: The Impact of Geolocation on Web Search Personalization. In *Proc. of IMC*.

59. Juhi Kulshrestha, Motahhare Eslami, Johnnatan Messias, Muhammad Bilal Zafar, Saptarshi Ghosh, Krishna P. Gummadi, and Karrie Karahalios. 2017. Quantifying Search Bias: Investigating Sources of Bias for Political Searches in Social Media. In *Proc. of CSCW*.
60. Shyong (Tony) K. Lam, Anuradha Uduwage, Zhenhua Dong, Shilad Sen, David R. Musicant, Loren Terveen, and John Riedl. 2011. WP:Clubhouse? An Exploration of Wikipedia's Gender Imbalance. In *Proc. of WikiSym*.
61. Mathias Lécuyer, Guillaume Ducoffe, Francis Lan, Andrei Papancea, Theofilos Petsios, Riley Spahn, Augustin Chaintreau, and Roxana Geambasu. 2014. XRay: Enhancing the Web's Transparency with Differential Correlation. In *Proc. of USENIX Security Symposium*.
62. Thomas Lee. 2017. LinkedIn, HiQ spat presents big questions for freedom, innovation. San Francisco Chronicle. (July 2017). <http://www.sfchronicle.com/business/article/LinkedIn-HiQ-spat-presents-big-questions-for-11274133.php>.
63. Lisa Madigan 2017. Madigan Probes National Job Search Sites Over Potential Age Discrimination. Illinois Attorney General Press Release. (March 2017). [http://www.illinoisattorneygeneral.gov/pressroom/2017\\_03/20170302.html](http://www.illinoisattorneygeneral.gov/pressroom/2017_03/20170302.html).
64. Wendy Liu and Derek Ruths. 2013. What's in a Name? Using First Names as Features for Gender Inference in Twitter. In *AAAI Spring Symposium Series*.
65. Ramona Ludolph, Ahmed Allam, and Peter J Schulz. 2016. Manipulating Google's Knowledge Graph box to counter biased information processing during an online search on vaccination: application of a technological debiasing strategy. *Journal of medical Internet research* 18, 6 (2016).
66. Binh Thanh Luong, Salvatore Ruggieri, and Franco Turini. 2011. k-NN As an Implementation of Situation Testing for Discrimination Discovery and Prevention. In *Proc. of KDD*.
67. Gabriel Magno, Camila Souza Araujo, Wagner Meira Jr., and Virgílio A. F. Almeida. 2016. Stereotypes in Search Engine Answers: Local or Global? *CoRR* abs/1609.05413 (2016).
68. Mediative 2015. Keeping an eye on Google - Eye tracking SERPs through the years. Mediative. (2015). <http://www.mediative.com/eye-tracking-google-through-the-years/>.
69. Jakub Mikians, László Gyarmati, Vijay Erramilli, and Nikolaos Laoutaris. 2012. Detecting Price and Search Discrimination on the Internet. In *Proc. of HotNets*.
70. Jakub Mikians, László Gyarmati, Vijay Erramilli, and Nikolaos Laoutaris. 2013. Crowd-assisted Search for Price Discrimination in e-Commerce: First Results. In *Proc. of ACM CoNEXT*.
71. Monster 2016. Monster About US. Monster. (2016). <http://www.monster.com/about/>.
72. Liza Mundy. 2017. Why Is Silicon Valley So Awful to Women? The Atlantic. (April 2017). <https://www.theatlantic.com/magazine/archive/2017/04/why-is-silicon-valley-so-awful-to-women/517788/>.
73. Alamir Novin and Eric Meyers. 2017. Making Sense of Conflicting Science Information: Exploring Bias in the Search Engine Result Page. In *Proceedings of the 2017 Conference on Human Information Interaction and Retrieval*. ACM, 175–184.
74. Devah Pager. 2008. *Marked: Race, crime, and finding work in an era of mass incarceration*. University of Chicago Press.
75. Devah Pager and Hana Shepherd. 2008. The sociology of discrimination: Racial discrimination in employment, housing, credit, and consumer markets. *Annual review of sociology* 34 (2008), 181.
76. Dino Pedreshi, Salvatore Ruggieri, and Franco Turini. 2008. Discrimination-aware Data Mining. In *Proc. of KDD*.
77. Joseph Reagle and Lauren Rhue. 2011. Gender Bias in Wikipedia and Britannica. *International Journal of Communication* 5 (2011), 1138–1158.
78. Peter A Riach and Judith Rich. 1991. Measuring discrimination by direct experimental methods: seeking gunsmoke. *Journal of Post Keynesian Economics* 14, 2 (1991), 143–150.
79. Matthew Richardson. 2007. Predicting clicks: Estimating the click-through rate for new ads. In *Proc. of WWW*.
80. Christian Sandvig, Kevin Hamilton, Karrie Karahalios, and Cedric Langbort. 2014. Auditing algorithms: Research methods for detecting discrimination on internet platforms. In *Proceedings of "Data and Discrimination: Converting Critical Concerns into Productive Inquiry", a preconference at the 64th Annual Meeting of the International Communication Association*.
81. Jennifer L. Skeem and Christopher T. Lowenkamp. 2016. Risk, Race, & Recidivism: Predictive Bias and Disparate Impact. (2016). <https://ssrn.com/abstract=2687339>.
82. Brittany N. Smith, Mamta Singh, and Vetle I. Torvik. 2013. A Search Engine Approach to Estimating Temporal Changes in Gender Orientation of First Names. In *Proceedings of the 13th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL '13)*. 199–208. DOI: <http://dx.doi.org/10.1145/2467696.2467720>
83. Social Security Administration 2016. Baby Names from Social Security Card Applications-National Level Data. data.gov. (2016). <https://catalog.data.gov/dataset/baby-names-from-social-security-card-applications-national-level-data>.

84. Gary Soeller, Karrie Karahalios, Christian Sandvig, and Christo Wilson. 2016. MapWatch: Detecting and Monitoring International Border Personalization on Online Maps. In *Proc. of WWW*.
85. Latanya Sweeney. 2013. Discrimination in Online Ad Delivery. (2013). <http://ssrn.com/abstract=2208240>.
86. Cong Tang, Keith W. Ross, Nitesh Saxena, and Ruichuan Chen. 2011. What's in a Name: A Study of Names, Gender Inference, and Gender Behavior in Facebook.. In *DASFAA Workshops*.
87. Jacob Thebault-Spieker, Loren G. Terveen, and Brent Hecht. 2015. Avoiding the South Side and the Suburbs: The Geography of Mobile Crowdsourcing Markets. In *Proc. of CSCW*.
88. Vette I. Torvik and Sneha Agarwal. 2016. Ethnea – an instance-based ethnicity classifier based on geo-coded author names in a large-scale bibliographic database. (2016). <http://hdl.handle.net/2142/88927>  
International Symposium on Science of Science March 22-23, 2016 - Library of Congress, Washington DC, USA.
89. Margery Austin Turner, Michael Fix, and Raymond J Struyk. 1991. *Opportunities denied, opportunities diminished: Racial discrimination in hiring*. The Urban Institute.
90. Sander van der Linden, Anthony Leiserowitz, Seth Rosenthal, and Edward Maibach. 2017. Inoculating the public against misinformation about climate change. *Global Challenges* 1, 2 (2017).
91. Claudia Wagner, David Garcia, Mohsen Jadidi, and Markus Strohmaier. 2015. It's a Man's Wikipedia? Assessing Gender Inequality in an Online Encyclopedia. In *Proc. of ICWSM*.
92. Doris Weichselbaumer. 2003. Sexual orientation discrimination in hiring. *Labour Economics* 10, 6 (2003), 629–642.
93. David L. Word, Charles D. Coleman, Robert Nunziata, and Robert Kominski. 2000. Demographic Aspects of Surnames from Census 2000. census.gov. (2000). <https://www2.census.gov/topics/genealogy/2000surnames/surnames.pdf>.
94. Ke Yang and Julia Stoyanovich. 2017. Measuring Fairness in Ranked Outputs. In *Proc. of SSDBM*.
95. Richard Zemel, Yu Wu, Kevin Swersky, Toniann Pitassi, and Cynthia Dwork. 2013. Learning Fair Representations. In *Proc. of ICML*.
96. Indre Zliobaite, Faisal Kamiran, and Toon Calders. 2011. Handling Conditional Discrimination. In *Proc. of ICDM*.