

# Characterizing Intersectional Group Fairness with Worst-Case Comparisons

Avijit Ghosh,<sup>1,2</sup> Lea Genuit,<sup>1</sup> Mary Reagan<sup>1</sup>

<sup>1</sup>Fiddler Labs <sup>2</sup>Northeastern University  
avijit@ccs.neu.edu, lea@fiddler.ai, mary@fiddler.ai

## Abstract

Machine Learning or Artificial Intelligence algorithms have gained considerable scrutiny in recent times owing to their propensity towards imitating and amplifying existing prejudices in society. This has led to a niche but growing body of work that identifies and attempts to fix these biases. A first step towards making these algorithms more fair is designing metrics that measure unfairness. Most existing work in this field deals with either a binary view of fairness (protected vs. unprotected groups) or politically defined categories (race or gender). Such categorization misses the important nuance of intersectionality - biases can often be amplified in subgroups that combine membership from different categories, especially if such a subgroup is particularly underrepresented in historical platforms of opportunity.

In this paper, we discuss why fairness metrics need to be looked at under the lens of intersectionality, identify existing work in intersectional fairness, suggest a simple worst case comparison method to expand the definitions of existing group fairness metrics to incorporate intersectionality, and finally conclude with the social, legal and political framework to handle intersectional fairness in the modern context.

## 1 Introduction

The use of machine learning algorithms is ubiquitous in the developed world. It has become an integral part of society, affecting the lives of millions of people. Algorithmic decisions vary from low-stakes determinations, like product or film recommendations, to high-impact like loan or credit approval (Mukerjee et al. 2002), hiring recommendations (Bogen and Rieke 2018), facial recognition (Vasilescu and Terzopoulos 2002) and prison recidivism (Corbett-Davies and Goel 2018). With this direct impact on people’s lives, the need for fair and unbiased algorithms is paramount. It is critical that algorithms do not replicate and enhance existing societal biases, including those rooted in differences of race, gender, or sexual orientation.

To tackle these problems, both fairness and bias need to be clearly defined. Currently, there does not exist a single universally agreed upon definition of fairness. Anti-discrimination legislation exists in various jurisdictions

around the world. In the US, anti-discrimination laws exist under the Civil Rights Act (Berg 1964), and under specific areas like credit lending<sup>1</sup> and housing<sup>2</sup>. There have also been efforts to introduce legislation combating algorithmic bias<sup>3</sup>. In the European Union, the GDPR law provides for regulations regarding digital profiling, data collection, and a “right to explanation” (Goodman and Flaxman 2017). Under Indian law, quotas for scheduled castes, scheduled tribes and other backward classes are mandated in public education and government employment.<sup>4</sup>

We begin with the broad definition of fairness as “the absence of prejudice or preference for an individual or group based on their characteristics”. Bias can also exist in a variety of forms. Mehrabi et al. (2019) provides an excellent overview on the differing types of bias and discrimination. In general, a fair machine learning algorithm is one that does not favor or make prejudice towards an individual or a group.

While most early fairness research focused on binary fairness metrics (protected vs. unprotected groups), newer methods to address fairness have begun to incorporate intersectional frameworks. These frameworks are derived from the third wave of feminist thought, which is rooted in the understanding of the interconnected nature of social categories, like race, gender, sexual orientation, and class (Crenshaw 1989). The intersection of these categories creates differing levels of privilege or disadvantage for the various possible subgroups. There exist legal precedents for discrimination under an intersectional lens : The Equal Employment Opportunity Commission (EEOC) describes some Intersectional Discrimination/Harassment examples<sup>5</sup>. Buolamwini and Gebru (2018) examined gender classification algorithms for facial image data and found that they performed substantially better on male faces than female faces. However, the largest performance drops came when both race and gender were considered, with

<sup>1</sup><https://www.justice.gov/crt/equal-credit-opportunity-act-3>

<sup>2</sup><https://www.justice.gov/crt/fair-housing-act-1>

<sup>3</sup><https://www.congress.gov/bill/116th-congress/house-bill/2231/all-info>

<sup>4</sup><http://www.legalservicesindia.com/article/1145/Reservations-In-India.html>

<sup>5</sup><https://www.eeoc.gov/initiatives/e-race/significant-eeoc-racecolor-casescovering-private-and-federal-sectors#intersectional>

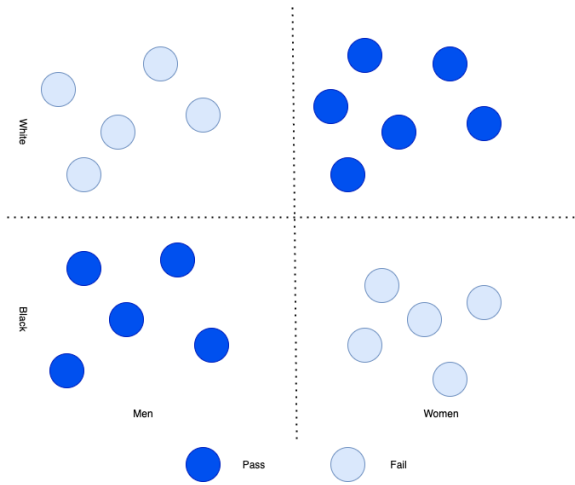


Figure 1: An example of "fairness gerrymandering"

darker skinned women disproportionately affected having a misclassification rate of  $\approx 30\%$ .

The example in Figure 1 describes the importance of intersectional fairness. In the figure, we observe equal numbers of black and white people pass. Similarly, there is an equal number of men and women passing. However, this classification is unfair because we don't have any black women and white men that passed, and all black men and white women passed. We observe the bias only while looking at the subgroups when we take race and gender as protected attributes. This phenomenon was called "*Fairness Gerrymandering*" by Kearns et al. (2018).

Additionally, there are minorities that have historically faced discrimination around the world, but due to their sparse population, empirical evidence of discrimination against them is difficult to trace, for example, the indigenous population (Paradies 2006; King, Smith, and Gracey 2009), or trans people (Feldman et al. 2016; Reisner et al. 2016; Bockting et al. 2016). This causes machine learning practitioners to either disinclude these groups from their training datasets due to statistical insignificance, or worse, conflate them with other minorities to create a general "protected" category, which leads to the same sort of neglected bias as shown in figure 1.

In this paper, we discuss the notion of intersectional group fairness. After introducing existing related work, we define a combinatorial approach giving subsets of the population. With this definition of subgroups, we introduce a measure of the worst case disparity using existing fairness metrics, to discover biases against underserved subgroups. We then show how this method can be applied to classification models, ranking models, and models with continuous output. We end the paper with a discussion about the limitations of our approach and future work.

## 2 Related Work

**Individual and Group Fairness** Fair machine learning differentiates *group* and *individual* fairness measures. While

group fairness metrics focus on treating two different groups equally, individual fairness metrics focus on treating similar individuals similarly. Binns (2019) introduces those two notions and discusses the motivations behind individual and group fairness. In this paper, we focus on group fairness metrics.

**Binary fairness metrics** A large majority of research in algorithmic fairness has covered fairness metrics for a single protected attribute (Corbett-Davies and Goel 2018). Hardt, Price, and Srebro (2016) introduces the definitions of Equalized odds and equal opportunity, two measures for discrimination against a binary sensitive attribute. Verma and Rubin (2018) collected some known binary fairness metrics for classification models and demonstrated each metric with a unique example on the German credit dataset. In their example, the protected class is *Gender*, which has two values *female* and *male*.

**Intersectional Fairness** More recently, however, some work has begun to address the issue of intersectionality in AI by providing statistical frameworks that control for bias within multiple subgroups. Hébert-Johnson et al. (2018) introduces the idea of *multi-calibration* which gives meaningful predictions for overlapping subgroups in a larger protected group. Kearns et al. (2018) developed an analogous method named *rich subgroup fairness* for false positive and negative constraints that hold over an infinitely large collection of subgroups. Kim, Ghorbani, and Zou (2019) extend these methods for classifiers to be equally accurate on a combinatorially large collection of all subgroups. Mary, Calauzènes, and El Karoui (2019) present the Renyi correlation coefficient as a fairness metric for datasets with continuous protected attributes. Finally, Foulds et al. (2020) introduce *differential fairness* (DF), as an intersectional fairness metric.

## 3 Intersectional group fairness metrics

In this section, we discuss our intersectional fairness metrics framework. We outline our definition of a subgroup of the population, define a *worst case* disparity metric that we call the *min-max ratio* and describe how we can operationalize the notion of *min-max ratio* to encompass intersectionality in existing metrics of fairness.

### Subgroup definition

For the purposes of this paper, similar to Kearns et al. (2018), we define a subgroup  $sg_{a_1 \dots a_n}$  as a set containing the intersection of all members who belong to groups  $g_{a_1}$  through  $g_{a_n}$ , where  $a_1, a_2 \dots a_n$  are marginal protected attributes, like race, gender, etc. Formally,

$$sg_{a_1 \times a_2 \times \dots \times a_n} = g_{a_1} \cap g_{a_2} \dots \cap g_{a_n} \quad (1)$$

Hence, for example, if  $g_1(\text{race}) \in \{\text{black}, \text{white}\}$  and  $g_2(\text{gender}) \in \{\text{man}, \text{woman}\}$ , then  $sg \in \{\text{black women}, \text{black men}, \text{white women}, \text{white men}\}$  and  $N = |sg| = 2 \times 2 = 4$ .

This combinatorial, or cartesian product of attributes approach gives us subsets of the original dataset, where in

each subgroup, the members have all the protected attributes of the groups they were composed of.

### Worst Case Disparity

We introduce a simple concept to measure the worst case disparity using existing fairness metrics to incorporate intersectionality - the *min-max ratio*. Essentially, the idea is to measure the value of the given fairness metric for every subgroup  $sg_i$  then take the ratio of the *minimum* and *maximum* values from this given list. The further this ratio is from 1, the greater the disparity is between subgroups. **This *min-max ratio* technique allows us to encompass the entire breadth of possible subgroups in a dataset, by considering the worst case scenario in terms of adverse impact.** For fairness metrics that are already comparative ratios of two groups, we redefine it by calculating the said ratio for all possible permutations of two subgroups and then simply take the minimum, also the worst possible case.

We discuss some of the most commonly used metrics in the literature below and show how we use the *worst possible case* framing to incorporate intersectionality.

### Fair Classification metrics

Several fair classification metrics exist in literature. Mehrabi et al. (2019) discuss four group fairness metrics that we introduce below.

**Demographic parity** According to demographic parity, the proportion of each segment of a protected class should receive positive outcomes at equal rates. Mathematically, demographic parity compares the pass rate (rate of positive outcome) of two groups. Demographic parity is satisfied for a predictor  $\hat{Y}$  and for a member  $A$  if:

$$P(\hat{Y}|A \in sg_i) = P(\hat{Y}|A \in sg_j); \forall i, j \in N, i \neq j \quad (2)$$

where  $N$  is the total number of subgroups. Demographic parity is also known as statistical parity (Dwork et al. 2012; Kusner et al. 2017).

Using our *worst case, min-max ratio* definition, *Demographic parity ratio* (DPR) would be defined as:

$$DPR = \frac{\min\{P(\hat{Y}|A \in sg_i) \forall i \in N\}}{\max\{P(\hat{Y}|A \in sg_i) \forall i \in N\}} \quad (3)$$

Disparate impact, as defined under the guideline by the Equal Employment Opportunity Commission et al. (1979) is similar to the demographic parity metric. It is intended as a way to measure indirect and unintentional discrimination in which certain decisions disproportionately affect members of a protected group. Disparate impact compares the pass rate of one group versus another. The Four-Fifths rule states that the ratio of the pass rate of group 1 to the pass rate of group 2 has to be greater than 80% (groups 1 and 2 interchangeable). Using our worst case definition, intersectional disparate impact (DI) is defined as the minimum disparate impact between all possible pairs of subgroups  $sg$ .

$$DI = \min \left\{ \frac{P(\hat{Y}|A \in sg_i)}{P(\hat{Y}|A \in sg_j)}; \forall i, j \in N, i \neq j \right\} \quad (4)$$

**Conditional statistical parity** Conditional statistical parity extends demographic parity by permitting a set of legitimate attributes to affect the outcome (Corbett-Davies et al. 2017). Conditional statistical parity is satisfied for a predictor  $\hat{Y}$ , a member  $A$  with a set of legitimate attributes  $L$  if:

$$P(\hat{Y}|L = 1, A \in sg_i) = P(\hat{Y}|L = 1, A \in sg_j) \quad \forall i, j \in N, i \neq j \quad (5)$$

Using the *worst case, min-max ratio* definition, *Conditional statistical parity ratio* (CSPR) would be defined just like equation 3:

$$CSPR = \frac{\min\{P(\hat{Y}|L = 1, A \in sg_i) \forall i \in N\}}{\max\{P(\hat{Y}|L = 1, A \in sg_i) \forall i \in N\}} \quad (6)$$

**Equal opportunity** Equal opportunity states that all members should be treated equally or similarly and not disadvantaged by prejudice or bias. Mathematically, it compares True Positive Rate (TPR) of the classifier between the protected group and the unprotected group<sup>6</sup>(Hardt, Price, and Srebro 2016). Equal opportunity for a binary predictor  $\hat{Y}$  and a member  $A$ , is satisfied if:

$$P(\hat{Y} = 1|A \in sg_i, Y = 1) = P(\hat{Y} = 1|A \in sg_j, Y = 1) \quad \forall i, j \in N, i \neq j \quad (7)$$

Using the *worst case, min-max ratio* definition, *Equal opportunity ratio* (EOppR) would be defined as:

$$EOppR = \frac{\min\{P(\hat{Y} = 1|A \in sg_i, Y = 1) \forall i \in N\}}{\max\{P(\hat{Y} = 1|A \in sg_i, Y = 1) \forall i \in N\}} \quad (8)$$

**Equalized odds** Equalized odds is a more complete version of equal opportunity as it compares both the True Positive Rate (TPR) and the False Positive Rate (FPR) between the protected group and the unprotected group<sup>7</sup>(Hardt, Price, and Srebro 2016). It essentially therefore measures the outcome of the classification irrespective of the ground truth. Equalized odds for a binary predictor  $\hat{Y}$  and a member  $A$  is satisfied if:

$$P(\hat{Y} = 1|A \in sg_i, Y = y) = P(\hat{Y} = 1|A \in sg_j, Y = y) \quad y \in \{0, 1\}, \forall i, j \in N, i \neq j \quad (9)$$

<sup>6</sup>TPR is the probability that a ground truth positive observation is correctly classified as positive.

<sup>7</sup>FPR is the probability that a ground truth negative observation is incorrectly classified as positive.

## Equal Opportunity

MIN-MAX ratio 0.53

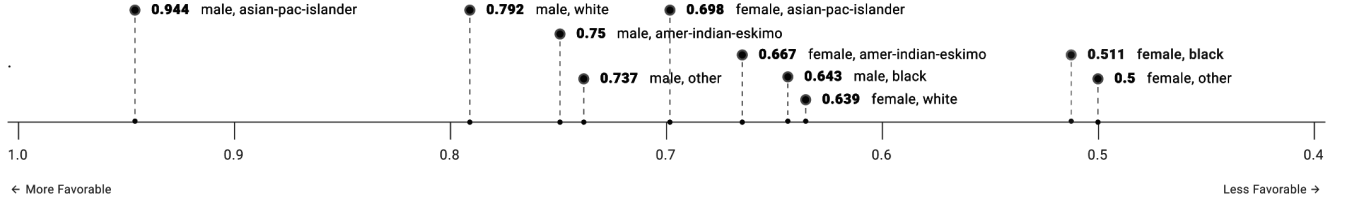


Figure 2: An example calculation of Equal Opportunity Ratio. The number line shows the True Positive Rates of the different subgroups in a sample classification output. The ratio of the minimum value to the maximum value in this range ( $0.5/0.944 = 0.53$ ) is the EOppR.

Using the *worst case, min-max ratio* definition, *Equal odds ratio* (EOddR) would be defined as:

$$\text{EOddR} = \frac{\min\{P(\hat{Y} = 1|A \in sg_i, Y = y)\}}{\max\{P(\hat{Y} = 1|A \in sg_i, Y = y)\}}; \quad (10)$$

$$y \in \{0, 1\}, \forall i \in N$$

## Multi-class classification models

For multiclass classification models, we present a modified version of the Equalized Odds metric, except instead of a binary *positive* or *negative* label, we measure the odds ratio for each possible discrete output, and then take the worst odds ratio among all outputs.

For instance, if a multiclass classifier has five possible output classes, we calculate the min-max ratio for each output class  $y$ , and then take the minimum of those five values as our final metric, since it is the *worst possible scenario*. Formally, Multiclass Equalized Odds Ratio (M-EOddR) is defined as:

$$\text{M-EOddR} = \min \left\{ \frac{\min\{P(\hat{Y} = y_k|A \in sg_i), \forall i \in N\}}{\max\{P(\hat{Y} = y_k|A \in sg_i), \forall i \in N\}} \right\} \quad \forall k \in K \quad (11)$$

where  $K$  is the set of all possible output classes. The closer the value of M-EOddR is to 1, the lower the disparity is of the classifier's performance among the various subgroups for all possible output classes.

## Models with continuous output

We can extend the worst possible case framing for models which produce a continuous output, like regression models, or recommendation models that provide relevance scores. The Kullback-Leibler (KL) divergence<sup>8</sup> between two

<sup>8</sup>Here we use the KL Divergence as our base metric, although this method would work for any distribution comparison metric.

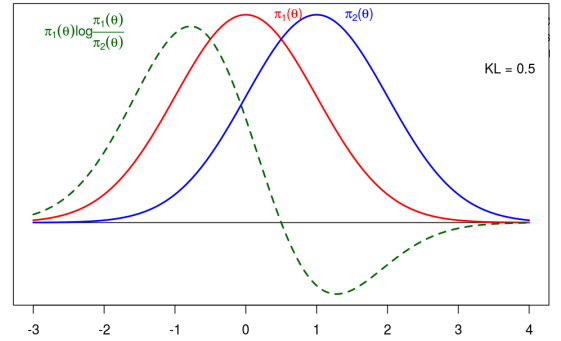


Figure 3: KL divergence example between two distributions adapted from Veen et al. (2018). In this example  $\pi_1$  is a standard normal distribution and  $\pi_2$  is a normal distribution with a mean of 1 and a variance of 1. The value of the KL divergence is equal to the area under the curve of the function (green line). The area under the green line above the x-axis adds to the divergence, while the area under the x-axis subtracts from the divergence.

distributions  $q$  and  $p$  is defined as the following:

$$D_{KL}(\pi_1||\pi_2) = \int_{-\infty}^{\infty} \pi_1(x) \log\left(\frac{\pi_1(x)}{\pi_2(x)}\right) dx \quad (12)$$

In the context of intersectional fairness, we compute the KL divergence between the model output distributions of all possible pairs of two subgroups, and we display the maximum KL divergence value obtained, since it is the *worst case scenario*. If this value is close to 0, the two subgroups have similar distributions, as well as the other subgroups.

Thus, Worst case KL Divergence ( $W-D_{KL}$ ) is formally defined as:

$$W-D_{KL} = \max \left\{ \int_{-\infty}^{\infty} \pi_{sg_i}(x) \log \left( \frac{\pi_{sg_i}(x)}{\pi_{sg_j}(x)} \right) dx \right\} \quad (13)$$

$\forall i, j \in N, i \neq j$

## Ranking metrics

Existing fair ranking metrics in the literature can be divided broadly into two classes - representation based (Yang and Stoyanovich 2017) and exposure based (Singh and Joachims 2018; Sapiezynski et al. 2019). We pick one of each kind and redefine them under the light of intersectionality.

**Skew** The *representation-based* metric we discuss is skew@k (Geyik, Ambler, and Kenthapadi 2019). For a ranked list  $\tau$ , the Skew for subgroup  $sg_i$  at the top  $k$  is defined as

$$\text{Skew}_{sg_i} @k(\tau) = \frac{p_{\tau^k, sg_i}}{p_{q, sg_i}} \quad (14)$$

where  $p_{\tau^k, sg_i}$  represents the fraction of members from subgroup  $sg_i$  among the top  $k$  items in  $\tau$ , and  $p_{q, sg_i}$  represents the fraction of members from subgroup  $sg_i$  in the overall population  $q$ . Ideally,  $\text{Skew}_{sg_i} @k$  should be close to one for each  $sg_i$  and  $k$ , to show that people from  $sg_i$  are represented in  $\tau$  proportionally relative to the overall population.

Using our *worst case* method, the skew ratio at K (SR@K) is defined as:

$$\text{SR@K} = \frac{\min\{\text{Skew}_{sg_i} @k(\tau), \forall i \in N\}}{\max\{\text{Skew}_{sg_i} @k(\tau), \forall i \in N\}} \quad (15)$$

**Attention** Attention is the *exposure-based* metric we discuss here. Ranking problems are unique from classification problems in the sense that the position of a ranked item, even within the top K results, can draw significantly different levels of visual attention. Previous research shows that people’s attention sharply drops off after the first few items in a ranked list (Mullick et al. 2019). Different papers have modeled visual attention as a function of the position K as a logarithmic distribution (Singh and Joachims 2018), a geometric distribution, or other sharply falling distributions with increasing rank (Sapiezynski et al. 2019). Assuming the attention distribution function of an item to be  $\text{Att}(k)$ , the mean attention per subgroup is defined as:

$$\text{MA}_{sg_i} = \frac{1}{|sg_i|} \sum_{k=1}^{|\tau|} \text{Att}(k) \text{ where } sg_k^\tau = sg_i \quad (16)$$

And, using our *worst case* method, the attention ratio (AR) is defined as:

$$\text{AR} = \frac{\min\{\text{MA}_{sg_i}, \forall i \in N\}}{\max\{\text{MA}_{sg_i}, \forall i \in N\}} \quad (17)$$

## 4 Discussion

**Conclusion** In this paper, we introduce the *worst-case* comparison as a simple, easily comprehensible method to surface hidden biases that commonly used fairness metrics may not be able to show. We establish the importance of introducing such modifications to better serve minorities with sparse populations and show how the method can be applied to a diverse range of model metrics, thereby being easy for practitioners and researchers to adapt without significantly changing their existing fairness monitoring systems.

### Limitations and Future Work

The idea of creating combinatorial subgroups has a couple of caveats: It does not take into account partial group membership (eg., an individual who is half white and half black, or in terms of gender grouping, a non-binary person), or continuous variables (for example, instead of treating age as an integer, we would convert the age attribute as discrete buckets). We encourage researchers to expand our method to include partial group membership and continuous attributes.

Secondly, creating a combinatorially large number of subgroups inevitably leads to subgroups which have a very small number of members, thereby demonstrating the effects of Simpson’s Paradox (Blyth 1972). A possible direction of research could be to introduce statistical significance measures for such small subgroups, and suggest thumb rules for subgroup creation via empirical measurements.

### Acknowledgments

The authors would like to thank Joshua Rubin, Amit Paka, Jack Reidy, Aalok Shanbhag and Christo Wilson for their helpful discussions and suggestions.

### References

- Berg, R. K. 1964. Equal employment opportunity under the Civil Rights Act of 1964. *Brook. L. Rev.* 31: 62.
- Binns, R. 2019. On the Apparent Conflict Between Individual and Group Fairness. *CoRR* abs/1912.06883. URL <http://arxiv.org/abs/1912.06883>.
- Blyth, C. R. 1972. On Simpson’s paradox and the sure-thing principle. *Journal of the American Statistical Association* 67(338): 364–366.
- Bockting, W.; Coleman, E.; Deutsch, M. B.; Guillamon, A.; Meyer, I.; Meyer III, W.; Reisner, S.; Sevelius, J.; and Ettner, R. 2016. Adult development and quality of life of transgender and gender nonconforming people. *Current opinion in endocrinology, diabetes, and obesity* 23(2): 188.
- Bogen, M.; and Rieke, A. 2018. Help wanted: An examination of hiring algorithms, equity, and bias .
- Buolamwini, J.; and Gebru, T. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, 77–91.
- Commission, E. E. O.; et al. 1979. Adoption of questions and answers to clarify and provide a common interpretation

- of the Uniform Guidelines on Employee Selection Procedures. *Federal Register* 44(43): 11996–12009.
- Corbett-Davies, S.; and Goel, S. 2018. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*.
- Corbett-Davies, S.; Pierson, E.; Feller, A.; Goel, S.; and Huq, A. 2017. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*, 797–806.
- Crenshaw, K. 1989. Demarginalizing the intersection of race and sex: A black feminist critique of antidiscrimination doctrine, feminist theory and antiracist politics. *u. Chi. Legal f.* 139.
- Dwork, C.; Hardt, M.; Pitassi, T.; Reingold, O.; and Zemel, R. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, 214–226.
- Feldman, J.; Brown, G. R.; Deutsch, M. B.; Hembree, W.; Meyer, W.; Meyer-Bahlburg, H. F.; Tangpricha, V.; T’Sjoen, G.; and Safer, J. D. 2016. Priorities for transgender medical and health care research. *Current opinion in endocrinology, diabetes, and obesity* 23(2): 180.
- Foulds, J. R.; Islam, R.; Keya, K. N.; and Pan, S. 2020. An intersectional definition of fairness. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, 1918–1921. IEEE.
- Geyik, S. C.; Ambler, S.; and Kenthapadi, K. 2019. Fairness-aware ranking in search & recommendation systems with application to linkedin talent search. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2221–2231.
- Goodman, B.; and Flaxman, S. 2017. European Union regulations on algorithmic decision-making and a “right to explanation”. *AI magazine* 38(3): 50–57.
- Hardt, M.; Price, E.; and Srebro, N. 2016. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, 3315–3323.
- Hébert-Johnson, Ú.; Kim, M.; Reingold, O.; and Rothblum, G. 2018. Multicalibration: Calibration for the (computationally-identifiable) masses. In *International Conference on Machine Learning*, 1939–1948.
- Kearns, M.; Neel, S.; Roth, A.; and Wu, Z. S. 2018. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In *International Conference on Machine Learning*, 2564–2572. PMLR.
- Kim, M. P.; Ghorbani, A.; and Zou, J. 2019. Multiaccuracy: Black-box post-processing for fairness in classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 247–254.
- King, M.; Smith, A.; and Gracey, M. 2009. Indigenous health part 2: the underlying causes of the health gap. *The lancet* 374(9683): 76–85.
- Kusner, M. J.; Loftus, J.; Russell, C.; and Silva, R. 2017. Counterfactual fairness. In *Advances in neural information processing systems*, 4066–4076.
- Mary, J.; Calauzènes, C.; and El Karoui, N. 2019. Fairness-aware learning for continuous attributes and treatments. In *International Conference on Machine Learning*, 4382–4391.
- Mehrabi, N.; Morstatter, F.; Saxena, N.; Lerman, K.; and Galstyan, A. 2019. A Survey on Bias and Fairness in Machine Learning. *CoRR* abs/1908.09635. URL <http://arxiv.org/abs/1908.09635>.
- Mukerjee, A.; Biswas, R.; Deb, K.; and Mathur, A. P. 2002. Multi-objective evolutionary algorithms for the risk–return trade-off in bank loan management. *International Transactions in operational research* 9(5): 583–597.
- Mullick, A.; Ghosh, S.; Dutt, R.; Ghosh, A.; and Chakraborty, A. 2019. Public Sphere 2.0: Targeted Commenting in Online News Media. In *European Conference on Information Retrieval*, 180–187. Springer.
- Paradies, Y. 2006. A systematic review of empirical research on self-reported racism and health. *International journal of epidemiology* 35(4): 888–901.
- Reisner, S. L.; Deutsch, M. B.; Bhasin, S.; Bockting, W.; Brown, G. R.; Feldman, J.; Garofalo, R.; Kreukels, B.; Radix, A.; Safer, J. D.; et al. 2016. Advancing methods for US transgender health research. *Current opinion in endocrinology, diabetes, and obesity* 23(2): 198.
- Sapiezynski, P.; Zeng, W.; E Robertson, R.; Mislove, A.; and Wilson, C. 2019. Quantifying the Impact of User Attention Fair Group Representation in Ranked Lists. In *Companion Proceedings of The 2019 World Wide Web Conference*, 553–562.
- Singh, A.; and Joachims, T. 2018. Fairness of exposure in rankings. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2219–2228.
- Vasilescu, M. A. O.; and Terzopoulos, D. 2002. Multilinear image analysis for facial recognition. In *Object recognition supported by user interaction for service robots*, volume 2, 511–514. IEEE.
- Veen, D.; Stoel, D.; Schalken, N.; Mulder, K.; and Van de Schoot, R. 2018. Using the data agreement criterion to rank experts’ beliefs. *Entropy* 20(8): 592.
- Verma, S.; and Rubin, J. 2018. Fairness Definitions Explained. In *Proceedings of the International Workshop on Software Fairness, FairWare ’18*, 1–7. New York, NY, USA: Association for Computing Machinery. ISBN 9781450357463. doi:10.1145/3194770.3194776. URL <https://doi.org/10.1145/3194770.3194776>.
- Yang, K.; and Stoyanovich, J. 2017. Measuring fairness in ranked outputs. In *Proceedings of the 29th International Conference on Scientific and Statistical Database Management*, 1–6.