

# Quantifying Query Fairness Under Unawareness

THOMAS JAENICH<sup>\*</sup>, University of Glasgow, UK  
ALEJANDRO MOREO<sup>†</sup>, Consiglio Nazionale delle Ricerche, IT  
ALESSANDRO FABRIS, University of Trieste, IT  
GRAHAM MCDONALD, University of Glasgow, UK  
ANDREA ESULI, Consiglio Nazionale delle Ricerche, IT  
IADH OUNIS, University of Glasgow, UK  
FABRIZIO SEBASTIANI, Consiglio Nazionale delle Ricerche, IT

Traditional ranking algorithms are designed to retrieve the most relevant items for a user’s query, but they often inherit biases from data that can unfairly disadvantage vulnerable groups. Fairness in information access systems (IAS) is typically assessed by comparing the distribution of groups in a ranking to a target distribution, such as the overall group distribution in the dataset. These fairness metrics depend on knowing the true group labels for each item. However, when groups are defined by demographic or sensitive attributes, these labels are often unknown, leading to a setting known as “fairness under unawareness.” To address this, group membership can be inferred using machine-learned classifiers, and group prevalence is estimated by counting the predicted labels. Unfortunately, such an estimation is known to be unreliable under dataset shift, compromising the accuracy of fairness evaluations. In this paper, we introduce a robust fairness estimator based on quantification that effectively handles multiple sensitive attributes beyond binary classifications. Our method outperforms existing baselines across various sensitive attributes and, to the best of our knowledge, is the first to establish a reliable protocol for measuring fairness under unawareness across multiple queries and groups.

**JAIR Track:** AI and Society Track

**JAIR Associate Editor:** Ulises Cortés

## JAIR Reference Format:

Thomas Jaenich, Alejandro Moreo, Alessandro Fabris, Graham McDonald, Andrea Esuli, Iadh Ounis, and Fabrizio Sebastiani. 2026. Quantifying Query Fairness Under Unawareness. *Journal of Artificial Intelligence Research* 85, Article 7 (January 2026), 19 pages. DOI: [10.1613/jair.1.17675](https://doi.org/10.1613/jair.1.17675)

## 1 Introduction

In addition to ensuring the relevance of search results, preventing unfairness and discrimination in ranking has become a fundamental objective in the development of information access systems (IAS) (Ekstrand, Das, et al. 2022; Zehlike, K. Yang, et al. 2022). With this in mind, there have been many approaches proposed in

<sup>\*</sup>Equal contribution

<sup>†</sup>Equal contribution

---

Authors’ Contact Information: Thomas Jaenich, ORCID: [0009-0009-6347-408X](https://orcid.org/0009-0009-6347-408X), [tjaenich.1@research.gla.ac.uk](mailto:tjaenich.1@research.gla.ac.uk), University of Glasgow, Glasgow, UK; Alejandro Moreo, ORCID: [0000-0002-0377-1025](https://orcid.org/0000-0002-0377-1025), [alejandro.moreo@isti.cnr.it](mailto:alejandro.moreo@isti.cnr.it), Consiglio Nazionale delle Ricerche, Pisa, IT; Alessandro Fabris, ORCID: [0000-0001-6108-9940](https://orcid.org/0000-0001-6108-9940), [alessandro.fabris@units.it](mailto:alessandro.fabris@units.it), University of Trieste, Trieste, IT; Graham McDonald, ORCID: [0000-0002-1266-5996](https://orcid.org/0000-0002-1266-5996), [graham.macdonald@glasgow.ac.uk](mailto:graham.macdonald@glasgow.ac.uk), University of Glasgow, Glasgow, UK; Andrea Esuli, ORCID: [0000-0002-5725-4322](https://orcid.org/0000-0002-5725-4322), [andrea.esuli@isti.cnr.it](mailto:andrea.esuli@isti.cnr.it), Consiglio Nazionale delle Ricerche, Pisa, IT; Iadh Ounis, ORCID: [0000-0003-4701-3223](https://orcid.org/0000-0003-4701-3223), [iadh.ounis@glasgow.ac.uk](mailto:iadh.ounis@glasgow.ac.uk), University of Glasgow, Glasgow, UK; Fabrizio Sebastiani, ORCID: [0000-0003-4221-6427](https://orcid.org/0000-0003-4221-6427), [fabrizio.sebastiani@isti.cnr.it](mailto:fabrizio.sebastiani@isti.cnr.it), Consiglio Nazionale delle Ricerche, Pisa, IT.



This work is licensed under a [Creative Commons Attribution International 4.0 License](https://creativecommons.org/licenses/by/4.0/).

© 2026 Copyright held by the owner/author(s).

DOI: [10.1613/jair.1.17675](https://doi.org/10.1613/jair.1.17675)

the literature for mitigating unfairness in the results of IAS (Biega et al. 2018; Geyik et al. 2019; Heuss et al. 2022; Jaenich et al. 2023, 2024; M. Morik et al. 2020; Singh and Joachims 2018). Providing fair search results is crucial, since ranking items that belong to groups identified by sensitive attributes can significantly impact real-world outcomes, such as economic opportunities (L. Chen et al. 2018; Pedreschi et al. 2008). When items that are associated with a specific demographic attribute are systematically ranked lower in the search results than items from other demographics, the low-ranked items will receive less attention from users, since items high up in the ranking are more likely to be examined by users (Craswell et al. 2008). This can be problematic in practical scenarios. For example, in job search, this positional bias means that recruiters may only notice candidates from the top-ranked applications, potentially overlooking qualified individuals who are ranked lower. While this will not always warrant a bias-mitigating intervention, it is imperative to at least *evaluate* and *monitor* the fairness of rankings in high-stakes domain (European Commission 2024; New York City Council 2021).

Another important factor influencing the fairness of search results is the query that is issued by the user. The degree of unfairness in the search results can vary across different queries, depending on how a query is formulated. Such unfairness can be introduced either directly by the user, e.g., through the replication of existing biases when formulating the query (Kopeinik et al. 2023), or automatically, e.g., through the auto-completion features of a search system (L. Chen et al. 2018). Therefore, assessing the fairness of search results related to a specific query is crucial to determine whether fairness interventions are needed. We introduce the terminology *query fairness estimation* (QFE) for the task of assessing the fairness of search results for a given query.

To perform QFE, one typically needs access to the group labels of the ranked items (Kuhlman et al. 2021; Raj and Ekstrand 2022; Zehlike, K. Yang, et al. 2022). These labels categorise items by sensitive attributes, such as race or gender. However, access to these group labels is often limited due to legal, ethical, or other data availability constraints (Bogen et al. 2020; Holstein et al. 2019). As a result, fairness evaluations must often occur under “unawareness,” where the labels are unknown.

One way to obtain the labels under unawareness is to use human annotators. However, this is costly and impractical in most scenarios. An alternative is to deploy classifiers to infer the document labels automatically. A classification method deployed in practice is the *Bayesian improved surname geocoding* (BISG) that is used to infer race from surnames and ZIP codes using Bayesian statistics and US Census data (Adjaye-Gbewonyo et al. 2014). While cost-efficient, deploying classifiers can introduce more unintended unfairness and bias by the classifiers themselves (Ghosh, Dutt, et al. 2021).

Moreover, even seemingly accurate classifiers can lead to unreliable results when performing QFE. For example, a good classifier may achieve high accuracy by focusing disproportionately on one class, thus minimising errors like false positives at the expense of increasing errors like false negatives (Esuli et al. 2023, §1.2). While technically accurate, this skewed performance fails to reflect the true distribution of groups, making it unreliable for assessing the proportions in a broader set of documents, and can lead to significant errors and misjudgments in the measurements of fairness (J. Chen et al. 2019).

In this work, we propose the use of quantification techniques, i.e., machine learning models specifically trained to estimate the relative frequencies of the classes in unlabelled data (Esuli et al. 2023) to improve QFE. Specifically, our main focus is QFE under unawareness of sensitive attributes with multiple classes. Fairness estimation under unawareness is not only of theoretical interest but also has practical relevance. For example, search engines and recommender systems are increasingly subject to fairness audits, where it is necessary to assess whether exposure across demographic groups is balanced without direct access to sensitive attributes. With this work, we provide a principled way to support such audits, aligning with recent efforts in fairness-aware evaluation of information access systems (Y. Wang et al. 2023). To the best of our knowledge, our proposed approach is the first to cover groups with non-binary protected attributes and to estimate ranking fairness across multiple queries. Our main contributions are as follows:

- We propose a new family of principled methods to perform Query Fairness Estimation across multiple queries and non-binary sensitive attributes.
- We introduce the first approach designed to make quantification methods robust against sample selection bias.
- Through extensive experiments on the TREC 2022 Fair Ranking Track collection (Ekstrand, McDonald, et al. 2022), we demonstrate that our quantification-based approach outperforms previous methods.

## 2 Related Work

In recent years, ensuring fairness in search results has emerged as a crucial objective alongside the traditional goal of relevance in the development of IAS (Ekstrand, Das, et al. 2022; Zehlike, K. Yang, et al. 2022). To measure the fairness of search results for a given query, i.e., for the task of QFE, several measures have been introduced (Biega et al. 2018; Diaz et al. 2020; Kirnap et al. 2021; Kuhlman et al. 2021; Raj and Ekstrand 2022; Sapiezynski et al. 2019; Singh and Joachims 2018; E. Yang et al. 2024; K. Yang and Stoyanovich 2017). While these measures cover different notions of fairness, for example *equality of opportunity* (Biega et al. 2018; Diaz et al. 2020; Hardt et al. 2016; Singh and Joachims 2018) or *statistical parity* (Geyik et al. 2019; Sapiezynski et al. 2019; Zehlike, Bonchi, et al. 2017), they all depend on accurate knowledge of the document labels in a ranking. In this work, we consider an “unawareness” scenario, where group labels are unavailable, requiring alternative solutions to ensure accurate fairness assessments.

Related to this, (Ghosh, Dutt, et al. 2021) have conducted an extensive study showing that inferring labels using standard classifiers can be problematic. Their work highlights the need for reliable methods to access accurate fairness labels. Although several studies have focused on enhancing the fairness of classifiers with noisy or incomplete group labels (Celis et al. 2021; Friedler et al. 2021; Ghosh, Kvitca, et al. 2023; Mozannar et al. 2020; S. Wang et al. 2020), they primarily address fairness metrics for classification, not for ranking tasks.

In this work, we focus specifically on QFE for rankings under unawareness, an area that has received less attention compared to classification. In the absence of group labels, (F. Chen and Fang 2023) proposed a distribution-based learning approach that leverages contextual features. Their approach uses a loss function that does not require explicit group labels but instead targets a fair distribution. Unlike their task of mitigating unfairness, our main objective is to obtain reliable estimates of group proportions in a ranking to accurately measure the fairness of search results.

(Kirnap et al. 2021) proposed a method using a small query-dependent subset of data annotated by human assessors for QFE. Moreover, they only focus on binary group fairness, comparing protected versus non-protected groups. Our work addresses the characteristics of multiclass groups and does not require impractical human annotations for each query.

Related to our work on QFE in rankings, (Ghazimatin et al. 2022) have proposed an approach that we term Post-Metric Correction (PMC). Both our approach and the PMC variants include a correction phase and rely on the outputs of an underlying classifier. However, there are significant differences between our quantification-based approaches and the PMC variants. First, the PMC variants are designed only for binary group fairness assessment; our method natively caters to multi-valued sensitive attributes. This is exceedingly rare and important in algorithmic fairness research (Simson et al. 2024). Moreover, each PMC variant is tied to a specific independence assumption, which may prove difficult to verify in general settings. Finally, the PMC methods are also tailored to one specific fairness metric, while our approach is more versatile. Our method applies a general pre-correction to the class prevalence estimates, which are then used to compute different fair ranking metrics.

In our work, we propose to use quantification methods to make QFE robust to the limitations when standard classifiers are applied. In a related effort, (Fabris et al. 2023) have shown that using quantification techniques is a useful way to assess the fairness of algorithms when the labels are unknown. However, their work focuses on classification problems, while our work focuses specifically on QFE.

The approaches introduced above highlight progress on both fairness estimation and quantification, but they leave important gaps when applied to fairness under unawareness. In particular, classifier-based methods such as PMC are restricted to binary sensitive attributes and rely on metric-specific adaptations. The existing work on quantification, in contrast, has primarily focused on prevalence estimation under dataset shift without considering its implications for fairness evaluation in information retrieval systems. However, a principled connection between quantification methods and the estimation of fairness metrics in multiclass settings under unawareness has not yet been established. The next section develops this connection by formulating fairness estimation as a quantification problem and adapting established quantification methods to the information retrieval context.

### 3 Proposed Approach

#### 3.1 Running Example

Consider an online hiring platform assisting recruiters to fill job openings with promising candidates. Recruiters query the platform with job descriptions and get ranked lists of candidates in return, in decreasing order of estimated job fitness. Platform developers want to ensure that their models are fair with respect to multiple attributes, including age, gender, race, and ethnicity. Given the sensitive nature of this information, most data subjects will be hesitant to disclose it (Bogen et al. 2020). Developers, therefore, obtain sensitive attributes for a subset of users through voluntary data disclosure (LinkedIn 2024; Wilson et al. 2021). This incomplete demographic data can be used to estimate platform fairness across all queries and users. Finding the optimal way to perform this estimate is an open research problem tackled in the remainder of this section.

#### 3.2 Learning to Quantify

The field of quantification emerges from the fundamental observation that counting over the labels predicted by a classifier tends to produce poor estimates of class prevalence (Esuli et al. 2023, §1.2), unless the classifier is a perfect one. The above naïve counting approach has come to be known as the “Classify & Count” (CC) method (Forman 2005), and nowadays represents the strawman baseline any proper quantification method is expected to beat.

More formally, a quantifier is a function  $\lambda : \mathbb{N}^X \rightarrow \Delta^{n-1}$  mapping bags (or multi-sets) of instances from the input space  $\mathcal{X} = \mathbb{R}^d$  to the probability simplex, so that  $\lambda(\mathbf{X}) = \mathbf{p}$  lies on the  $(n - 1)$ -simplex defined as  $\Delta^{n-1} = \{p_1, \dots, p_n : p_i \geq 0, \sum_i p_i = 1\}$ , in which  $n$  is the number of classes  $\mathcal{Y} = \{1, \dots, n\}$  and  $p_i$  is the prior probability (a.k.a. “class frequency”, or “class prevalence”) of class  $i$  in bag  $\mathbf{X}$ . Given a classifier  $\phi : \mathcal{X} \rightarrow \mathcal{Y}$ , and a bag  $\mathbf{X}$ , CC is defined as

$$\text{CC}(\mathbf{X})_i = \hat{p}_i = \frac{1}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} \mathbb{1}[\phi(\mathbf{x}) = i] \quad (1)$$

**A caveat on terminology.** This paper integrates concepts from different disciplines (quantification, information retrieval, and fairness), each of which employs its own consolidated terminology. Throughout this paper, we will interchangeably use the terms “classes” (here denoted by  $\mathcal{Y}$ ) and “groups” (often denoted by  $\mathcal{A}$  in the fairness literature), “labels” and “sensitive attributes”. The reader should also note that, despite referring to different concepts, we might interchangeably use “bags” (or “multi-sets”) and “ranked lists of items”, since our quantifiers regard the latter as unordered objects.

**Dataset shift.** The essence of quantification is that of tackling a situation in which there is a change (“shift”, or “drift”) between the distribution  $P_{tr}$  from which instances used to train the quantifier have been drawn and the distribution  $P_{te}$  from which the test data are drawn. In online hiring, this corresponds to a realistic setting where the training set, consisting of candidates who disclose their sensitive data  $y$ , is not sampled IID from

the same distribution as the test set. The scenario in which  $P_{tr}(X, Y) \neq P_{te}(X, Y)$  is generally known as *dataset shift* (Storkey 2009).

Although  $P_{tr}$  and  $P_{te}$  are rather standard notation in machine learning for referring to the training and test distributions, throughout this paper we will consider more than two such distributions. For this reason, we will use the nomenclature  $P_A$  to refer to the distribution from which an empirical sample of data items  $A$  has been drawn. In this way, we use  $P_L$  to denote the distribution from which labelled documents are drawn, and  $P_U$  to denote the distribution from which the unlabelled documents are drawn. Further distributions will be introduced when needed.

Among the main types of shift that have been described in the literature, quantification has traditionally focused on *prior probability shift* (PPS). This type of shift is characteristic of *anti-causal learning* (Schölkopf et al. 2012) – i.e., learning problems in which the covariates represent symptoms of the phenomenon we want to predict, and that are typically modelled via the factorization  $P(X, Y) = P(X|Y)P(Y)$  – and is characterized by the fact that  $P_L(Y) \neq P_U(Y)$  while  $P_L(X|Y) = P_U(X|Y)$ .

**A simple quantifier.** Arguably, the simplest quantification method devised to counter PPS is the so-called *Adjusted Classify & Count* (ACC) (Forman 2005). ACC is better described in the binary case  $\mathcal{Y} = \{0, 1\}$  (the multiclass extension is straightforward), by observing that

$$P_U(\hat{Y} = 1) = P_U(\hat{Y} = 1|Y = 1)P_U(Y = 1) + P_U(\hat{Y} = 1|Y = 0)P_U(Y = 0) \quad (2)$$

where  $P_U(\hat{Y} = 1)$  corresponds to the fraction of predicted positives (that we can estimate via CC) and  $P_U(\hat{Y} = 1|Y = 1)$  and  $P_U(\hat{Y} = 1|Y = 0)$  are the *true positive rate* (tpr) and *false positive rate* (fpr) of the classifier  $\phi$ . These two quantities can be estimated using the training data given that the class-conditional distributions of the training and test data are assumed invariant. ACC is thus defined as

$$\text{ACC}(\mathbf{X})_1 = \frac{\text{CC}(\mathbf{X})_1 - \hat{\text{fpr}}}{\hat{\text{tpr}} - \hat{\text{fpr}}} \quad (3)$$

In online hiring, ACC can estimate prevalence of different ethnicities in sets of candidates. To aid intuition, we illustrate this point in the context of our running example of online hiring. Suppose a recruiter issues a query for candidates and the system retrieves 100 applications. A classifier predicts that 40 candidates belong to a minority demographic group and 60 to the majority group. In reality, however, the true distribution is 30 versus 70. Although the classifier may still achieve high overall accuracy, the misestimation of group proportions distorts the fairness evaluation: fairness metrics that rely on exposure across groups would incorrectly indicate that minority candidates receive greater visibility than they actually do. Quantification methods are designed to correct such systematic biases in prevalence estimates, thereby providing a more reliable basis for assessing fairness in information retrieval. This distinction clarifies why the Classify & Count method is inadequate in our setting and motivates the use of dedicated quantifiers.

### 3.3 Countering Sample Selection Bias

**Origin.** Consider the random variable  $Q$  that takes on values 1 (“the item is relevant”) and 0 (“the item is irrelevant”) with respect to a specific query. Note that, for a generic distribution  $P$ , the class-conditional distribution is a mixture of relevant and irrelevant items, i.e.,  $P(X, Q|Y) = P(X|Y, Q = 1)P(Q = 1) + P(X|Y, Q = 0)P(Q = 0)$ . However, if  $U_q$  are the test documents retrieved for a query, and we denote  $P_{U_q}$  the distribution from which  $U_q$  is drawn, note that we might expect  $P_L(Q = 1) \ll P_{U_q}(Q = 1)$ , since it is likely that the vast majority of the training data used for learning our quantifier is irrelevant to a specific query, while the majority of the items retrieved for the query are indeed relevant to it. The class-conditional distributions of  $P_L$  and  $P_{U_q}$  are thus different, and this clashes with the PPS assumptions.

Note that the random variable  $Q$  might be regarded as the “selection variable”, which is representative of a type of dataset shift known as *sample selection bias* (SSB).<sup>1</sup> While different quantification methods have shown varying degrees of effectiveness in addressing different types of shift (González et al. 2024), we are unaware of any quantification method robust to SSB.

**Mitigation.** Let us turn back to the ACC method (Equation (3)) to illustrate the problem (a similar rationale applies to other quantification algorithms as well, as we discuss in Section 3.4). Recall that ACC replaces  $P_U(\hat{Y} = 1|Y)$  with  $P_L(\hat{Y} = 1|Y)$  in Equation (2) on the grounds that the class-conditional distributions of training and test datapoints are invariant. That  $P_L(X|Y) = P_U(X|Y)$  implies  $P_L(\hat{Y}|Y) = P_U(\hat{Y}|Y)$  follows from the fact that  $\hat{Y}$  depends uniquely on  $X$  by means of a function (the classifier); inasmuch as the classifier is a measurable function this equivalence holds; see Lemma 1 in (Lipton et al. 2018).

In general,  $P_U(\hat{Y} = i|Y = j)$  represents the *classification rates* of the classifier in the *test set*, and is given by

$$P_U(\hat{Y} = i|Y = j) = \mathbb{E}_{\mathbf{x} \sim P_U(X|Y=j)} [\mathbb{1}[\phi(\mathbf{x}) = i]] \quad (4)$$

Of course, we do not have access to the true distribution  $P_U(X|Y = j)$  of the expectation, but if we could assume  $P_U(X|Y) = P_L(X|Y)$ , then this expectation could be estimated by means of an empirical distribution  $\mathbf{x}_1, \dots, \mathbf{x}_m \sim P_L(X|Y = j)$ , as

$$P_U(\hat{Y} = i|Y = j) \approx \frac{1}{m} \sum_{k=1}^m \mathbb{1}[\phi(\mathbf{x}_k) = i] \quad (5)$$

Although we know that this assumption is flawed in the presence of SSB (Section 3.3), one fundamental observation arises: the pitfall stems from the choice of the empirical distribution used to characterize the classifier  $\phi$ , rather than from the classifier itself.

This is important since most quantifiers use a training set  $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ ,  $\mathbf{x}_i \in \mathcal{X}$ ,  $y_i \in \mathcal{Y}$  to both learn a classifier  $\phi$ , and estimate, via cross-validation,<sup>2</sup> its classification rates (in our binary example, this reduces to estimating tpr and fpr). Note also that training a classifier is costly, whereas learning the classification rates is rather inexpensive as it only involves issuing predictions and rearranging counts.

**Main idea.** We propose to disentangle the classifier-training phase from the correction learned by the quantifier. We therefore assume to have access to two sets of labelled data,  $L_\phi = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$  that we use to train our classifier (offline since it is costly), and  $L_{corr} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{m'}$  that we use to learn the correction (at query time since it is inexpensive). What remains to make quantification robust to SSB is to counter the selection bias between  $L_q$  and  $L_{corr}$ .

Figure 1 summarizes our approach, with  $\mathbf{M}$  and  $\mathbf{t}$  denoting item representations explained Section 3.4. Reducing the sampling bias from the test data is impossible; the test items are retrieved by the search engine precisely to guarantee that items are relevant to a specific query. The main idea we propose in this paper is *to use the same search engine with the same query* to retrieve, from an *auxiliary* pool of training documents  $L_{corr}$  (with  $L_{corr} \cap L_\phi = \emptyset$ ), a subset of training items  $L_q \subset L_{corr}$  that are biased towards the query similarly to items in  $U_q$ . This way, the empirical distribution  $L_q$  that we retrieve from the auxiliary pool ( $L_{corr}$ ) can now be regarded as a sample from a query-biased distribution  $P_{L_q}$  and, since we can now assume  $P_{L_q}(Q) \approx P_{U_q}(Q)$  (i.e., both distributions are biased towards the query) then we can also assume  $P_{L_q}(X, Q|Y) \approx P_{U_q}(X, Q|Y)$  and restore the fundamental PPS assumption. Given that the SSB in  $L_q$  mimics the SSB that affects the test data, *the sampling bias shift vanishes*.

<sup>1</sup>Sample selection bias is often defined differently, since the selection variable has an effect on the way the training instances (and not the test instances, as in our case) are selected (Storkey 2009).

<sup>2</sup>This is in order to avoid the same datapoint being classified to take part in the training of the classifier.

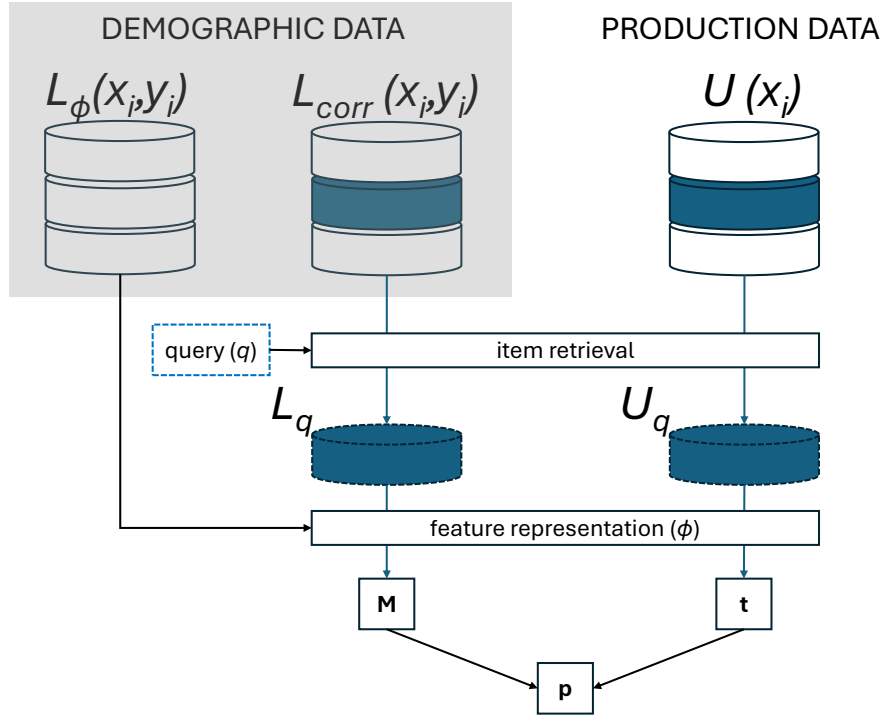


Fig. 1. Schematic illustration of our proposed approach for Query Fairness Evaluation. Demographic data is split into  $L_\phi$ , used to learn a feature representation function  $\phi$ , and  $L_{corr}$ , employed for query-specific corrections. For each query  $q$ , the prevalence  $p$  of sensitive attributes in the retrieved production data  $U_q$  is estimated by representing ranked items according to  $\phi$  and leveraging a subset  $L_q$ , containing the top-ranked items of  $L_{corr}$ , to learn a correction factor  $M$  specific to  $q$ .

We thus consider an auxiliary set of labelled items  $L_{corr}$  containing pairs  $(x_i, y_i)$  labelled by sensitive attributes (hereafter called the “correction pool”). From this set  $L_{corr}$ , we select, using the same retrieval model and query that we issue on the test pool  $(U)$ , a ranked list of (labelled) items  $L_q$  that we use to learn a per-query quantification correction to estimate group prevalence in the top- $k$  prefixes of (unlabelled) rankings  $U_q$ . Considering ACC applied to online hiring as an example, to estimate the prevalence of different ethnicities in search results  $U_q$  for a given query, we compute the classification rates for Equation (3) from a subset  $L_q \subset L_{corr}$  of top-ranked items, instead of using the whole labelled set  $L_{corr}$ . To make the main idea concrete, consider a query whose top- $k$  results in reality contain about one third of documents from a particular group. When a query retrieves a set of documents, its composition differs from that of the training data. As a result, the classifier tends to make different errors on this set. If we then estimate group proportions by simply counting predictions, the result can be biased and may overstate the true share. This is an instance of sample-selection bias: the evaluation set is not representative of the data used to estimate error rates. As a consequence, fairness metrics computed from such distorted estimates can indicate disparities that are not actually present. To counter this, we estimate a per-query correction on a small, query-matched labelled subset from the correction pool, which produces proportions closer to the true distribution before computing the fairness metrics.

### 3.4 Quantifying Query Fairness

In the previous section, we have explained the intuitions behind our method with respect to (binary) ACC, a relatively simple quantifier. In this section, we generalize the rationale to more sophisticated multiclass quantification methods.

In the modern perspective of multiclass quantification (Bunse and K. Morik 2022), most quantifiers can be framed as the problem of solving for  $\mathbf{p} \in \Delta^{n-1}$  (unknown prevalences) the system of linear equations

$$\mathbf{t} = \mathbf{M}\mathbf{p} \quad (6)$$

where  $\mathbf{t} = \Phi(U)$  is the representation of the test bag  $U$ , and  $\mathbf{M} = [\Phi(L_{corr}^1), \dots, \Phi(L_{corr}^n)]$  is the matrix containing the class-wise representations of the class-specific correction sets  $L_{corr}^i = \{\mathbf{x}_k : (\mathbf{x}_k, y_k) \in L_{corr}, y_k = i\}$ , for a given representation function  $\Phi : \mathbb{N}^X \rightarrow \mathbb{R}^z$  that embeds bags into  $z$ -dimensional vectors, for some  $z$ .

Most quantifiers rely on a representation function of the form:

$$\Phi(\mathbf{X}) = \frac{1}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} \phi(\mathbf{x}) \quad (7)$$

in which a surrogate instance-wise representation function  $\phi : \mathcal{X} \rightarrow \mathbb{R}^z$  is invoked, thus effectively computing a mean embedding. Here, we deliberately use  $\phi$  to denote both the representation function and a classifier, since most quantification methods define the representation function upon the output generated by a classifier.

Different choices for  $\Phi$  and  $\phi$  give rise to different instances of quantification methods. For example, ACC comes down to choosing, as our representation function  $\phi$ , the output of a crisp classifier encoded as a one-hot vector. The columns of  $\mathbf{M}$  thus represent the classification rates of the classier (as estimated on correction data), and the problem comes down to reconstructing the class counts of the test examples as a linear combination of the columns in  $\mathbf{M}$ , by solving  $\mathbf{p} = \mathbf{M}^{-1}\mathbf{t}$ .

More sophisticated methods exist. For example, *Probabilistic Adjusted Classify and Count* (PACC) (Bella et al. 2010) (that we use in our experiments), defines  $\phi$  as a probabilistic classifier returning the posterior probabilities for each class. When  $\mathbf{M}$  is not invertible (Bunse 2022), the variant we employ frames Equation (6) as the minimization problem

$$\mathbf{p}' = \arg \min_{\mathbf{p} \in \Delta^{n-1}} |\mathbf{t} - \mathbf{M}\mathbf{p}|^2 \quad (8)$$

We also use KDEy (Moreo, González, et al. 2025), a state-of-the-art multiclass quantification method that defines  $\Phi$  as a Gaussian mixture model, obtained via kernel density estimation, of the posterior probabilities returned by a probabilistic classifier. We use the maximum-likelihood variant that solves Equation (6) as the minimization problem

$$\mathbf{p}' = \arg \min_{\mathbf{p} \in \Delta^{n-1}} \mathcal{D}_{\text{KL}}(\mathbf{t} || \mathbf{M}\mathbf{p}) \quad (9)$$

where  $\mathcal{D}_{\text{KL}}$  is the well-known Kullback-Leibler divergence.

Note that PACC and KDEy both rely on a probabilistic classifier that is trained in advance on  $L_\phi$ . Equations (8) and (9) only require converting the items in  $L_q$  and the test items  $U_q$  (both retrieved for the same query but from different pools,  $L_{corr}$  and  $U$ , respectively) into posterior probabilities and solving the optimization problem. Both operations are rather fast given modern optimization routines (Section 5.4).

## 4 Experimental Setup

In this work, we aim to provide answers for the following research questions:

- **RQ1:** Does our new method improve over existing baselines?
- **RQ2:** Do quantification techniques allow for accurate QFE in a multiclass setting?

- **RQ3:** To what extent does the performance of our algorithm depend on the size of the correction pool and the rank exposure?

The code that implements all our proposed and all baseline methods, and reproduces our experimental results, is publicly available.<sup>3</sup>

#### 4.1 Dataset and Retrieval

To answer our research questions, we use the TREC 2022 Fair Ranking Track collection (Ekstrand, McDonald, et al. 2022). This dataset was established for the TREC Fair Ranking Track, in which the fairness of rankings produced by IR systems across different queries and various groups of documents with multiple classes is evaluated. It includes 6.5M English-language Wikipedia articles labelled with group information for various sensitive attributes. We choose this dataset over alternatives from the hiring domain since it is publicly available, it is large, and it encodes several multi-valued sensitive attributes.

To analyse the generalisation of our approach, we select multiple attributes for our evaluation, namely geographic location, gender, and age of topic. Standard classifiers such as logistic regression and SVM demonstrate reasonable accuracy ( $> 0.75$ ) for these attributes.

In addition to the documents and their demographic labels, the collection includes 97 queries. We index the documents using a Porter stemmer and stop word removal with the help of PyTerrier (Macdonald et al. 2021). As our retrieval model, we use BM25 (Robertson et al. 1994) with its standard parameters.

#### 4.2 Evaluation Measures

In our experiments, we assess the accuracy of QFE on multiclass groups. We concentrate on metrics that implement an exposure drop-off based on rank position, thereby reflecting the bias by which users tend to pay more attention to documents ranked higher (Craswell et al. 2008). As our metric of query fairness we rely on the normalized discounted Kullback-Leibler divergence (rKL) (Zehlike, K. Yang, et al. 2022) which, for a given set of retrieved documents  $U_q$  and different ranking levels  $k \in K$  (we consider  $K = \{50, 100, 500, 1000\}$ ), is given by

$$\text{rKL}(U_q) = \frac{1}{Z} \sum_{k \in K} \frac{1}{\log_2 k} \mathcal{D}_{\text{KL}}(\mathbf{p}^k \| \mathbf{p}^*) \quad (10)$$

where  $\mathbf{p}^k$  is the group distribution for the top- $k$  documents, and  $\mathbf{p}^*$  is the group distribution in all judged-relevant documents in the test collection  $U$  from which the ranked list  $U_q$  is retrieved.  $Z$  is simply the normalization factor computed as  $Z = \sum_{k \in K} 1/\log_2 k$ .

In Section 5.1 we will also test our methods against other baseline methods that are binary-only and that are tailored to one specific fairness metric called normalized discounted difference (rND) (Ghazimatin et al. 2022) defined by

$$\text{rND}(U_q) = \frac{1}{Z} \sum_{k \in K} \frac{1}{\log_2 k} \left| p_1^k - p_1^* \right| \quad (11)$$

Note that rND only considers the prevalence of the positive class, which is taken to be the prevalence of the protected or disadvantaged group.

Since we work under the unawareness assumption,  $\mathbf{p}^k$  is unknown and needs to be estimated. We thus denote by  $\hat{\text{rKL}}(U_q)$  (resp.,  $\hat{\text{rND}}(U_q)$ ) the score obtained using predicted distributions  $\hat{\mathbf{p}}^k$  in place of the true ones. In order to assess the accuracy on the prediction of the fairness metric  $M$  (be it rKL or rND), we report the *absolute error*

<sup>3</sup><https://github.com/AlexMoreo/query-fairness-estimation>

```

Input :: data
  • Classifier learner CLS
  • Quantification method QUANT
  • Group labels
Output :: rKL of the fairness estimates
  • RAE of the group prevalence estimates
1  $L, U \leftarrow \text{split}(\text{data}; 50\%, 50\%)$ 
2  $L_\phi \leftarrow \text{draw}(L; 500 \text{ documents per group})$ 
3  $L \leftarrow L - L_\phi$ 
4  $\phi \leftarrow \text{CLS}(L_\phi)$ 
5 for  $size \in \{10K, 50K, 100K, 500K, 1M, 3.25M\}$  do
6    $L_{size} \leftarrow \text{undersample}(L; size)$ 
7    $L_{corr} := L_{size}$  # alias
8   for  $query \in \text{Queries}$  do
9      $U_q \leftarrow \text{Retrieve}(U, query)@1000$ 
10     $L_q \leftarrow \text{Retrieve}(L_{corr}, query)@1000$ 
11     $L_q \leftarrow \text{keep}(L_q; \text{max 200 most relevant documents per group})$ 
12     $\lambda_h \leftarrow \text{QUANT}(L_q, h)$ 
13    for  $k \in \{50, 100, 500, 1000\}$  do
14       $\hat{p}^k \leftarrow \lambda_h(U_q^k)$ 
15       $p^k \leftarrow \text{prevalence}(U_q^k)$ 
16      compute RAE( $p^k, \hat{p}^k$ )
17    end
18  end
19  compute AE( $\text{Queries}$ , rKL)
20 end

```

**Algorithm 1:** Experimental protocol.

averaged across all ranked lists  $U_q$  retrieved for all queries ( $\text{Queries}$ ), which is defined as

$$\text{AE}(\text{Queries}, M) = \frac{1}{|\text{Queries}|} \sum_{U_q \in \text{Queries}} |M(U_q) - \hat{M}(U_q)| \quad (12)$$

We also report the relative absolute error (RAE), a quantification-specific measure that confronts a predicted distribution with the true distribution, and is defined as

$$\text{RAE}(\mathbf{p}, \hat{\mathbf{p}}) = \frac{1}{n} \sum_{i=1}^n \frac{|\hat{p}_i - p_i|}{p_i} \quad (13)$$

We choose RAE as it caters to minority classes by highlighting estimation errors that are small in absolute terms but proportionally large (Sebastiani 2020).

### 4.3 Experimental Protocol

In this section, we turn to describe the experimental protocol we have designed to provide answers to our RQs. The pseudocode describing our protocol can be consulted in Algorithm 1. The experimental variables we consider are listed below:

- *size*: The size of the correction pool  $L_{corr}$  from which the documents  $L_q$  are retrieved. This is important since there is an evident trade-off between cost and performance: a larger pool size implies higher labelling cost, while at the same time helps in reducing the discrepancy between the distribution of the documents retrieved

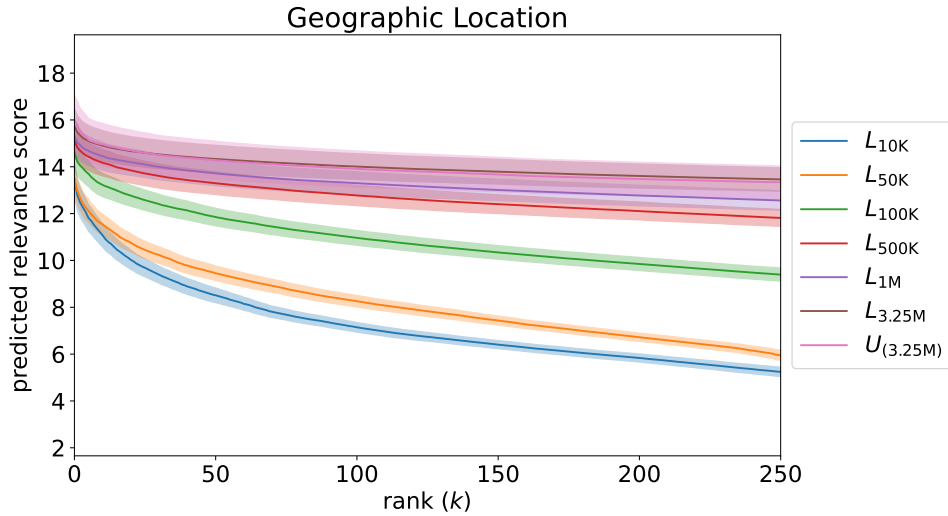


Fig. 2. Distribution of predicted relevance score per document rank across all queries. As  $|L_{corr}|$  increases its distribution aligns more closely with  $U$ .

from the training and test pool (Figure 2). In Line 5 we let *size* vary from a more realistic setting in which 10K labelled documents are assumed available, to the more optimistic (and unrealistic) scenario in which the correction pool is of the same size as the test pool (here corresponding to roughly 3.25M documents). We explore  $size \in \{10K, 50K, 100K, 500K, 1M, 3.25M\}$ .

- $k$ : The rank cutoff. We investigate the impact the rank of the top- $k$  examined has in the accuracy of QFE. We let  $k$  vary in the range  $K = \{50, 100, 500, 1000\}$  (Line 13).
- *Queries*: We assess the performance of our methods in QFE across all 97 queries available in the TREC collection (Line 8).

**4.3.1 Protocol Viewed from the Quantification Literature.** In quantification research, experimental evaluations often use a *sampling generation protocol* to simulate shifts in class distributions, providing a stress test for assessing a quantifier’s performance. Typically, these shifts are applied to the test set. In our case, it is not clear how to achieve this without interfering with SSB. Still, we deem it important to impose such a shift in the prior since, in our case, the labelled and unlabelled pools are obtained by partitioning one pre-existing collection. This results in the true underlying training and test distributions being nearly identical, which is an oversimplification of the problem. To avoid this oversimplification and test our quantifier under more realistic conditions, we introduce a shift by limiting the labelled data to a fixed number of documents per class (500 for  $L_\phi$  in Line 2, and 200 for  $L_q$  in Line 11).

For clarity within our protocol, we briefly summarise the quantification step. We first run a probabilistic classifier over the retrieved items to obtain, for each item, the probabilities of belonging to each group. Using a small labelled sample retrieved with the same query (the correction pool), we then summarise how these probabilities typically look for each group. Next, we select the mix of groups whose combined probabilities best matches those observed in the unlabelled set; this yields corrected estimates of the group proportions. Finally, we use these corrected proportions to compute the fairness metrics at the chosen top- $k$  cut-off.

#### 4.4 Methods

We experiment with our proposed variants:

- PACC (Bella et al. 2010): the “Probabilistic Adjusted Classify & Count” quantification method described in Section 3.4. Our proposed variant of PACC models the classification rates matrix of  $\phi$ , using  $L_q$ .
- KDEy (Moreo, González, et al. 2025): the KDE-based quantification method described in Section 3.4. Our proposed variant of KDEy models the class-wise densities of the posterior probabilities of  $\phi$ , using  $L_q$ .

We compare our methods to the following baselines:

- Naive@ $k$ : a method that does not inspect the search results at all to estimate the prevalence of  $U_q^k$  but simply reports the prevalence of  $L_q^k$ , i.e., of the top- $k$  training documents retrieved for each query from  $L_{corr}$ .
- CC (Forman 2005): the “Classify & Count” method described in Section 3 and Equation (1). This method relies on the predictions issued by  $\phi$ , and does not learn any correction based on  $L_q$ .
- $PMC_b$  (Ghazimatin et al. 2022): a binary-only method that corrects a preliminary estimate  $r\hat{ND}_\phi$  obtained using “proxy labels” (i.e., labels predicted by  $\phi$ ), by applying equation:

$$rND(U_q) = \frac{r\hat{ND}_\phi}{(1-p) - w} \quad (14)$$

where  $p = P_L(\hat{Y} = 1|Y = 0)$  and  $w = P_L(\hat{Y} = 0|Y = 1)$  (Ghazimatin et al. 2022, Figure 1(b)).

- $PMC_b^+$ : a variant we propose for  $PMC_b$  in which the correction factors  $p$  and  $w$  are not modelled on the same dataset  $L$  used to train the classifier (as proposed by the inventors of the method), but from the items  $L_q$  retrieved for each query.
- $PMC_d$  (Ghazimatin et al. 2022): a binary-only method that corrects a preliminary estimate  $r\hat{ND}_h$  by applying equation:

$$rND(U_q) = r\hat{ND}_h \cdot \left( \frac{(1-w) \cdot \beta}{x} - \frac{w \cdot \beta}{y} \right) \quad (15)$$

where  $p$  and  $w$  are defined as for  $PMC_b$ ,  $\beta = P_L(Y = 1)$ , and  $x = (1-w) \cdot \beta + p \cdot (1-\beta)$  and  $y = w \cdot \beta + (1-p) \cdot (1-b)$  (Ghazimatin et al. 2022, Figure 1(d)).

- $PMC_d^+$ : a variant we propose for  $PMC_d$  in which  $p$ ,  $w$ ,  $\beta$ ,  $x$ , and  $y$  are not modelled on  $L$  but on  $L_q$ .

Our implementations of quantification methods CC, PACC, and KDEy rely on the implementations available in QuaPy (Moreo, Esuli, et al. 2021), while the implementations of Naive@ $k$ ,  $PMC_b$ ,  $PMC_b^+$ ,  $PMC_d$ , and  $PMC_d^+$  are our own.

The choice of baselines reflects the need to evaluate our approach against both elementary strategies and the most relevant prior work. Naïve@ $k$  and CC are included as simple counting-based strategies that establish a lower bound on performance in the absence of correction. In contrast, PMC and its variants constitute the strongest available comparisons, as they also introduce corrections based on classifier outputs but are restricted to binary attributes and require metric-specific adaptations. Together, these baselines span the range from simple reference methods to state-of-the-art approaches, ensuring that the empirical evaluation is both comprehensive and properly contextualised

**Classifier training.** All the methods we consider in this work, with the exception of Naive@ $k$ , rely on the outputs of a classifier. For the sake of a fair comparison, we use the same classifier in all cases. Our classifier of choice is Logistic Regression (LR), which is arguably becoming a *de facto* choice in the field of quantification, due to the fact that many methods require probabilistic decisions and LR is known to deliver reasonably well-calibrated posterior probabilities (Moreo, Esuli, et al. 2021; Schumacher et al. 2025). LR is trained offline, once and for all, for each category, on  $L_\phi$ .

**Model selection.** All quantification methods in our evaluation require a probabilistic base classifier. We adopt logistic regression (LR), as it is widely used in quantification research due to its interpretability and efficiency at scale. Hyperparameters are tuned via 5-fold cross-validation (5-FCV) on the labelled pool  $L_\phi$ . We explore the regularisation strength  $C \in \{10^i\}$  for  $-4 \leq i \leq 4$ , a logarithmic grid that balances coverage with computational tractability, and the `CLASS_WEIGHT` parameter in `{BALANCED, NONE}`, in order to assess whether explicit reweighting improves prevalence estimation. The configuration achieving the lowest cross-validated log-loss is retained, ensuring well-calibrated posterior probabilities for downstream quantification.

The only other method with additional hyperparameters is KDEy, which depends on the kernel “bandwidth”. Since query-time optimisation is infeasible, we pre-select this parameter using a development set of 100 queries from the TREC 2021 Fair Ranking Track (Ekstrand, McDonald, et al. 2021), a precursor of the TREC 2022 collection employed in our main experiments. This choice ensures that bandwidth selection is informed by queries similar in nature to those used in evaluation, while avoiding overlap with the test set. To ensure sufficient coverage of all attributes and classes during this tuning step, we issue the queries on a 100K-sized correction pool, which provides robust sampling for each group. We evaluate bandwidth values in the range  $\{0.01, 0.02, \dots, 0.10\}$ , which represents a suitable range when kernel density estimation is computed on a probability simplex and that captures both narrow and moderately smooth kernels. The selected bandwidth is then fixed for all subsequent experiments, providing a consistent and reproducible parametrisation.

## 5 Results

In this section, we report and discuss the experimental results we have obtained. Section 5.1 presents a comparison against the PMC variants (discussed in Section 2); this experimental comparison is discussed separately because the above-mentioned models are binary-only (RQ1). In Section 5.2 we show our main set of experiments, in which we assess the effectiveness of QFE considering multiclass groups (RQ2). Section 5.3 analyses the extent to which the performance of our proposed methods depends on the rank  $k$  and the correction pool size (RQ3). Finally, Section 5.4 reports averaged time measurements of our methods.

### 5.1 Binary Protected Attributes

In this section, we compare our proposed approach to previous related work (RQ1). The results we discuss in the next paragraphs correspond to the more realistic scenario in which  $size = 10K$ . We compare our proposed methods (PACC and KDEy) against the PMC variants (Ghazimatin et al. 2022) Section 4.4.

The PMC methods apply a post-correction to the fairness metric score, making them metric-specific and binary-only. In order to allow for an experimental comparison against the PMC variants, we produce binary versions of our datasets. In the binary setting, the positive class ( $Y = 1$ ) traditionally represents the minority or disadvantaged group. We thus binarize our datasets towards the following groups: “Africa” for Geographic Location, “Female” for Gender, and “Pre-1900s” for Age of Topic. the rest of the groups are merged into the negative class ( $Y = 0$ ).

Table 1 reports the absolute error in the estimation of rND. The displayed values are averaged scores of the absolute error on the prediction of rKL (lower is better) across all 97 queries. Boldface indicates the best method for a given category. Superscripts † and ‡ indicate the methods (if any) whose scores are *not* statistically significantly different from the best one at different confidence levels: symbol † indicates  $0.001 < p < 0.01$ , while symbol ‡ indicates  $p \geq 0.01$ . As the test for statistical significance, we rely on the non-parametric Wilcoxon signed-rank test. We use colour coding to facilitate the interpretation of the results, with green indicating the best result and red indicating the worst one per category. From the analysis of the results the following observations can be drawn:

- KDEy is the best-performing approach for the three categories.

Table 1. Results of QFE in terms of AE (lower is better) of rND prediction for  $L_{10K}$  (binary case). KDEy beats all methods.

	Naive@ $k$	CC	PACC (ours)	KDEy (ours)	PMC <sub>b</sub>	PMC <sub>b</sub> <sup>+</sup>	PMC <sub>d</sub>	PMC <sub>d</sub> <sup>+</sup>
Geographic Location	.014 <sup>†</sup> ± .023	.013 <sup>‡</sup> ± .029	.043 <sup>‡</sup> ± .143	<b>.012 ± .027</b>	.020 <sup>†</sup> ± .049	.030 ± .087	.020 <sup>†</sup> ± .049	.031 ± .083
Gender	.047 ± .048	.014 <sup>‡</sup> ± .042	.016 <sup>‡</sup> ± .040	<b>.011 ± .026</b>	.014 <sup>‡</sup> ± .041	.015 <sup>‡</sup> ± .037	.014 <sup>‡</sup> ± .041	.015 <sup>‡</sup> ± .037
Age of Topic	.040 ± .040	.025 <sup>‡</sup> ± .040	.019 <sup>†</sup> ± .023	<b>.017 ± .021</b>	.028 ± .040	.042 ± .068	.028 ± .040	.044 ± .063

- Although CC performs consistently worse than KDEy, the statistical test reveals these differences are not significant. This may be due to high classifier accuracy in the binary case, leaving small room for improvement for the correction phase of other methods. Indeed CC fares significantly worse than KDEy in the multi-class settings (Section 5.2).
- PMC models perform worse than KDEy in a statistically significant sense in most cases. This is an indication that the assumptions upon which PMC models are built do not apply in QFE.
- The variants PMC<sup>+</sup> we propose fare consistently worse than the original methods. Intuitively, these variants should perform better, since the documents on which the correction is modelled are more similar to the test data for which the correction is required. This may be an indication that the post correction implemented by the PMC variants do not align with the characteristics of the distributions under consideration in QFE.
- Naive@ $k$  performs badly in two out of three cases. This speaks in favour of the ability of KDEy to correct the class prevalence values, since the class prevalence of the top- $k$  training documents is not a good estimate for the top- $k$  test documents *per se*.
- PACC falls short in terms of performance. A plausible reason for this failure is the relatively low number of training documents used to model the classification rates (more on this in Section 5.3).

## 5.2 Multiclass Protected Attributes

We now turn to query fairness estimation for multiclass sensitive attributes. Table 2 reports the absolute error of rKL estimates (Equation 10) for a realistic scenario where the dataset size is set to  $size = 10K$ . Notational conventions are as in Table 1. Since ours is the first multiclass method, we compare our estimator against quantification baselines. The following observations emerge from our results:

- KDEy consistently achieves the best performance compared to all other baselines, with statistically significant differences in the majority of cases across the evaluated attributes. Furthermore, KDEy also exhibits the smallest standard deviation, indicating consistent performance.
- Naive@ $k$  performs erratically. In Geographic Location, it achieves results similar to the best-performing method, KDEy. However, in the Gender category, it produces significantly higher errors compared to the best performer. This inconsistent performance is reflected in a high standard deviation across all categories.
- CC performs consistently worse than KDEy in all cases. The differences are statistically significant at  $p = 0.01$ ; in one out of three cases they are significant at  $p = 0.001$ .
- PACC is, as in the binary case, not competitive with KDEy.

Table 2. Results of QFE in terms of AE (lower is better) of rKL prediction for  $L_{10K}$  (multiclass case). KDEy beats all methods.

	Naive@ $k$	CC	PACC (ours)	KDEy (ours)
Geo. Location	.188 <sup>†</sup> ± .255	.212 ± .246	.298 ± .241	<b>.132 ± .181</b>
Gender	.305 ± .356	.068 <sup>†</sup> ± .143	.064 ± .108	<b>.037 ± .060</b>
Age of Topic	.213 <sup>†</sup> ± .265	.173 <sup>†</sup> ± .219	.285 ± .321	<b>.129 ± .158</b>

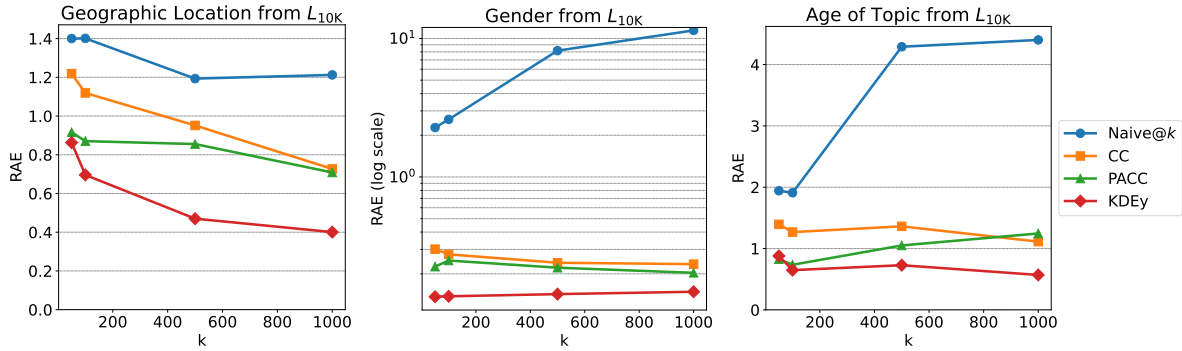


Fig. 3. Variations in quantification performance (measured in terms of RAE – lower is better) at different values of  $k$ . Note the log-scale in Gender. KDEy and PACC outperform all baselines.

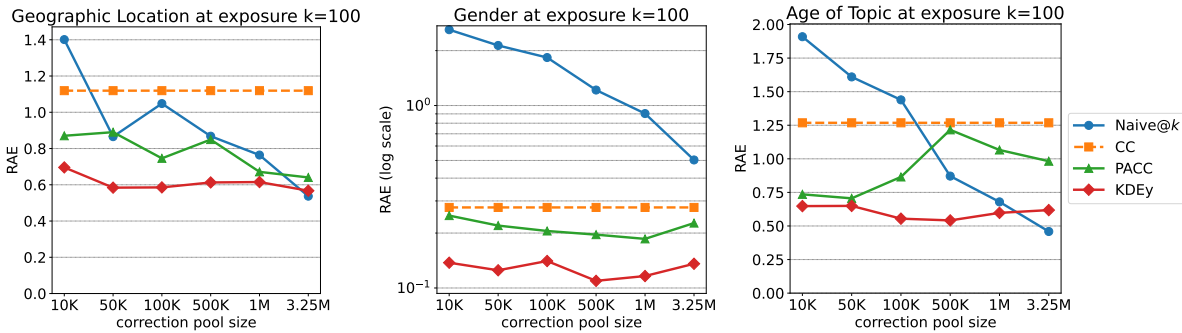


Fig. 4. Variations in quantification performance (measured in terms of RAE – lower is better) at different correction pool sizes. Note the log-scale in Gender. KDEy and PACC are more stable than Naive@ $k$  and dominate CC.

The disparate outcomes we have obtained for PACC and KDEy deserve further analysis. Both methods rely on the same principle of deferring the correction-training phase at query time. As we will see in Section 5.3, though, PACC still performs decently in terms of quantification performance. Concerning RQ2 and in light of our observations, we can conclude that quantification techniques are indeed suitable for accurate QFE in a multiclass scenarios.

### 5.3 Variations of $k$ and Size

In this experiment, we analyse the variations in quantification performance at different exposure levels  $k$  and variations in the correction pool size (RQ3). We evaluate these experiments in terms of RAE (a quantification-specific measure) between the true distribution and the predicted distribution. Figure 3 displays variations in performance at different rank levels ( $k \in \{50, 100, 500, 1000\}$ ) for the case  $L_{corr} := L_{10K}$ , while Figure 4 displays variations in performance due to variations in the correction pool size at rank  $k = 100$ .

These plots reveal some interesting findings. First, PACC performs consistently better than CC in terms of quantification. This comes as a surprise, since the experiments reported in Table 1 and Table 2 seemed to indicate PACC performs worse than CC.

The key difference with respect to our previous experiments is the evaluation measure under consideration. We argue that RAE is an appropriate measure in multiclass scenarios like the one we are facing here, given its ability to reflect the importance of an error with respect to the proportion of the true prevalence. This can be especially important in cases where one of the classes (the disadvantaged group) tends to display very low prevalence with respect to other classes (especially the privileged group).

However, other evaluation metrics for QFE (e.g., rKL and rND) do not seem to align with the intuitions behind RAE.<sup>4</sup> In the future, it will be interesting to analyse the adequacy of a normalized-discounted-variant of RAE for fairness evaluation. We leave these considerations to future investigations.

Figure 3 shows KDEy performs better than all other methods. Almost all methods show a tendency to improve as the exposure level increases. This is expected, as the quality of an estimated descriptive statistic (in this case: the prevalence) is known to depend on the size of the population under investigation. The only exception to this trend is Naive@ $k$ . The reason is that we actively tried to hide the original distribution (not only for Naive@ $k$ , but for all methods) by keeping no more than 200 documents per group in  $L_q$  (Line 11 in Algorithm 1). Moreover, Naive@ $k$  and CC performing consistently worse than PACC and KDEy proves that our quantifiers are effectively learning a correction for the class counts of the classifier which is not spuriously based on the class distribution of  $L_q$ .

Figure 4 also shows that Naive@ $k$  has a clear dependency on the distributional similarity between  $L_{corr}$  and  $U$ , as witnessed by its drastic improvement when adding labelled data to the correction pool  $L_{corr}$ , which makes it converge towards the same distribution as  $U$ . Conversely, the quantification performance of PACC and KDEy is relatively stable, and does not improve markedly when the amount of labelled data available in the pool increases. This indicates that our quantification-based methods are robust to drifts between  $L_{corr}$  and  $U$ . Moreover, this implies that the amount of labelling effort required to achieve reliable QFE scores is kept under reasonable bounds. Note also that the CC method is represented as a flat curve. The reason is that CC does not leverage the data available in  $L_{corr}$  by any means.

## 5.4 Time Measurements

We measured the training and testing times of our Python implementations on a desktop computer equipped with a 12th Gen Intel(R) i9-12900K processor and 64GB of RAM, running Ubuntu 22. Our methods consist of a per-query correction learning phase followed by a phase of prevalence predictions. On average, PACC required 0.907 ms for learning and 3.316 ms for every prediction, while KDEy took 1.586 ms for learning and 6.395 ms for every prediction. Despite the requirement of a learning phase at query time, our methods demonstrate relatively fast performance and can be integrated into the standard IAS pipeline.

## 6 Discussion and Conclusion

In this work, we have investigated how to reliably assess the fairness of search results when demographic labels for ranked items are unavailable. We have demonstrated that simply counting over the predictions of a classifier leads to unreliable fairness assessments. To address this limitation, we have proposed novel quantification-based methods to accurately estimate the prevalence of different groups in a ranking. The experimental evaluation has shown that our approach can successfully predict query fairness, including in the previously unaddressed multiclass case, and it does so more accurately than existing methods in the binary case. While most quantification techniques are designed to counter prior probability shift, the problem at hand is instead mainly affected by sample selection bias. To the best of our knowledge, our approach is the first attempt towards making quantification robust to this type of shift naturally occurring in QFE.

<sup>4</sup>Note that the KL-divergence in rKL also considers the ratio with respect to the target distribution. However, note that it also scales this factor by the predicted prevalence, which might simply be very close to 0, thus cancelling out the term in the aggregation; see Equation (10).

**Limitations.** QFE may be challenging to integrate into learning-to-rank pipelines. First, leveraging protected attributes correction pools comes with infrastructural challenges of data integration. Second, conveying uncertainty of fairness estimates is an important problem we did not address in this work. Third, correction pools may exhibit particular properties (e.g., self-selection effects) that we did not explicitly model. However, the correction pools can easily be maintained as moderate-sized labelled resources, in line with how IR systems already employ benchmark collections for evaluation. Because they are only required for prevalence estimation, their computational footprint is low, and our results indicate that reliable estimates can be obtained with pools of manageable size. This makes correction pools a practical, though necessarily limited, instrument for fairness, for example when auditing in real-world systems

**Future work.** In future work, we aim to investigate the suitability of a normalised-discounted variant of RAE for fairness evaluation. We are also interested in exploring different collections where the group prevalence may have naturally varied across training and test conditions. Finally, we aim to endow our QFE solutions with the ability to provide confidence intervals for the estimated values.

## Acknowledgments

The work of AM, AE, FS, was partially supported by projects “Future Artificial Intelligence Research” (FAIR), project “Quantification under Dataset Shift” (QuaDaSh), and project “Strengthening the Italian RI for Social Mining and Big Data Analytics” (SoBigData.it), all funded by the European Union under the NextGenerationEU funding scheme (CUP B53D22000980006, CUP B53D23026250001, CUP B53C22001760006, respectively). The work of GM and IO has been supported in part by the Engineering and Physical Sciences Research Council grant number EP/Y009800/1, through funding from Responsible Ai UK (KP0011).

## References

- D. Adjaye-Gbewonyo, R. A. Bednarczyk, R. L. Davis, and S. B. Omer. 2014. “Using the Bayesian improved surname geocoding method (BISG) to create a working classification of race and ethnicity in a diverse managed care population: A validation study.” *Health Services Research*, 49, 1, 268–283.
- A. Bella, C. Ferri, J. Hernández-Orallo, and M. J. Ramírez-Quintana. 2010. “Quantification via probability estimators.” In: *Proceedings of the 11th IEEE International Conference on Data Mining (ICDM 2010)*. Sydney, AU, 737–742. doi:10.1109/icdm.2010.75.
- A. J. Biega, K. P. Gummadi, and G. Weikum. 2018. “Equity of attention: Amortizing individual fairness in rankings.” In: *Proceedings of the 41st ACM Conference on Research and Development in Information Retrieval (SIGIR 2018)*. Ann Arbor, US, 405–414.
- M. Bogen, A. Rieke, and S. Ahmed. 2020. “Awareness in practice: Tensions in access to sensitive attribute data for antidiscrimination.” In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT\* 2020)*. Barcelona, ES, 492–500.
- M. Bunse. 2022. “On multi-class extensions of adjusted classify and count.” In: *Proceedings of the 2nd International Workshop on Learning to Quantify (LQ 2022)*. Grenoble, IT, 43–50.
- M. Bunse and K. Morik. 2022. “Unification of algorithms for quantification and unfolding.” In: *Proceedings of the 2nd International Workshop on Learning to Quantify (LQ 2022)*. Grenoble, IT, 1–10.
- L. E. Celis, L. Huang, V. Keswani, and N. K. Vishnoi. 2021. “Fair classification with noisy protected attributes: A framework with provable guarantees.” In: *Proceedings of the 38th International Conference on Machine Learning (ICML 2021)*. Virtual Event, 1349–1361.
- F. Chen and H. Fang. 2023. “Learn to be fair without labels: A distribution-based learning framework for fair ranking.” In: *Proceedings of the 46th ACM Conference on Research and Development in Information Retrieval (SIGIR 2023)*. Taipei, TW, 23–32.
- J. Chen, N. Kallus, X. Mao, G. Svacha, and M. Udell. 2019. “Fairness under unawareness: Assessing disparity when protected class is unobserved.” In: *Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency (FAT\* 2019)*. Atlanta, US, 339–348.
- L. Chen, R. Ma, A. Hannák, and C. Wilson. 2018. “Investigating the impact of gender on rank in resume search engines.” In: *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI 2018)*. Montreal, CA, 651.
- N. Craswell, O. Zoeter, M. Taylor, and B. Ramsey. 2008. “An experimental comparison of click position-bias models.” In: *Proceedings of the 1st ACM Conference on Web Search and Web Data Mining (WSDM 2008)*. Palo Alto, US, 87–94.
- F. Diaz, B. Mitra, M. D. Ekstrand, A. J. Biega, and B. Carterette. 2020. “Evaluating stochastic rankings with expected exposure.” In: *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM 2020)*. Virtual Event, 275–284.
- M. D. Ekstrand, A. Das, R. Burke, and F. Diaz. 2022. “Fairness in information access systems.” *Foundations and Trends in Information Retrieval*, 16, 1-2, 1–177.

- M. D. Ekstrand, G. McDonald, A. Raj, and I. Johnson. 2021. "Overview of the TREC 2021 Fair Ranking Track." In: *Proceedings of the 30th Text REtrieval Conference (TREC 2021)*. Virtual Event].
- M. D. Ekstrand, G. McDonald, A. Raj, and I. Johnson. 2022. "Overview of the TREC 2022 Fair Ranking Track." In: *Proceedings of the 31st Text REtrieval Conference (TREC 2022)*. Virtual Event.
- A. Esuli, A. Fabris, A. Moreo, and F. Sebastiani. 2023. *Learning to quantify*. Springer Nature, Cham, CH. doi:10.1007/978-3-031-20467-8.
- European Commission. 2024. *Regulation on harmonized rules on artificial intelligence (Artificial Intelligence Act)*. Official Journal of the European Union. Accessed: 2024-08-06. (2024).
- A. Fabris, A. Esuli, A. Moreo, and F. Sebastiani. 2023. "Measuring fairness under unawareness of sensitive attributes: A quantification-based approach." *Journal of Artificial Intelligence Research*, 76, 1117–1180. doi:10.1613/jair.1.14033.
- G. Forman. 2005. "Counting positives accurately despite inaccurate classification." In: *Proceedings of the 16th European Conference on Machine Learning (ECML 2005)*. Porto, PT, 564–575. doi:10.1007/11564096\_55.
- S. A. Friedler, C. Scheidegger, and S. Venkatasubramanian. 2021. "The (im)possibility of fairness: Different value systems require different mechanisms for fair decision making." *Communications of the ACM*, 64, 4, 136–143.
- S. C. Geyik, S. Ambler, and K. Kenthapadi. 2019. "Fairness-aware ranking in search and recommendation systems with application to LinkedIn talent search." In: *Proceedings of the 25th ACM International Conference on Knowledge Discovery and Data Mining (KDD 2019)*. Anchorage, US, 2221–2231.
- A. Ghazimatin, M. Kleindessner, C. Russell, Z. Abedjan, and J. Golebiowski. 2022. "Measuring fairness of rankings under noisy sensitive information." In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAcT 2022)*. Seoul, KR, 2263–2279.
- A. Ghosh, R. Dutt, and C. Wilson. 2021. "When fair ranking meets uncertain inference." In: *Proceedings of the 44th ACM Conference on Research and Development in Information Retrieval (SIGIR 2021)*. Virtual Event, 1033–1043.
- A. Ghosh, P. Kvitca, and C. Wilson. 2023. "When fair classification meets noisy protected attributes." In: *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society (AIES 2023)*. Montreal, CA, 679–690.
- P. González, A. Moreo, and F. Sebastiani. 2024. "Binary quantification and dataset shift: An experimental investigation." *Data Mining and Knowledge Discovery*, 38, 4, 1670–1712. doi:10.1007/s10618-024-01014-1.
- M. Hardt, E. Price, and N. Srebro. 2016. "Equality of opportunity in supervised learning." In: *Proceedings of the 29th Annual Conference on Neural Information Processing Systems (NIPS 2016)*. Barcelona, ES, 3315–3323.
- M. Heuss, F. Sarvi, and M. de Rijke. 2022. "Fairness of exposure in light of incomplete exposure estimation." In: *Proceedings of the 45th ACM Conference on Research and Development in Information Retrieval (SIGIR 2022)*. Madrid, ES, 759–769.
- K. Holstein, J. Wortman Vaughan, H. Daumé III, M. Dudik, and H. Wallach. 2019. "Improving fairness in machine learning systems: What do industry practitioners need?" In: *Proceedings of the 2019 ACM Conference on Human Factors in Computing Systems (CHI 2019)*. Glasgow, UK, 600.
- T. Jaenich, G. McDonald, and I. Ounis. 2023. "CoBERT-FairPRF: Towards fair pseudo-relevance feedback in dense retrieval." In: *Proceedings of the 45th European Conference on Information Retrieval (ECIR 2023)*. Dublin, IE, 457–465.
- T. Jaenich, G. McDonald, and I. Ounis. 2024. "Fairness-aware exposure allocation via adaptive reranking." In: *Proceedings of the 47th ACM Conference on Research and Development in Information Retrieval (SIGIR 2024)*. Washington, US, 1504–1513.
- Ö. Kirnap, F. Diaz, A. Biega, M. D. Ekstrand, B. Carterette, and E. Yilmaz. 2021. "Estimation of fair ranking metrics with incomplete judgments." In: *Proceedings of the 2021 Web Conference (WWW 2021)*. Virtual Event, 1065–1075.
- S. Kopeinik, M. Mara, L. Ratz, K. Krieg, M. Schedl, and N. Rekabsaz. 2023. "Show me a "male nurse"! How gender bias is reflected in the query formulation of search engine users." In: *Proceedings of the 2023 ACM Conference on Human Factors in Computing Systems (CHI 2023)*. Hamburg, DE, 137:1–137:15.
- C. Kuhlman, W. Gerych, and E. Rundensteiner. 2021. "Measuring group advantage: A comparative study of fair ranking metrics." In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (AIES 2021)*. Virtual Event, 674–682.
- LinkedIn. 2024. *LinkedIn settings*. New York City Council Legislation. Accessed: 2024-08-06. (2024). <https://www.linkedin.com/mypreferences/d/demographic-info>.
- Z. C. Lipton, Y. Wang, and A. J. Smola. 2018. "Detecting and correcting for label shift with black box predictors." In: *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*. Stockholm, SE, 3128–3136.
- C. Macdonald, N. Tonello, S. MacAvaney, and I. Ounis. 2021. "PyTerrier: Declarative experimentation in Python from BM25 to dense retrieval." In: *Proceedings of the 30th ACM International Conference on Knowledge Management (CIKM 2021)*. Gold Coast, AU, 4526–4533.
- A. Moreo, A. Esuli, and F. Sebastiani. 2021. "QuaPy: A Python-based framework for quantification." In: *Proceedings of the 30th ACM International Conference on Knowledge Management (CIKM 2021)*. Gold Coast, AU, 4534–4543. doi:10.1145/3459637.3482015.
- A. Moreo, P. González, and J. J. del Coz. 2025. "Kernel density estimation for multiclass quantification." *Machine Learning*, 114, 4, 38 pages. doi:10.1007/s10994-024-06726-5.
- M. Morik, A. Singh, J. Hong, and T. Joachims. 2020. "Controlling fairness and bias in dynamic learning-to-rank." In: *Proceedings of the 43rd ACM Conference on Research and Development in Information Retrieval (SIGIR 2020)*. Virtual Event, 429–438.

- H. Mozannar, M. Ohanessian, and N. Srebro. 2020. "Fair learning with private demographic data." In: *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*. Virtual Event, 7066–7075.
- New York City Council. 2021. *Local law 144 of 2021*. New York City Council Legislation. Accessed: 2024-08-06. (2021). <https://legistar.council.nyc.gov/LegislationDetail.aspx?ID=4344524&GUID=B051915D-A9AC-451E-81F8-6596032FA3F9>.
- D. Pedreschi, S. Ruggieri, and F. Turini. 2008. "Discrimination-aware data mining." In: *Proceedings of the 14th ACM International Conference on Knowledge Discovery and Data Mining (KDD 2008)*. Las Vegas, US, 560–568.
- A. Raj and M. D. Ekstrand. 2022. "Measuring fairness in ranked results: An analytical and empirical comparison." In: *Proceedings of the 45th ACM Conference on Research and Development in Information Retrieval (SIGIR 2022)*. Madrid, ES, 726–736.
- S. E. Robertson, S. Walker, S. Jones, M. M. Hancock-Beaulieu, and M. Gatford. 1994. "Okapi at TREC-3." In: *Proceedings of the 2nd Text Retrieval Conference (TREC 1993)*. Gaithersburg, US, 109–126.
- P. Sapiezynski, W. Zeng, R. E. Robertson, A. Mislove, and C. Wilson. 2019. "Quantifying the impact of user attention on fair group representation in ranked lists." In: *Companion Proceedings of the 2019 World Wide Web Conference (WWW 2019)*. San Francisco, US, 553–562. doi:10.1145/308560.3317595.
- B. Schölkopf, D. Janzing, J. Peters, E. Sgouritsa, K. Zhang, and J. M. Mooij. 2012. "On causal and anticausal learning." In: *Proceedings of the 29th International Conference on Machine Learning (ICML 2012)*. Edinburgh, UK.
- T. Schumacher, M. Strohmaier, and F. Lemmerich. 2025. "A comparative evaluation of quantification methods." *Journal of Machine Learning Research*, 26, 55, 1–54.
- F. Sebastiani. 2020. "Evaluation measures for quantification: An axiomatic approach." *Information Retrieval Journal*, 23, 3, 255–288. doi:10.1007/s10791-019-09363-y.
- J. Simson, A. Fabris, and C. Kern. 2024. "Lazy data practices harm fairness research." In: *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency (FAccT 2024)*. Rio de Janeiro, BR, 642–659. doi:10.1145/3630106.3658931.
- A. Singh and T. Joachims. 2018. "Fairness of exposure in rankings." In: *Proceedings of the 24th ACM International Conference on Knowledge Discovery and Data Mining (KDD 2018)*. London, UK, 2219–2228.
- A. Storkey. 2009. "When training and test sets are different: Characterizing learning transfer." In: *Dataset shift in machine learning*. Ed. by J. Quiñero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence. The MIT Press, Cambridge, US, 3–28.
- S. Wang, W. Guo, H. Narasimhan, A. Cotter, M. Gupta, and M. Jordan. 2020. "Robust optimization for fairness with noisy protected groups." In: *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems (NIPS 2020)*. Virtual Event, 5190–5203.
- Y. Wang, W. Ma, M. Zhang, Y. Liu, and S. Ma. 2023. "A survey on the fairness of recommender systems." *ACM Transactions on Information Systems*, 41, 3, 1–43. doi:10.1145/3547333.
- C. Wilson, A. Ghosh, S. Jiang, A. Mislove, L. J. Baker, J. Szary, K. Trindel, and F. Polli. 2021. "Building and auditing fair algorithms: A case study in candidate screening." In: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT 2021)*. Toronto, CA, 666–677. doi:10.1145/3442188.3445928.
- E. Yang, T. Jänich, J. Mayfield, and D. Lawrie. 2024. "Language fairness in multilingual information retrieval." In: *Proceedings of the 47th ACM Conference on Research and Development in Information Retrieval (SIGIR 2024)*. Washington, US, 2487–2491.
- K. Yang and J. Stoyanovich. 2017. "Measuring fairness in ranked outputs." In: *Proceedings of the 29th International Conference on Scientific and Statistical Database Management (SSDBM 2017)*. Chicago, US, 22:1–22:6.
- M. Zehlike, F. Bonchi, C. Castillo, S. Hajian, M. Megahed, and R. Baeza-Yates. 2017. "FA\*IR: A fair top-k ranking algorithm." In: *Proceedings of the 26th ACM International Conference on Knowledge Management (CIKM 2017)*. Singapore, SN, 1569–1578.
- M. Zehlike, K. Yang, and J. Stoyanovich. 2022. "Fairness in ranking, Part I: Score-based ranking." *ACM Computing Surveys*, 55, 6, 1–36.

Received 20 November 2024; accepted 30 August 2025