

---

# Intersectional Fairness Score : the overlooked but far-reaching choice of aggregation design

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Fairness assessment in AI is essential for building responsible models. Tradition-  
2 ally, it focuses on two demographic groups situations, but real-world complexity  
3 requires considering intersectionality [2]. This work explores how to aggregate  
4 bias measurements across multiple subgroups into one single score—a critical  
5 step often overlooked. We first show that the choices made in aggregation design  
6 (norm, maximum, probabilistic approaches, etc.) can significantly influence results,  
7 leading to divergent or even conflicting conclusions. We identify and analyze the  
8 various possible methods, highlighting their ethical implications and providing a  
9 first framework of criteria to guide their selection based on context. Our goal is to  
10 foster interdisciplinary discussion on this often-neglected step, aiming for a fairer,  
11 more informed and transparent evaluation of intersectional biases in AI.

## 12 1 Introduction

13 Fairness assessment in artificial intelligence (AI) models is a critical step toward developing re-  
14 sponsible, trustworthy, and ethically sound systems. As AI increasingly influences high-stakes  
15 decisions across diverse domain, the need for rigorous evaluation of biases and fairness becomes  
16 paramount. Existing literature has established a broad framework for fairness metrics [7], highlight-  
17 ing the complexity of defining and operationalizing fairness. These efforts often focus on binary  
18 scenarios, where fairness notions are straightforwardly applied to two groups distinguished by a  
19 unique attribute. However, real-world applications frequently involve multiple sensitive attributes  
20 requiring an intersectional approach of bias to consider the combinatorial way discrimination occurs.

21 Intersectionality, originally conceptualized within social sciences, emphasizes the interconnected  
22 nature of social categorizations such as race, gender, and class, which produce overlapping systems  
23 of discrimination or privilege. In the context of AI fairness, addressing intersectionality involves  
24 moving from a binary two-dimensional problem to a multi-dimensional one. This transition comes  
25 with new particularities and additional steps in the process of bias assessment. We claim that these  
26 steps have been overlooked due to insufficiently meticulous reuse of existing tools developed for the  
27 binary case.

28 This work aims to shed light on the specific step of subgroups measure aggregation strategy to  
29 produce a comprehensive fairness measure. We argue that this step implies a series of technical  
30 choices that are embedded with ethical considerations. After an identification and an analysis of  
31 various fairness score design possibilities (Section 3), we demonstrate how different choices can  
32 lead to markedly divergent fairness assessments (Section 4). Practical and intuitively understandable  
33 examples illustrate the profound impact of these decisions, emphasizing the importance of informed,  
34 context-sensitive design rather than default or arbitrary selections. Section 5 opens a discussions on  
35 some guidelines that could help choose the appropriate technical tool depending on the results of the  
36 ethical analysis of the model building’s context.

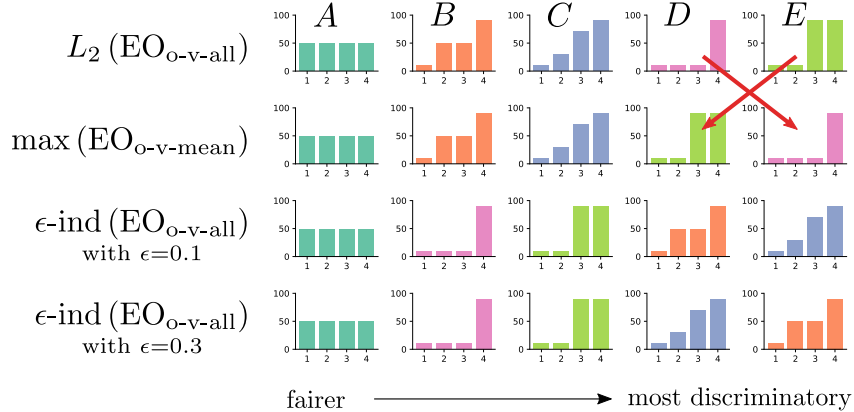


Figure 1: Rankings obtained by different aggregation methods on the same set of values for subgroups. By choosing a different aggregation method, we can reverse what is considered fair or discriminatory.

37 Our contribution underscores the necessity of making explicit the design choices involved in inter-  
 38 sectional fairness assessment. By raising awareness of this critical stage, we aim to foster more  
 39 transparent, equitable, and contextually appropriate AI systems, ultimately advancing the field toward  
 40 more responsible AI development.

## 41 2 Related Work

42 **Fairness as part of AI ethics** Fairness in AI require choices in defining the assessment framework  
 43 [9]. These are not only technical, but ethical in essence. As ethical choices, they address questions  
 44 that do not have a universal correct and obvious answer, but rather require careful consideration  
 45 in order to ultimately determine the best possible solution. These elements encompass the context  
 46 in which the AI model will be used, among which cultural habits, values, political ideologies and  
 47 societal choices. These should be collaboratively discussed, in working groups comprising individuals  
 48 from a range of social backgrounds and different academic disciplines (law, sociology, philosophy,  
 49 formal sciences, etc.), informed by comprehensive technical assessments of potential fairness issues  
 50 [10].

51 **The various levels of technical choices in bias assessment that involve ethical considerations**  
 52 Aggregation design is not the first technical and ethical choice. First, human characteristics on which  
 53 discrimination must not be based have to be identified [4]. In AI fairness, the term *sensitive attribute*  
 54 is used to refer to these characteristics [8]. Whilst this step may appear to be an immediate result, it is  
 55 rare for a justification to be provided for the reasons behind their selection.

56 Secondly, If we agree that discrimination is prejudicial differential treatment, how does this translate  
 57 in practical terms in the case of an AI model prediction? [1] If we focus on the case of binary  
 58 classification models, is it to obtain the disadvantageous prediction? is it to obtain this prediction  
 59 knowing that we should have obtained the opposite one? The extensive range of possibilities has  
 60 resulted in the generation of numerous *notions of fairness* in the state of the art [7]. Among the  
 61 best known and most frequently used : Statistical Parity (SP), Equalized Odds (EOd) and Equal  
 62 Opportunity (EO).

## 63 3 Problem Setting

64 Let  $D$  be a dataset of instances  $(\mathbf{X}, \mathbf{A}, Y) = (x_1, \dots, x_m, a_1, \dots, a_{n-m}, y)$ , each corresponding  
 65 to an individual. Among the features,  $\mathcal{A} = \prod_{i=m+1}^n \mathcal{A}_i$  denotes the *sensitive attributes*, i.e the  
 66 features for which discrimination should not occur. We denote the number of sensitive attributes  
 67 considered as  $|\mathcal{A}|$ . When  $|\mathcal{A}| > 2$ , intersectionality occurs. Based on the sensitive attributes, we  
 68 define *demographic groups* or *subgroups* in intersectional setups, denoted as  $g \in G$ , where each  
 69 group is characterized by individuals sharing the same values for the sensitive attributes. The total

70 number of groups considered is denoted by  $|G|$ . We distinguish two main cases:  $|G| = 2$ , the binary  
 71 scenario on which existing literature mostly focuses, and  $|G| > 2$ , the non-binary scenario.

72 For a classifier  $M$  learned using a machine learning algorithm, group fairness is usually measured by  
 73 choosing a metric (e.g. the probability of being correctly classified, or the probability of being assigned  
 74 the advantageous outcome), corresponding to the previously chosen *fairness notion*, that quantifies  
 75 how much its prediction differs depending on whether an instance belongs to a protected group or not.  
 76 Even though our work is generic and can be extended to different choices of fairness notion/metric, for  
 77 ease of exposure, we focus on the Equalized Odds metric  $EO(M) = P(M(x) = 1|Y = 1, A = 1) -$   
 78  $P(M(x) = 1|Y = 1, A = 0)$  that quantifies the difference in the odds of obtaining the positive  
 79 outcome  $M(x) = 1$  provided that the true class  $Y$  is 1, depending on the protected attribute  $A$ .

80 **Binary case** Provided that the outcome  $M(x) = 1$  is judged favorably, the metric directly compares  
 81 the (potential) benefit of belonging to subgroup  $A = 1$ :  $P(M(x) = 1|Y = 1, A = 1)$  over subgroup  
 82  $A = 0$ :  $P(M(x) = 1|Y = 1, A = 0)$  by a simple subtraction since there are only 2 quantities to  
 83 compare, which immediately results in a single fairness score. When comparing any 2 models  $M_1$   
 84 and  $M_2$ , we can tell which model is fairer than the other by comparing the score (a scalar).

85 **Intersectionality: the need for summarizing many measurements.** Intersectionality implies  
 86 a setup where the number of subgroups  $|G| > 2$ . We obtain a benefit score for each subgroup.  
 87 Comparison of benefit for every pair of subgroups results in  $\frac{|G|(|G|-1)}{2}$  measurements, that we would  
 88 like to *summarize in a global fairness score* by aggregating measurements into a single scalar. The  
 89 choice of aggregation method is essential, as the subtle differences between aggregation methods  
 90 can imply different conclusions. Model  $M_1$  can be deemed more fair than model  $M_2$  by a choice of  
 91 aggregation method, or less fair by a different choice.

92 While there are certain applications (e.g., healthcare) where improving the utility of the worst-case  
 93 groups is an important goal, other applications can be required by law to ensure more parity for all  
 94 subgroups. In the following, we explicitly list the possible design choices when applying fairness  
 95 assessment to an intersectional setup, and describe different mathematical tools to summarize  $|G|$   
 96 measurements into a single scalar value.

### 97 3.1 First design choice: comparison type

In the **one-vs-all** approach, we form the vector of differences between all  $\frac{|G|(|G|-1)}{2}$  possible pairs  
 ( $i, j$ ) of subgroups.

$$EO_{\text{one-vs-all}}(M, i, j) := P(M(x) = 1|Y = 1, A = A_i) - P(M(x) = 1|Y = 1, A = A_j)$$

In the **one-vs-mean** approach (e.g. in [3]), the value measured on one subgroup is instead compared  
 to the others through their mean value and we use the  $|G|$  vector of differences to the averaged benefit.

$$EO_{\text{one-vs-mean}}(M, i) := P(M(x) = 1|Y = 1, A = A_i) - P(M(x) = 1|Y = 1)$$

98 In these first two approaches, subtraction is used but we could instead use ratios.

The **all-in-one** approach is slightly different because it is not founded on an individual analysis but  
 directly encompasses all measures immediately leading to the final result of a single score. It studies  
 the dependence of the variable of interest, i.e. the variable considered by the chosen fairness metric,  
 on the random vector of sensitive attributes, as measured by an independence criterion (IC) such as  
 the Mutual Information [5, 6] or the Hilbert Schmidt Independence Criterion.

$$EO_{\text{all-in-one}}(M) := IC(M(x) = 1, A | Y = 1)$$

### 99 3.2 Second design choice: aggregation method

100 In the one-vs-all and one-vs-mean approaches, it is then necessary to aggregate the multiple values  
 101 obtained, which can be achieved using either:

102

$L_q$  **norms**: Let  $q \in \mathbb{R}^+$ , a  $q$ -**norm** is defined for any vector  $(x_1, x_2, \dots, x_n)$  by

$$L_q = \left( \sum_{i=1}^n |x_i|^q \right)^{1/q}$$

103  $L_q$ -norms provide a flexible measure of the magnitude of a vector, where larger values of  $q$  penalize  
 104 the largest components, until the limit  $q \rightarrow \infty$  where  $L_\infty := \max_i(x_i)$ :

- 105 • **For  $q = 1$** : a uniform weight is given to all subgroups [3];
- 106 • **For  $q = 2$** : the norm highlights significantly discriminated subgroups more strongly, as it  
 107 emphasizes higher values;
- 108 • **For  $q \rightarrow \infty$** : the norm is dominated by the most discriminated subgroup, indicating the  
 109 **worst case** scenario.

**Ordered Weighted Averaging**: OWA consists of calculating the weighted average of a vector by weighting its elements according to their rank. Let  $(x'_1, x'_2, \dots, x'_n)$  be values of  $(x_1, x_2, \dots, x_n)$  ranked by increasing order, and  $(w_1, w_2, \dots, w_n)$  be a set of weights.

$$OWA := \sum_{i=1}^n w_i x'_i$$

110 This method also allows to adjust the sensitivity of the aggregate measurement to extreme values or  
 111 overall distribution, with even greater flexibility than that allowed by  $L_q$  norms.  $(w_1, w_2, \dots, w_n)$  are  
 112 hyperparameters that broaden the possibilities and make the method fully customizable. It notably  
 113 encompasses measure of **minimum**, **maximum** and **mean**.

**Threshold indicators**: Let  $\varepsilon \in [0, 1]$ . For any vector  $(x_1, x_2, \dots, x_n)$ , the  $\varepsilon$ -*threshold indicator* counts the number of values that are greater than  $\varepsilon$ :

$$\varepsilon - ind := \frac{1}{n} \sum_{i=0}^n \mathbb{1}(x_i > \varepsilon)$$

114 The hyperparameter  $\varepsilon$  sets the threshold height. For small  $\varepsilon$  values, the priority is to avoid gaps, no  
 115 matter how large or small. For large  $\varepsilon$  values, the focus is more on avoiding very large gaps. This  
 116 focuses attention on whether a certain threshold has been exceeded in the measured gaps. Unlike  
 117 previous methods, it does not consider the magnitude of the gaps. The salient consideration is the  
 118 frequency with which biases exceed a substantial magnitude to be considered. This places greater  
 119 emphasis on minimizing the occurrence of bias rather than reducing it to the lowest possible level.

### 120 3.3 Variations

121 The possibilities listed in the previous sections lead to methods which can be adapted to specific  
 122 needs by slightly modifying them in different ways. One option is to weight measurements made by  
 123 subgroups according to the size of their sample. This is useful when samples are too small and may  
 124 not be representative. Weighting will enable the exclusion of outliers that could have a significant  
 125 impact even though they are in fact incorrectly estimated.

126 On the contrary, there is in fairness a propensity to focus on small samples, which often correspond  
 127 to subgroups less visible and therefore discriminated against. It may be desirable to give more weight  
 128 to these groups. In this case, weighting by the inverse of the sample size might be considered.

## 129 4 Illustrative cases

130 **Different aggregation methods implicitly imply different rankings** As a first illustrative case,  
 131 for 4 choices of aggregation methods, we rank the same 5 distributions of benefits over subgroup as  
 132 measured by their aggregated score (Figure 1). All 4 methods deem distribution  $A$  fairer among all  
 133 since all subgroups receive the exact same benefit, but even though these 5 distributions all span the  
 134 same range of benefit, and have the same mean benefit = 50 (except for distribution  $D$ ), changing  
 135 aggregation alters the overall ranking. For instance (red arrows), distribution  $D$  and  $E$  are exchanged  
 136 by switching from aggregation method  $L_2$  ( $EO_{0-v-all}$ ) to  $\max$  ( $EO_{0-v-mean}$ ).

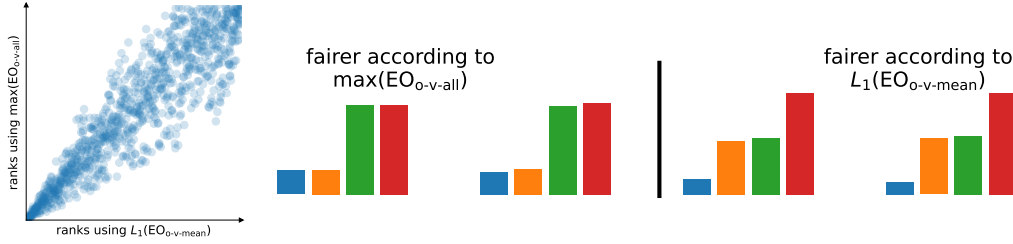


Figure 2: Pairwise comparison of aggregation methods  $L_1(\text{EO}_{o-v-\text{mean}})$  and  $\max(\text{EO}_{o-v-\text{all}})$ . **(left)** rank correlation and **(right)** distributions that produce most different rankings.

137 **Systematically quantifying pairwise difference between aggregation methods** We randomly  
 138 sample 1000 distributions of benefits for 4 subgroups with equal mean = 50. We measure their  
 139 fairness score using aggregation methods  $L_1(\text{EO}_{o-v-\text{mean}})$  and  $\max(\text{EO}_{o-v-\text{all}})$ , and compare the  
 140 rankings. In Figure 2 (left), we observe that overall, these produce quite correlated rankings (Kendall’s  
 141  $\tau$  rank correlation = 0.75). We then identify the distributions for which the rankings change the most  
 142 (Figure 2 right). Highlighting these discrepancies reveals the implications of the chosen aggregation  
 143 method. As we argue, this choice should be emphasized to encourage an informed discussion.

## 144 5 Discussions and Conclusions

145 In Sections 3 and 4, we highlighted the impact of the fairness score design on the valued patterns of  
 146 biases distributions. This intersectional fairness score can ultimately serve two purposes: first, as a  
 147 simple measure of bias in a dataset or model (fairness assessment); and second, as an indicator to  
 148 guide bias mitigation, for example by regularizing the learning objective function. In the second case,  
 149 this means that the bias mitigation protocol will cause the distribution of biases in the model to tend  
 150 towards those that obtain the best rankings by applying the chosen aggregation design. Consequently,  
 151 the score must be designed in line with how stakeholders define fairness.

### 152 5.1 Some guidelines

153 To facilitate the understanding and implementation of these fairness tools in AI models training, it  
 154 is essential to link analytical studies and observations with practical model user expectations. To  
 155 this end, we present preliminaries sets of guidelines that may assist in this decision. These initial  
 156 elements will require further study in future work.

157 **Trade-off between focusing on amplitude of gaps or frequencies of gaps.** In the first case,  
 158  $q$ -norms and OWA are recommended while  $\varepsilon$ -threshold are better suited for second option.

159 **Worst case or global amount of bias.** This criterion is achieved through correct adjustment of  
 160 hyperparameters. High  $q$ - and  $\varepsilon$ -values will both lead to focus on worst cases, and vice versa.

161 **Distribution of biases leading to isolated cases, or formation of blocks of groups treated in the  
 162 same way.** The choice of either *one-vs-mean* or *one-vs-all* approach enables the transition towards  
 163 the first situation or the second one, respectively.

### 164 5.2 Multidisciplinary discussion is key

165 Addressing the ethical challenges posed by advanced technological tools such as AI models requires  
 166 more than technical solutions alone. The role of the researcher in fairness in machine learning  
 167 cannot go beyond providing these technical tools for fairness assessment in intersectional setups, and  
 168 analyzing their implications to be discussed transparently. Subsequently, a multidisciplinary dialog is  
 169 required to bring together expertise from formal sciences and social sciences, as well as collective  
 170 societal discussion in the general public in cases that require political choices. This would lead to the  
 171 best informed option to answer the ethical goal of assessing and mitigating biases, and integrate the  
 172 corresponding tools in our mathematical models and AI systems.

## References

- 173
- 174 [1] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning: Limita-*  
175 *tions and opportunities*. MIT press, 2023.
- 176 [2] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in  
177 commercial gender classification. In *Conference on fairness, accountability and transparency*,  
178 pages 77–91. PMLR, 2018.
- 179 [3] Jaewoong Cho, Gyeongjo Hwang, and Changho Suh. A fair classifier using mutual information.  
180 In *2020 IEEE international symposium on information theory (ISIT)*, pages 2521–2526. IEEE,  
181 2020.
- 182 [4] Maryam Amir Haeri and Katharina Anna Zweig. The crucial role of sensitive attributes in fair  
183 classification. In *2020 IEEE symposium series on computational intelligence (SSCI)*, pages  
184 2993–3002. IEEE, 2020.
- 185 [5] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-aware classifier  
186 with prejudice remover regularizer. In *Joint European conference on machine learning and*  
187 *knowledge discovery in databases*, pages 35–50. Springer, 2012.
- 188 [6] Jian Kang, Tiankai Xie, Xintao Wu, Ross Maciejewski, and Hanghang Tong. Infofair:  
189 Information-theoretic intersectional fairness. In *2022 IEEE international conference on big*  
190 *data (big data)*, pages 1455–1464. IEEE, 2022.
- 191 [7] Karima Makhlof, Sami Zhioua, and Catuscia Palamidessi. Machine learning fairness no-  
192 tions: Bridging the gap with real-world applications. *Information Processing & Management*,  
193 58(5):102642, 2021.
- 194 [8] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A  
195 survey on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35,  
196 2021.
- 197 [9] Shira Mitchell, Eric Potash, Solon Barocas, Alexander D’Amour, and Kristian Lum. Algorithmic  
198 fairness: Choices, assumptions, and definitions. *Annual review of statistics and its application*,  
199 8(1):141–163, 2021.
- 200 [10] Deirdre K Mulligan, Joshua A Kroll, Nitin Kohli, and Richmond Y Wong. This thing called  
201 fairness: Disciplinary confusion realizing a value in technology. *Proceedings of the ACM on*  
202 *Human-Computer Interaction*, 3(CSCW):1–36, 2019.