




Metrics for Dataset Demographic Bias: A Case Study on Facial Expression Recognition

Iris Dominguez-Catena , *Student Member, IEEE*, Daniel Paternain , *Member, IEEE*,
and Mikel Galar , *Member, IEEE*

Abstract—Demographic biases in source datasets have been shown as one of the causes of unfairness and discrimination in the predictions of Machine Learning models. One of the most prominent types of demographic bias are statistical imbalances in the representation of demographic groups in the datasets. In this article, we study the measurement of these biases by reviewing the existing metrics, including those that can be borrowed from other disciplines. We develop a taxonomy for the classification of these metrics, providing a practical guide for the selection of appropriate metrics. To illustrate the utility of our framework, and to further understand the practical characteristics of the metrics, we conduct a case study of 20 datasets used in Facial Emotion Recognition (FER), analyzing the biases present in them. Our experimental results show that many metrics are redundant and that a reduced subset of metrics may be sufficient to measure the amount of demographic bias. The article provides valuable insights for researchers in AI and related fields to mitigate dataset bias and improve the fairness and accuracy of AI models.

Index Terms—AI fairness, artificial intelligence, deep learning, demographic bias, facial expression recognition.

I. INTRODUCTION

GENERAL advancements in technology, compounded with the widespread adoption of personal computers of all sorts, have led to an ever increasing exposure of society and non-expert users to autonomous systems. This interaction has also led to an accelerated deployment speed of state-of-the-art systems. Complex systems, such as general-purpose language and image models, conversational chatbots, or automatic face recognition systems, to name a few, are now deployed within months of their creation directly into the hands of nonexpert

Manuscript received 28 March 2023; revised 18 January 2024; accepted 31 January 2024. Date of publication 5 February 2024; date of current version 2 July 2024. This work was supported in part by a predoctoral fellowship from the Research Service of the Universidad Pública de Navarra through open access funding, in part by the Spanish MICIN under Grants PID2019-108392GB-I00, PID2020-118014RB-I00, and PID2022-136627NB-I00/AEI/10.13039/501100011033 FEDER, UE, and in part by the Government of Navarre under Grant 0011-1411-2020-000079 - Emotional Films. Recommended for acceptance by T. Hassner. (*Corresponding author: Iris Dominguez-Catena.*)

The authors are with the Department of Statistics, Computer Science and Mathematics, Public University of Navarre (UPNA), 31006 Pamplona, Spain, and also with the Institute of Smart Cities (ISC), Public University of Navarre (UPNA), 31006 Pamplona, Spain (e-mail: iris.dominguez@unavarra.es; daniel.paternain@unavarra.es; mikel.galar@unavarra.es).

The code is available at https://github.com/irisdominguez/dataset_bias_metrics.

This article has supplementary downloadable material available at <https://doi.org/10.1109/TPAMI.2024.3361979>, provided by the authors.

Digital Object Identifier 10.1109/TPAMI.2024.3361979

users. These systems, by their very nature, are difficult to evaluate and test, raising safety concerns. As systems interact with users in new and unpredictable ways, how can we ensure that no harm of any type is done to the user?

This general question is answered through the field of AI ethics [1]. This field, in turn, takes shape in several other aspects, focusing on issues such as the integration of robotics in society [2], issues of digital privacy [3], and many others. One particularly interesting concept is algorithmic fairness [4], which focuses on how systems can replicate human biases, discriminating people based on protected characteristics such as sex, gender, race, or age. Even if the concept of algorithmic fairness is broad and multifaceted, this notion of unwanted bias as the unwanted patterns learned by the machine makes them easier to characterize. In turn, the characterization and measurement of fairness favors the methodological mitigation of unfair behavior in trained models.

Although the development of bias is a complex phenomenon, deep learning techniques are especially susceptible to bias in datasets [5]. These techniques learn patterns autonomously and can often get confused between correlated patterns. When certain demographic characteristics are correlated with the target class of a problem, it is possible for the models to incorporate and amplify that correlation. This ends up resulting in a biased and differentiated prediction for certain individuals and demographic groups.

To recognize and solve these issues, it is crucial to measure bias, both in the final models and in the datasets. Although different metrics have been proposed [6], most of them focus only on the bias exhibited by the trained models. The measurement of bias in the source datasets has not received the same attention, although it enables the validation of new bias mitigation methods [7], the explanation of bias transference throughout the training process [8], and the demographic description of the application environment where a dataset or model can be safely used [9].

In this article, we explore metrics that can be used to measure demographic bias in datasets. Most previous works that have focused on this issue [10], [11] study biases only from an intuitive notion, without using bias metrics. A few works [7], [12] have employed metrics specific to bias, although only considering a single metric and without taking the different types of demographic bias into account. This work aims to serve as a unifying framework for the few metrics that have been already used for this purpose and those that can be adapted

from other fields with equivalent problems, such as population diversity metrics used in ecology. This wider variety of metrics allows us to cover specific types of biases that were previously unmeasured. Accordingly, we propose a taxonomy for dataset bias metrics based on the type of bias measured, facilitating the selection of appropriate metrics. Based on this taxonomy, we aim to find a concrete, expressive, and interpretable set of metrics to facilitate the work of analyzing the demographic and bias properties of ML datasets. To our knowledge, no such taxonomy or metric selection has been previously proposed.

An interesting case study to evaluate our taxonomy and set of metrics is Facial Expression Recognition (FER). FER is a problem in which photographs of people are classified according to the emotion they appear to express. The applications of this problem are varied, including healthcare [13] and assistive robotics [14] among others, and in most cases involve direct interaction with non-expert users. Unlike other problems in which bias is usually studied, mainly based on tabular data and with explicit demographic information [15], in FER the demographics of the person are rarely known. However, these demographic factors directly influence the person's appearance, generating demographic proxies that are completely embedded in the input images, making FER an interesting problem from the AI bias perspective. Even if this information is masked, it has been shown that both emotion expression and emotion identification are conditioned by a person's demographics [16], [17], so differential treatment based on demographics may be unavoidable.

More specifically, we gather twenty FER datasets according to a clear set of criteria, obtain a demographic profile of each of them, and apply the reviewed dataset bias metrics. We then employ these results to explore both the characteristics and limitations of each metric. We use this information to select a set of non-redundant interpretable metrics that can summarize the demographic biases in a dataset. Additionally, we employ the metrics to assess the types of biases found in FER datasets, where we observe differences between datasets created from different data sources. Identifying the different bias profiles of specific datasets can improve both the choice of training datasets and the choice of mitigation techniques.

Although this article focuses on dataset bias, in the Supplementary Material, we also provide the results of a series of experiments on the downstream propagation of such biases to the trained model, showing the importance of the appropriate characterization of the different types of dataset demographic bias.

The following sections are as follows. First, Section II recalls some related work in the field. Subsequently, in Section III we review and present a taxonomy of dataset demographic bias metrics. Section IV then presents the experimental framework for the FER case study, while Section V gathers the results found in these experiments. Finally, Section VI summarizes our findings and potential future work.

II. RELATED WORK

This section introduces some relevant background for our work. In Section II-A, we provide a brief overview of fairness,

its relationship to bias, and the methods used to measure it. Subsequently, in Section II-B, we explore the application of these concepts in the context of FER.

A. Fairness

The advances in ML and Artificial Intelligence (AI) in the last decades have led to an explosion of real-world applications involving intelligent agents and systems. This has inevitably led researchers to consider the social implications of these technologies and study what fairness means in this context [6], [18], [19], [20], [21]. The efforts to develop technical standards and best practices have also generated an increasing number of guidelines for ethical AI and algorithmic fairness [1].

Most definitions of fairness revolve around the concept of unwanted bias [20], also known as prejudice or favoritism, where particular individuals or groups defined by certain protected attributes, such as age, race, sex, and gender, receive an undesired differentiated treatment. In this sense, an algorithm or system is defined as fair if it is free of unwanted bias. It is important to note that although the concept of demographic bias is related to bias in machine learning and many results can be adapted to both [22], the particularities and potential harm resulting from demographic bias require an independent study. To this end, different metrics and mathematical definitions have been designed to characterize both fairness and demographic bias [6], [20], [23]. It is important to note that together with these metrics and definitions, criticism has also arisen [21], [24], since an excessive optimization of any given quantitative metric can lead to a loss of meaning, resulting in a false impression of fairness. As the fairness definitions and metrics proposed in the literature [20], [23] are mostly concerned with the social impact of the deployed systems, they deal only with the presence of bias in the final trained model, regardless of the source of that bias. These definitions, such as *equalized odds* [25], *equal opportunity* [25] or *demographic parity* [26], are usually designed to detect a disparate treatment between a single *privileged* demographic group and a single *protected* demographic group, in classification problems where one of the classes is considered preferable (usually the *positive* class). Despite this classical perspective, recent works [20], [21], [27] have focused on the multiple complementary sources that can lead to unwanted bias in the final model. These sources of bias originate in different phases of the AI pipeline [27], such as data collection, model training, and model evaluation. Regarding practical applications, these definitions and taxonomies of bias have been applied to multiple domains, where demographic biases have been found in facial [28], [29] and gender [11], [30] recognition, to name a few.

The bias detected in the source data has been a topic of particular interest over the years [5], [31]. Large public datasets have become the basic entry point for many AI projects, where developers often use them with limited knowledge of the origin of the data and its biases [31]. Some of the most popular ML datasets, such as ImageNet [32], [33] and COCO [34], have been revealed to include severe demographic biases [10] and even direct examples of racism, such as racial slurs [31]. To identify

such biases, some auxiliary datasets have been proposed, either by annotating previous datasets for apparent demographic characteristics [34], [35], or developing new demographically annotated datasets [36], [37], which can then be used to evaluate models and datasets. However, few works [7], [10], [11], [12], [38] have focused on the measurement and mathematical characterization of the bias present in the source data. Some works [10], [11] focus on non-systematic approaches, calculating the proportion of different demographic subgroups and manually looking for imbalances. Other works [7], [12] employ metrics based on information theory and statistics, such as the Mutual Information, to quantify bias in the source datasets, usually as part of a bias mitigation methodology, or focus on specific types of dataset, such as object detection [38], where specific metrics relating to the nature of the problem can be developed.

To the best of our knowledge, no previous work has unified these approaches, systematically comparing the properties of the different metrics. In this work, we explore and classify the full array of metrics already in use to measure dataset demographic bias and propose the application of existing metrics from equivalent problems in other fields, especially in ecology.

B. Facial Expression Recognition

FER is the problem of automatic emotion recognition based on facial images. Although several variants exist, the most common implementation employs static images to identify a specific set of possible emotions. These range from smile [39] or pain [13] recognition, to the most widely used emotion classification proposed by Ekman [40] (angry, disgust, fear, sad, surprise, and happy), with most publicly available datasets labeled with this codification. Although research has raised some concerns about the universality of both the underlying emotions [41] and their associated facial expressions [17], [42], the simplicity of the discrete codification and its labeling make it the most popular.

Recent developments in Deep Learning (DL) and deep convolutional networks [43], technical advances such as the training of DL models on GPUs, together with the surge of larger datasets, have allowed the end-to-end treatment of FER as a simple classification problem, trained with supervised techniques based on labeled datasets. For this reason, numerous facial expression datasets have been published to aid in the development of FER. Several reviews have focused on collecting and analyzing available datasets over the years [44], [45], but to our knowledge, none of them has reviewed their demographic properties and potential biases. The FER datasets available differ in many aspects, including the source of data (internet, media, artificial, or gathered in laboratory conditions), the image or video format and technologies (visible spectrum, infrared, ultraviolet and 3D), the type of expressions registered (micro- and macro-expression), the emotion codification (continuous, discrete, and their variants) and the elicitation method (acted, induced, or natural). Demographically speaking, some datasets openly focus on facial expressions of specific demographic groups, such as JAFFE [46], [47] (Japanese women) and iSAFE [48] (Indian

people), but for the most part the datasets have not been collected taking diversity into account.

Some recent works have already found specific biases around gender, race, and age in both commercial FER systems [49], [50] and research models [51], [52], [53], [54], [55], [56]. From these works, Kim et al. [49] focus on the age bias of commercial models in an age-labeled dataset. Ahmad et al. [50] also study commercial models, but extend the research to age, gender, and race (considering two racial categories) by employing a custom database of politician videos. Regarding research models, Xu et al. [52] study age, gender, and race bias in models trained on an Internet search-gathered dataset and evaluated on a different dataset with known demographic characteristics. Two works [53], [54] focus on gender bias in trained models. Deuschel et al. [55] analyzes biases in the prediction of action units with respect to gender and skin color in two popular datasets. Poyiadzy et al. [56] work on age bias in a dataset collected from Internet searches, performing additional work to generate apparent age labels for it. Additionally, some works [52], [54], [56] have explored mitigation strategies applicable to this problem. According to the flaws and biases found in previous studies, Hernandez et al. [57] propose a set of guidelines to assess and minimize potential risks in the application of FER-based technologies.

Notwithstanding these previous works on specific biases and the resulting general guidelines, no other work has comparatively analyzed the demographic bias of a large selection of FER datasets or used multiple metrics to account for representational and stereotypical bias. We hope that the work presented here motivates new approaches to bias detection and mitigation in FER datasets.

III. DATASET BIAS METRICS

Many specific metrics have been proposed to quantify bias and fairness [58]. Unfortunately, most of these metrics only consider the disparity in the treatment of demographic groups in the trained model predictions. This type of bias, directly related to discrimination in the legal sense, disregards the source of the disparity. Additionally, the few metrics that have focused on dataset bias [23] are only defined for binary classification problems, making the measurement of demographic biases in the source dataset an unexplored problem for more general multiclass classification problems.

In this section, our aim is to fill this gap by collecting metrics that are applicable to demographic bias in datasets, especially in multiclass problems. For this, we include both a few metrics that have previously been used in this context and, for the most part, metrics from other disciplines and contexts that can be adapted for this purpose, such as metrics from information theory (such as metrics based on *Shannon entropy*) and ecology (such as *Effective number of species*).

A. Taxonomy of Demographic Bias Metrics

We propose a taxonomy of demographic bias metrics based on the two main types of statistical demographical bias, that is, representational and stereotypical bias. This coarse classification

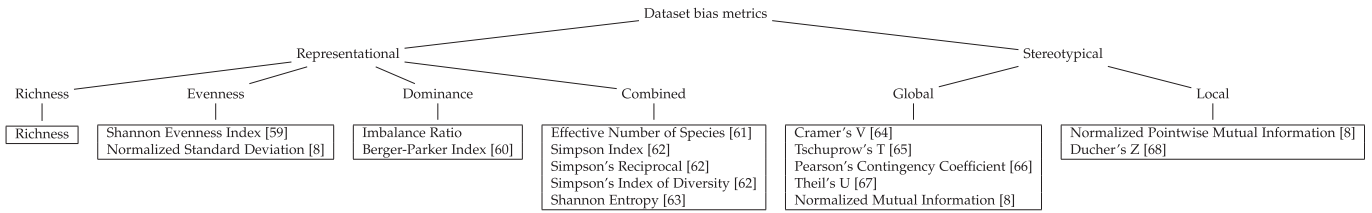


Fig. 1. Taxonomy of dataset demographic bias metrics.

is then further refined into the families of metrics that can measure each bias. The proposed taxonomy is outlined with the associated metrics in Fig. 1. The following families of metrics can be identified.

Representational Bias: Representational bias refers to a lack of general demographic diversity in the dataset, that is, an unequal representation of demographic groups, such as having more samples from male presenting people than female presenting ones. This type of bias is not related to the target variable and, therefore, can be applied to any dataset, not only to classification tasks. A similar concept can be found in the field of ecology, where diversity, and more specifically the local-scale diversity of species in an ecosystem (α -diversity [69]), is used as the opposite of representational bias. As the two concepts are directly related, one being the opposite of the other, any metric designed for diversity can be borrowed as a representational bias metric. Furthermore, in ecology diversity (and consequently representational bias) can be refined into three related components, namely *Richness*, *Evenness* and *Dominance*. Several metrics target specific components, while others measure multiple components at the same time. We classify the latter as *Combined* metrics.

- **Richness:** This category has a single metric, since it refers to the raw number of groups represented and constitutes the most direct and simple metric of representational bias. A dataset that has a representation of only a small number of groups (low richness) will be biased in favor of those groups, and susceptible to differentiated treatment in the form of inaccurate predictions in the trained model towards any group not represented. An example of a bias from lack of richness in a FER dataset would be having a dataset with only one racial group represented.
- **Evenness:** Unfortunately, even if a dataset represents many groups, this representation may not be homogeneous, having a global overrepresentation of certain groups and an underrepresentation of others. The homogeneity of the group representation is also known as *evenness*. For example, a FER dataset representationally biased by a lack of evenness would be one that, having representation of white, black and indian people, has an uneven composition of 45% each of white and black people and only 10% of Indian people.
- **Dominance:** Dominance refers to the population quota of the largest or *dominant* group in the dataset. Dominance is not independent from Richness and Evenness, but it is a robust and easy to interpret notion, making its metrics common choices for representational bias measurement. Unfortunately, in exchange, it loses a lot of information

related to the rest of the groups present, making them insufficient metrics when employed alone. For example, a FER dataset could be dominated by an age group if 95% of the population is in the range 20–25, with the remaining 5% shared among the rest of the age groups.

- **Combined:** Several metrics are not directly related to any of the other components and instead measure combinations of them. These metrics can usually summarize the amount of representational bias in a dataset, at the cost of not being able to distinguish the specific components of that bias.

Stereotypical Bias: This type of bias can be identified when working on any labeled dataset, as an association between sample variables. These variables are usually a demographic property and the target variable in classification and regression tasks, although it can also be applied to spurious associations between several demographic properties. In stereotypical bias, the under- or overrepresentation is not directly found in a global view of the dataset, but in the specific demographic composition of each of the target classes. In FER, for example, a common stereotypical bias is an overrepresentation of female presenting people in the happy class [8]. Stereotypical bias can be measured at two levels:

- **Global:** Global stereotypical bias refers to the global grade of association between a demographic component and the target classes, as a single measure of the whole dataset. For example, in FER, how related is the gender of the subjects and the target facial expressions.
- **Local:** Local stereotypical bias refers to how specific combinations of the demographic group and the target class are over- or underrepresented in the dataset. For example, in FER, if the proportion of female presenting subjects in the happy class in particular is above or below the expected proportion.

B. Considerations for the Metrics

In this work, we consider classification problems where both the demographic groups of interest and the target classes are given as nominal variables. In regression problems, where the target variable is ordinal or numerical, or when one of the demographic groups can be codified as an ordinal or numerical variable, other metrics can provide more accurate information. Despite this, ordinal and numerical variables can be reinterpreted as nominal, making it a useful first approach. A good example is found in our case study, where the age variable is codified into age groups that can be treated as a nominal variable, although at the cost of some information loss. Regarding the

TABLE I
SUMMARY OF DATASET DEMOGRAPHIC BIAS METRICS AND THEIR CHARACTERISTICS

Name	Symbol	Type	Subtype	Range	Relation to bias
Richness	R	representational	richness	$[0, \infty)$	inverse
Shannon Evenness Index [59]	SEI	representational	evenness	$(0, 1]$	inverse
Normalized Standard Deviation [8]	NSD	representational	evenness	$(0, 1]$	direct
Imbalance Ratio	IR	representational	dominance	$[1, \infty)$	direct
Berger-Parker Index [60]	BP	representational	dominance	$[1/R, 1]$	direct
Effective Number of Species [61]	ENS	representational	combined	$[1, R]$	inverse
Simpson Index [62]	D	representational	combined	$(0, 1]$	direct
Simpson's Reciprocal [62]	$1 / D$	representational	combined	$[1, R]$	inverse
Simpson's Index of Diversity [62]	$1 - D$	representational	combined	$[0, 1)$	inverse
Shannon Entropy [63]	H	representational	combined	$[0, \ln R]$	direct
Cramer's V [64]	ϕ_c	stereotypical	global	$[0, 1]$	direct
Tschuprow's T [65]	T	stereotypical	global	$[0, 1]$	direct
Pearson's Contingency Coefficient [66]	C	stereotypical	global	$[0, 1]$	direct
Theil's U [67]	U	stereotypical	global	$[0, 1]$	direct
Normalized Mutual Information [8]	NMI	stereotypical	global	$[0, 1]$	direct
Normalized Mutual Pointwise Information [8]	NPMI	stereotypical	local	$[-1, 1]$	direct
Ducher's Z [68]	Z	stereotypical	local	$[-1, 1]$	direct

target variable, although we focus on classification problems with a nominal target variable, the presented bias metrics can also be applied in other problems. Representational bias metrics do not consider the target variable, and as such can be directly applied in unlabeled datasets. Stereotypical bias metrics can also be applied in the same datasets, although in these cases they can only be used to measure the correlation between several demographic variables, rather than a demographic variable and the target variable, as is presented here.

An important consideration is the application of sampling corrections to the metrics. Many of the metrics, such as those from the field of ecology, are intended to be applied to a discrete sample of a larger population and subsequently corrected for sample size. In the case of AI datasets, we can either understand them as a sample from a global population and keep these corrections, or as complete populations and employ the uncorrected formulas. In practice, as models are mostly trained on a single dataset or a small set of them, we care more about the properties of the specific dataset than about their relationship to the real world population from which they were obtained. Thus, we do not employ the sample size correction variations, as they can hide and lower the bias of the smaller datasets. In certain techniques based on variable size datasets, such as the use of generative adversarial networks to generate datasets on demand [70], sampling corrections must be applied when evaluating biases. We consider this case outside of the scope of this work and focus on fixed-size datasets.

Taking these considerations into account, the reviewed metrics and their properties are summarized in Table I. In the table, each row includes the information corresponding to a metric, namely, the full name and references, the symbol to be used in the rest of this work, the type and subtype of the metric according to the taxonomy presented in Section III-A, the upper and lower bounds of the metric, and whether it directly (or inversely) measures bias. The metrics are further discussed in the following sections. As the metrics come from various sources, we present a unified mathematical formulation. The following unifying notation is employed:

- We define G as the set of demographic groups defined by the value of the protected attribute. For example, if G stands for *gender presentation*, a possible set of groups would be {masculine, feminine, androgynous}.
- We define Y as the set of classes of the problem.
- We define X as a population of n samples.
- We define n_a as the number of samples from the population X that have the property a . For example, n_g with $g \in G$, represents the number of samples in X that corresponds to subjects belonging to the demographic group g . a will usually be a demographic group $g \in G$, a target class $y \in Y$, or the combination of both properties $g \wedge y, g \in G, y \in Y$.
- Similarly, we define p_a as the proportion of samples from the population X that have the property a :

$$p_a = \frac{n_a}{n} .$$

C. Metrics for Representational Bias

First, let us consider the **richness** metrics.

Richness (R) [69]: Richness is the simplest and most direct metric of diversity, which can be understood as the opposite of representational bias. Mathematically, richness is defined as:

$$R(X) = |\{g \in G | n_g > 0\}| . \quad (1)$$

Although Richness is a highly informative and interpretable metric, it disregards *evenness* information on potential imbalances between group populations. In this sense, it can only assert that there are diverse examples through the dataset, but nothing about the proportions at which these examples are found. This metric is still vital for the interpretation of many other metrics, as they are either mathematically bounded by Richness or best interpreted when accompanied by it.

The following metrics focus on the *evenness* component of representational bias.

Shannon Evenness Index: (SEI) [59]. The most common example of *evenness* metrics is the Shannon Evenness Index,

a normalized version of the *Shannon entropy* designed to be robust to *Richness* variations. Due to this, only *evenness*, the homogeneity of the groups present in the population, is taken into account. It is defined as:

$$\text{SEI}(X) = \frac{H(X)}{\ln(R(X))}, \quad (2)$$

where $H(X)$ is the Shannon entropy, defined later in (7).

As the entropy is divided by its theoretical maximum, corresponding to an even population of $R(X)$ groups, the metric is bounded between 0 for uneven populations and 1 for perfectly balanced ones. This value is independent for different number of represented groups, so that datasets with different R can achieve the same SEI.

Normalized Standard Deviation (NSD) [8]: Another metric for *evenness* is the Normalized Standard Deviation of the population distribution. It is defined as:

$$\text{NSD}(X) = \frac{|G|}{\sqrt{|G|-1}} \sqrt{\frac{\sum_{g \in G} (p_g - \bar{p})^2}{|G|}}, \quad (3)$$

where \bar{p} stands for the arithmetic mean of the population profile, which for a normalized profile is $\bar{p} = 1/|G|$.

The normalization used in this metric produces the same upper and lower bounds as those of SEI. In this case, the metric is designed to target representational bias, so the meaning of the bounds is inverted, 0 for the more balanced and even datasets, and 1 for the extremely biased ones.

The third component of representational bias, *dominance*, is measured by the following metrics.

Imbalance Ratio (IR): The most common metric for class imbalance in AI is the *Imbalance Ratio*. Although its most common application is measuring target class imbalance, the same metric can be used for any partitioning of a dataset, such as the one defined by a demographic component. It is defined as the population ratio between the most represented class and the least represented class:

$$\text{IR}(X) = \frac{\max_{g \in G} n_g}{\min_{g \in G} n_g}. \quad (4)$$

This definition leads to a metric that ranges from 1 for more balanced populations to infinity for more biased ones. The inverse of this definition can be used to limit the metric between 0 (exclusive) and 1 (inclusive), with values close to 0 indicating strongly biased populations and 1 indicating unbiased ones. In this work, we employ this alternative $\text{IR}^{-1}(X)$ formulation.

This metric is commonly used for binary classification problems. When applied to more than two classes or groups, the metric simply ignores the rest of the classes, losing information in these cases.

Berger-Parker Index (BP) [60]: A metric closely related to the *Imbalance Ratio* is the Berger-Parker Index. This metric measures the relative representation of the more abundant group relative to the whole population. As IR, it does not use all information of the population distribution, as imbalances between

minority classes are not taken into account. It is defined as:

$$\text{BP}(X) = \frac{\max_{g \in G} n_g}{n}. \quad (5)$$

This metric is bounded between $1/R(X)$ and 1, with values close to $1/R(X)$ indicating representationally unbiased datasets, and 1 indicating biased ones.

Finally, some metrics measure representational bias as a *combination* of several components simultaneously.

Effective Number of Species (ENS) [61]: The Effective Number of Species is a robust measure that extends the *Richness*, keeping the same bounds but integrating additional information about *evenness*. This metric is upper bounded by $R(X)$, being equal to it for a totally balanced population, and smaller for increasingly biased populations, down to a lower bound of 1 for populations with total dominance of a single group. It is defined as:

$$\text{ENS}(X) = \exp\left(-\sum_{g \in G} p_g \ln p_g\right). \quad (6)$$

The *ENS* has several alternative formulations. The specific formula presented here is based on the *Shannon entropy*. This means that this formulation is equivalent to the *Shannon entropy* in their resulting ordering of populations. The difference lies only in the interpretability of the results, which are scaled to fit into the $[1, R(X)]$ range, with higher values indicating less representationally biased populations or datasets. The result is intended to follow the notion of an effective or equivalent number of equally represented groups. For example, a population with $\text{ENS}(X) = 1.5$ is more diverse than one with one represented group ($\text{ENS} = 1$) and less than one with two equally represented groups ($\text{ENS} = 2$).

Shannon Entropy (H) [63]: The Shannon entropy, also known as *Shannon Diversity Index* and *Shannon-Wiener Index*, can also be directly used to measure diversity. In this case, diversity is measured by the amount of uncertainty, as defined by the entropy, with which we can predict to which group a random sample belongs. It is defined as:

$$H(X) = -\sum_{g \in G} p_g \ln(p_g). \quad (7)$$

This metric lies in the range $[0, \ln(R(X))]$, where 0 identifies a population with a single represented group, and a value of $\ln(R(X))$ corresponds to a perfectly balanced dataset composed of $R(X)$ different groups.

Simpson Index (D) [62], Simpson's Index of Diversity ($1 - D$), and Simpson's Reciprocal ($1/D$). The Simpson Index is another metric influenced by both *Richness* and *evenness*. Mathematically, it is defined as:

$$D(X) = \sum_{g \in G} p_g^2. \quad (8)$$

This metric ranges from 1 for extremely biased populations with a single represented group, and approaches 0 for increasingly diverse populations, with a lower bound of $1/R(X)$. Two variants are commonly employed, the *Simpson's Index of*

Diversity defined by $1 - D(X)$, and the *Simpson's Reciprocal* defined as $1/D(X)$. Of these three, Simpson's Reciprocal index is of particular interest in this context, as it shares the same range of ENS, from 1 in populations representing biased in favor of a single group to an upper limit of $R(X)$ for more diverse populations. This metric is more influenced by the *evenness* of the population compared to ENS, especially when there are more groups present.

D. Metrics for Stereotypical Bias

The following metrics can be used to measure stereotypical bias from a *global* perspective.

Cramer's V (ϕ_C) [64], *Tschuprow's T* (T) [65] and *Pearson's Contingency Coefficient* (C) [66]: These three metrics are directly based on the Pearson's chi-squared statistic of association ($\chi^2(X)$), employed in the popular *Pearson's chi-squared test*. The $\chi^2(X)$ statistic is defined as:

$$\chi^2(X) = \sum_{g \in G} \sum_{y \in Y} \frac{(n_{g \wedge y} - \frac{n_{g \wedge y}}{n})^2}{\frac{n_{g \wedge y}}{n}}. \quad (9)$$

As the $\chi^2(X)$ calculates the difference between the real number of samples of a subgroup $n_{g \wedge y}$ and the expected number of samples of the subgroup $\frac{n_{g \wedge y}}{n}$, it detects the under or overrepresentation of specific subgroups independently of the potential representational bias in the distribution of target classes $y \in Y$ or demographic groups $g \in G$. Unfortunately, the result is both dependent on the total number of samples n and does not have clear units or intuitive bounds, making it difficult to interpret as a bias metric. Due to this, several corrections with defined bounds are available. In particular, $\phi_C(X)$, $T(X)$, and $C(X)$ are defined as:

$$\phi_C(X) = \sqrt{\frac{\chi^2(X)/n}{\min(|G| - 1, |Y| - 1)}}, \quad (10)$$

$$T(X) = \sqrt{\frac{\chi^2(X)/n}{\sqrt{(|G| - 1) \cdot (|Y| - 1)}}}, \text{ and} \quad (11)$$

$$C(X) = \sqrt{\frac{\chi^2(X)/n}{1 - \chi^2(X)/n}}. \quad (12)$$

The three metrics share the same bounds, from 0, which represents no bias or association between the demographic component and the target class, up to 1, a maximum bias or association. This difference makes them generally more meaningful and interpretable than the original statistic. Both $T(X)$ and $C(X)$ can only achieve their theoretical maximum of 1 when both nominal variables have the same number of possible values, $|G| = |Y|$. This restriction does not apply to $\phi_C(X)$, whose notion of correlation can be maximized even when $|G| \neq |Y|$. This difference makes $\phi_C(X)$ more widely used in practice.

Additionally, thresholds of significance have been provided for $\phi_C(X)$ [71], and can be also used when measuring bias. These thresholds depend on the degrees of freedom of the metric, calculated as $\text{DoF}(X) = \min(|G| - 1, |Y| - 1)$ for our application. In particular, for $\text{DoF}(X) = 1$, $\phi_C < 0.1$ is considered a

small or weak association or bias, $\phi_C < 0.3$ a medium bias, and $\phi_C < 0.5$ a large or strong bias. For $\text{DoF}(X) > 1$, the thresholds are corrected to $0.1/\sqrt{\text{DoF}(X)}$, $0.3/\sqrt{\text{DoF}(X)}$, and $0.5/\sqrt{\text{DoF}(X)}$, respectively.

Theil's U (U) [67] or *Uncertainty Coefficient* is a measure of association based on Shannon entropy, that can therefore be employed to measure stereotypical bias. It is defined as:

$$U(X, P_1 \rightarrow P_2) = \frac{H(X, P_1) - H(X, P_1|P_2)}{H(X, P_2)}, \quad (13)$$

where P_1 and P_2 stand for G and Y in any order. In this article, we will treat $P_1 = G, P_2 = Y$ as the *default* order (denoted by $U(X)$) and $P_1 = Y, P_2 = G$ as the *reverse* order (denoted by $U^R(X)$). Additionally, $H(X, P)$ with $P \in \{P_1, P_2\}$ is defined as:

$$H(X, P) = - \sum_{i \in P} p_i \ln(p_i), \quad (14)$$

and $H(X, P_1|P_2)$ is defined as:

$$H(X, P_1|P_2) = - \sum_{i \in P_1} \sum_{j \in P_2} p_{i \wedge j} \ln \left(\frac{p_{i \wedge j}}{p_j} \right). \quad (15)$$

This metric has a lower bound of 0, no bias or association between the demographic component and the target class, and an upper bound of 1, a maximum bias or association.

The definitions of stereotypical bias presented up to this point are all direction-agnostic, meaning that they produce the same result for any pair of variables, regardless of which one is provided first. A key characteristic of the Theil's U metric is that it is instead asymmetric, measuring the proportion of uncertainty reduced in one of the variables (target) when the other one (source) is known. Thus, the application of this metric could potentially establish a differentiation between forward and backward stereotypical bias.

Normalized Mutual Information (NMI) [8], [72]: A different approach to measuring the association between two variables is the use of Mutual Information based variables. In particular, a previous work [8] used a normalized variant to measure stereotypical bias in a dataset.

$$\text{NMI}(X) = \frac{\sum_{g \in G} \sum_{y \in Y} p_{g \wedge y} \ln \frac{p_{g \wedge y}}{p_g p_y}}{\sum_{g \in G} \sum_{y \in Y} p_{g \wedge y} \ln p_{g \wedge y}}. \quad (16)$$

The value of $\text{NMI}(X)$ is in the range $[0, 1]$, with 0 being no bias and 1 being total bias.

Finally, the stereotypical bias can also be measured using the following *local* metrics.

Normalized Pointwise Mutual Information (NPMI) [8], [72]: The $\text{NMI}(X)$ metric has a local variant, NPMI . This metric has a different application than the previous metrics, as it is not intended for the analysis of the stereotypical bias in the dataset as a whole. Instead, NPMI is a local stereotypical bias metric capable of highlighting the particular combination of demographic groups and the target class in which bias is found.

TABLE II
SUMMARY OF THE FER DATASETS AND THEIR CHARACTERISTICS

Abbreviation	Full Name	Year	Collection	Images	Videos	Subjects	Labelling ^a
JAFFE [46], [47]	Japanese Female Facial Expression	1998	LAB	213	—	10	6 + N
KDEF [73]	Karolinska Directed Emotional Faces	1998	LAB	4,900	—	70	6 + N
CK [74]	Cohn-Kanade Dataset	2000	LAB	8,795	486	97	7 + N
Oulu-CASIA [75]	Oulu-CASIA	2008	LAB	66,000	480	80	6
CK+ [76]	Extended Cohn-Kanade Dataset	2010	LAB	10,727	593	123	7 + N + FACS
GEMEP [77]	Geneva Multimodal Emotion Portrayals	2010	LAB	2,817	1,260	10	17
MUG [78]	Multimedia Understanding Group Facial Expression Database	2010	LAB	70,654	—	52	6 + N
SFEW [79]	Static Facial Expressions In The Wild	2011	ITW-M	1,766	—	330	6 + N
FER2013 [80]	Facial Expression Recognition 2013	2013	ITW	32,298	—	—	6 + N
WSEFEP [81]	Warsaw Set of Emotional Facial Expression Pictures	2014	LAB	210	—	30	6 + N + FACS
ADFES [82]	Amsterdam Dynamic Facial Expressions Set	2016	LAB	—	648	22	9
FERPlus [83]	FER2013 Plus	2016	ITW	32,298	—	—	6 + N
AffectNet [84]	AffectNet	2017	ITW	291,652	—	—	7 + N
ExpW [85]	Expression in-the-Wild	2017	ITW	91,793	—	—	6 + N
RAF-DB [86], [87]	Real-world Affective Faces Database	2017	ITW	29,672	—	—	6 + N
CAER-S [88]	CAER Static	2019	ITW-M	70,000	—	—	6 + N
LIRIS-CSE [89]	LIRIS Children Spontaneous Facial Expression Video Database	2019	LAB	26,000	208	12	6 + U
iSAFE [48]	Indian Semi-Acted Facial Expression	2020	LAB	—	395	44	6 + N + U
MMAFEDB [90]	Mahmoudi MA Facial Expression Database	2020	ITW	128,000	—	—	6 + N
NHFIER [91]	Natural Human Face Images for Emotion Recognition	2020	ITW	5,558	—	—	7 + N

^a 6: angry, disgust, fear, sad, surprise, and happy. 7: 6 + contempt. N: Neutral. U: Uncertain. FACS: Facial Action Coding System.

Mathematically, it is defined as:

$$\text{NPMI}(X, g, y) = -\frac{\ln \frac{p_{g \wedge y}}{p_g p_y}}{\ln p_{g \wedge y}}, \quad (17)$$

where $g \in G$ is the demographic group of interest and $y \in Y$ is the target class. The values of NPMI are in the range $[-1, 1]$, with 1 being the maximum overrepresentation of the combination of group and class, 0 being no correlation and -1 being the maximum underrepresentation of the combination.

Ducher's Z (Z) [68]: The Z measure of local association, originally developed in the field of biology, can also be employed to measure local stereotypical bias. It is defined as:

$$Z(X, g, y) = \begin{cases} \frac{p_{g \wedge y} - p_g p_y}{\min[p_g, p_y] - p_g p_y} & \text{if } p_{g \wedge y} - p_g p_y > 0 \\ \frac{p_{g \wedge y} - p_g p_y}{p_g p_y - \max[0, p_g + p_y - 1]} & \text{if } p_{g \wedge y} - p_g p_y < 0 \\ 0 & \text{otherwise,} \end{cases} \quad (18)$$

where $g \in G$ is the demographic group of interest and $y \in Y$ is the target class. The values of $Z(X)$ are also in the range $[-1, 1]$, with 1 being the maximum overrepresentation of the combination of group and class, 0 being no correlation and -1 being the maximum underrepresentation of the combination.

IV. CASE STUDY: FER DATASETS AND DEMOGRAPHIC INFORMATION

In this section, we present the experimental framework used to observe the real-world behavior of bias metrics in the FER case study. First, Section IV-A presents the selection of datasets used in this work. Then Section IV-B details the steps taken to preprocess and homogenize the different datasets. Finally, Section IV-C explains the demographic profiling of the samples in the datasets, as a necessary step to enable the application of the bias metrics.

A. Datasets

For this work, we initially considered a total of 55 datasets used for FER tasks, collected from a combination of dataset lists provided in previous reviews [44], [45], [45], datasets

cited in various works, and datasets discovered through Internet searches. This list was reduced to a final list of 20 datasets, presented in Table II, according to the following criteria:

- 1) 2D image-based datasets, or video-based datasets with per-frame labeling. This is the most extended approach to FER.
- 2) Datasets based on real images. Although some artificial face datasets are available, the demographic relabeling process can be unreliable in these contexts.
- 3) Datasets that include labels for the six basic emotions (anger, disgust, fear, happy, sadness, and surprise), and optionally neutral. This codification is the most popular in FER datasets, and the use of a unified label set makes the stereotypical biases comparable across datasets.
- 4) Availability of the datasets at the time of request.

These datasets can be categorized into three groups, depending on the source of the images:

- *Laboratory-gathered (Lab)*, which usually includes a limited selection of subjects whose images or sequences of images are taken under controlled conditions. The images in these datasets are intended for FER from inception, so the images are usually high-quality and taken in consistent environments.
- *In The Wild from Internet queries (ITW-I)*. These datasets are created from images not intended for FER learning, with varied quality. These datasets usually have a larger number of images, as their sourcing is relatively cheap.
- *In The Wild from Motion Pictures (ITW-M)*. These datasets try to improve the inconsistent quality of *ITW-I* datasets by sampling their images from motion pictures (including video from films, TV shows, and other multimedia), while retaining the advantages of a relatively high number of samples.

B. Data Preprocessing

To enable the comparison of the studied datasets, we preprocessed them to make the data as homogeneous as possible and

to ensure accurate demographic labeling. For every dataset, we performed the following steps:

- 1) *Frame extraction*: For datasets based on video, namely, ADFES, CK, CK+, GEMEP, iSAFE, LIRIS-CSE, and Oulu-CASIA, we either extracted the frames or used the per frame version when available.
- 2) *Face extraction*: Although most of the datasets provide extracted face images, the resolution and amount of margin around the face often vary considerably. To facilitate demographic prediction (see Section IV-C), we used the same face extraction methodology of FairFace [92], namely a Max-Margin (MMOD) CNN face extractor [93] implemented in DLIB¹. Face extraction was performed with a target face size of 224×224 (resized if necessary) with a margin around the face of 0.25 (56 pixels). When needed, a zero padding (black border) was used when the resized image includes portions outside the original image. On the EXPW dataset, where images have several faces, we applied the same process to each individual face extracted from the face bounding boxes provided.
- 3) *Emotion classification relabeling*: For each dataset, we consider the images corresponding to the six basic emotions [40]: angry, disgust, fear, sad, surprise, and happy, plus a seventh category for neutrality. For some datasets, the original emotion names differ (such as *angry* and *fury*). In these cases, we only rename the emotion if it is a clear synonym of the intended one. Additionally, it must be noted that not all included datasets provide examples of all emotions, with some (Oulu-CASIA and GEMEP) missing examples of neutrality.

C. Inferring Demographic Labels

For the analysis of demographic bias in datasets, it is indispensable to have demographic information of the depicted subjects. In the context of FER, the large majority of datasets do not provide or gather this information, or when they do, it is often partial information (such as in ADFES, which only refers to race and gender) or global statistics for the whole dataset and not to each sample (such as CK+). In most ITW datasets, in particular those based on Internet queries, this information is unavailable even to the original developers, as the subjects are mostly anonymous (such as in FER2013).

To overcome this limitation, we propose, following previous work [8], [94], the study of biases with respect to a proxy demographic prediction instead of the original unrecoverable information. This can be achieved through a secondary demographic model, such as the FairFace [92] face model, based on the homonymous dataset. The FairFace model is trained with the FairFace dataset, made up of 108,501 images of faces from Flickr (an image hosting service) hand-labeled by external annotators according to apparent race, gender, and age. The dataset is designed to be balanced across seven racial groups, namely White, Black, Indian, East Asian, Southeast Asian, Middle Eastern, and Latino. The dataset is also labeled with both binary gender

(Male or Female) and age group (9 age groups), although not fully balanced across these characteristics. The trained model is publicly available² and was studied against other demographic datasets, namely UTKFace, LFWA+ and CelebA, improving previous results and showing accurate classification in gender and race categories, and more modest results in the age category.

It is important to note that even for FairFace, the demographic data comes from external annotators and not self-reported demographic characteristics. Furthermore, even self-reported race and gender identities are extremely subjective and for many individuals complicated to specify [95]. This labeling is also limited to a small set of fixed categories, leaving out many possible race and age descriptors [30]. The lack of descriptors means that any classification based on these categories is fundamentally incomplete, leaving out individuals from potentially impacted minorities and misrepresenting the true identity of others.

Despite these limitations, in problems like FER, the system can only discriminate based on apparent age, gender, and race, as no true demographic information is available. Therefore, the same apparent categories can be used in our bias analysis. In such cases, FairFace provides a useful prediction of these categories, enabling a study that would otherwise be impossible.

V. RESULTS

The main objective of this section is to evaluate the behavior of dataset bias metrics when evaluating real-world datasets, in our case targeted to the FER task. To this end, we aim to answer the following research questions:

- 1) Do the different representational dataset bias metrics experimentally agree with each other? What metrics constitute the minimal and sufficient set that can characterize representational bias?
- 2) Do the different global stereotypical dataset bias metrics experimentally agree with each other? What metrics constitute the minimal and sufficient set that can characterize global stereotypical bias?
- 3) Do the local stereotypical dataset bias metrics experimentally agree? What is the most interpretable local stereotypical bias metric?
- 4) How much representational and stereotypical bias can be found using the aforementioned metrics in FER datasets? Which FER datasets are the most and less biased?

In this section, we present the main results of the case study, over the selected datasets and the three demographic components, namely age, race, and gender. First, Sections V-A, V-B and V-C focus on the agreement results of the metrics in the categories of representational, global stereotypical, and local stereotypical bias, respectively. Next, Section V-D presents the bias analysis of the datasets with respect to representational and stereotypical bias.

A. Agreement Between Representational Bias Metrics

Overview: Fig. 2 shows the values obtained for the different representational bias metrics presented in Section III-C when

¹[Online]. Available: <http://dlib.net/>

²[Online]. Available: <https://github.com/joojs/fairface>

		MMAFEDB ITW_{-1}	RAF-DB ITW_{-1}	FER2013 ITW_{-1}	EXPW ITW_{-1}	NHFIER ITW_{-1}	FER+ ITW_{-1}	AFFECTNET ITW_{-1}	GEMEP LAB	SFEW ITW_{-M}	MUG LAB	CAER-S ITW_{-M}	Outli-CASIA LAB	ADFS LAB	WSEFEP LAB	KDEF LAB	CK+ LAB	CK LAB	ISAFE LAB	LIRIS-CSE LAB	JAFFE LAB
Age	R	9	9	9	9	9	9	9	5	8	2	9	4	2	3	2	4	3	3	2	1
	ENS	5.86	5.74	5.59	5.98	5.05	5.59	5.85	4.6	4.44	1.84	2.46	2.62	1.36	1.6	1.65	1.72	1.7	1.9	1.27	1
	1/D	4.41	4.36	3.98	4.57	3.59	3.99	4.27	4.34	3.69	1.71	1.88	2.03	1.2	1.31	1.47	1.37	1.37	1.57	1.14	1
	1-D	0.77	0.77	0.75	0.78	0.72	0.75	0.77	0.77	0.73	0.42	0.47	0.51	0.17	0.24	0.32	0.27	0.27	0.36	0.12	0
	H	1.77	1.75	1.72	1.79	1.62	1.72	1.77	1.53	1.49	0.61	0.9	0.96	0.3	0.47	0.5	0.54	0.53	0.64	0.24	0
	SEI	0.8	0.8	0.78	0.81	0.74	0.78	0.8	0.95	0.72	0.88	0.41	0.7	0.44	0.43	0.72	0.39	0.48	0.58	0.34	-
	1-NSD	0.64	0.64	0.6	0.65	0.57	0.6	0.63	0.8	0.59	0.59	0.31	0.43	0.18	0.2	0.4	0.2	0.23	0.32	0.13	-
	IR ⁻¹	0.03	0.02	0.03	0.03	0.04	0.03	0.03	0.31	0.01	0.42	< 0.01	0.08	0.1	0.04	0.25	< 0.01	0.09	0.06	0.07	1
	1-BP	0.62	0.6	0.56	0.63	0.55	0.56	0.59	0.71	0.65	0.3	0.31	0.33	0.09	0.13	0.2	0.15	0.15	0.22	0.06	0
	Race	R	7	7	7	7	7	7	7	1	7	3	7	3	3	1	1	7	6	3	3
ENS		3.56	3.68	2.99	3.67	2.9	2.99	3.15	1	1.85	2.11	2.01	2.19	2.48	1	1	2.86	2.51	1.42	1.9	1.68
1/D		2.44	2.38	1.94	2.49	1.87	1.94	2.06	1	1.34	1.96	1.66	2.07	2.12	1	1	1.94	1.74	1.21	1.52	1.51
1-D		0.59	0.58	0.48	0.6	0.47	0.48	0.51	0	0.26	0.49	0.4	0.52	0.53	0	0	0.48	0.42	0.18	0.34	0.34
H		1.27	1.3	1.09	1.3	1.06	1.09	1.15	0	0.61	0.75	0.7	0.79	0.91	0	0	1.05	0.92	0.35	0.64	0.52
SEI		0.65	0.67	0.56	0.67	0.55	0.56	0.59	-	0.32	0.68	0.36	0.72	0.83	-	-	0.54	0.51	0.32	0.59	0.75
1-NSD		0.44	0.43	0.34	0.45	0.32	0.34	0.37	-	0.16	0.49	0.27	0.53	0.55	-	-	0.34	0.3	0.14	0.3	0.43
IR ⁻¹		0.03	0.05	0.02	0.03	0.02	0.02	0.02	1	< 0.01	0.03	< 0.01	0.05	0.29	1	1	0.01	0.02	0.01	0.11	0.27
1-BP		0.39	0.37	0.29	0.4	0.28	0.29	0.32	0	0.14	0.39	0.26	0.46	0.36	0	0	0.3	0.25	0.1	0.2	0.21
Gender		R	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
	ENS	1.99	1.99	2	1.99	1.99	2	2	2	1.98	1.99	1.87	1.77	1.93	2	1.9	1.95	1.96	1.99	2	1
	1/D	1.98	1.98	2	1.97	1.99	2	2	2	1.96	1.97	1.77	1.62	1.86	2	1.82	1.9	1.92	1.98	2	1
	1-D	0.49	0.5	0.5	0.49	0.5	0.5	0.5	0.5	0.49	0.49	0.44	0.38	0.46	0.5	0.45	0.47	0.48	0.49	0.5	0
	H	0.69	0.69	0.69	0.69	0.69	0.69	0.69	0.69	0.68	0.69	0.63	0.57	0.66	0.69	0.64	0.67	0.67	0.69	0.69	0
	SEI	0.99	0.99	1	0.99	1	1	1	1	0.99	0.99	0.9	0.82	0.95	1	0.93	0.96	0.97	0.99	1	-
	1-NSD	0.9	0.91	1	0.89	0.93	1	0.99	0.98	0.86	0.88	0.64	0.51	0.73	1	0.69	0.77	0.8	0.89	0.96	-
	IR ⁻¹	0.82	0.83	1	0.8	0.86	1	0.98	0.95	0.76	0.79	0.47	0.34	0.57	1	0.52	0.62	0.67	0.8	0.93	1
	1-BP	0.45	0.45	0.5	0.44	0.46	0.5	0.5	0.49	0.43	0.44	0.32	0.26	0.36	0.5	0.34	0.38	0.4	0.45	0.48	0
	Label	R	7	7	7	7	7	7	7	7	7	6	6	6	6	7	7	7	7	7	6
ENS		5.17	5.13	6.06	4.48	6.52	4.76	4.22	5.74	6.54	6.38	7	6	6	7	7	4.77	4.41	6.07	4.44	7
1/D		4.24	4.17	5.68	3.6	6.06	4	3.24	5.54	6.24	6.19	7	5.99	6	7	7	3.44	3.14	5.3	4.11	6.99
1-D		0.76	0.76	0.82	0.72	0.83	0.75	0.69	0.82	0.84	0.84	0.86	0.83	0.83	0.86	0.86	0.71	0.68	0.81	0.76	0.86
H		1.64	1.64	1.8	1.5	1.88	1.56	1.44	1.75	1.88	1.85	1.95	1.79	1.79	1.95	1.95	1.56	1.48	1.8	1.49	1.95
SEI		0.84	0.84	0.93	0.77	0.96	0.8	0.74	0.98	0.97	0.95	1	1	1	1	1	0.8	0.76	0.93	0.83	1
1-NSD		0.67	0.66	0.8	0.6	0.84	0.65	0.56	0.87	0.86	0.85	1	0.98	1	1	1	0.58	0.55	0.77	0.7	0.99
IR ⁻¹		0.11	0.06	0.06	0.03	0.31	0.02	0.03	0.45	0.29	0.12	1	0.91	1	1	0.99	0.08	0.07	0.17	0.01	0.9
1-BP		0.68	0.61	0.74	0.62	0.74	0.63	0.53	0.78	0.8	0.82	0.86	0.83	0.83	0.86	0.86	0.51	0.48	0.68	0.7	0.85

Fig. 2. Representational bias metrics for the three demographic components and the target label. The metrics are calculated as diversity metrics, with higher values corresponding to lower representational bias. The graphical representations of the values are normalized to the maximum value of the row. The datasets are sorted by the average of the normalized metrics.

applied to each of the datasets considered, in each of the three potential demographic axes. Additionally, we apply the metrics to the distribution of the target label, as an example of their use on non-demographic components. As most metrics in this category are diversity metrics, bias metrics, such as D, NSD, and BP, are computed in their complementary form based on their upper limit 1, that is, $1 - D$, $1 - NSD$ and $1 - BP$, respectively, to allow for a more direct comparison. IR is computed in its inverse form, IR^{-1} , as it has a closer range to the other metrics. For a better visualization, the plotted values are normalized by the maximum value of each line (metric in each component). The datasets are sorted by the average of these normalized values in decreasing order, from higher average values, which indicate less representational bias, to lower values which indicate more representational bias.

Globally, we can observe that the ranking of the datasets according to the different metrics is mostly consistent, suggesting a high agreement between the metrics even if the scales vary.

A marked exception occurs in the combinations of dataset and component where a single group is represented ($R(X) = 1$), such as JAFFE in the age and gender components, and GEMEP, WSEFEP, and KDEF, in the race component, where IR^{-1} reports high values, corresponding to a situation of high diversity and low bias, contrary to the intuitive notion that these datasets are strongly biased in these components. These results are caused by a strict implementation of the IR metric, since in these datasets and components the most and least represented groups are the same. Additionally, it can be observed that IR^{-1} has a different behavior in the gender component, where the values are similar to those of the other metrics, compared to the other components. This can be explained by the number of groups in each component, with only two groups in the gender component while the rest of the components have more, as well as by the balance in the components, with the gender component showing less representational bias overall. In these cases, the limitations of the IR^{-1} metric have less impact in this context, when applied

to components with few groups represented in roughly equal proportions.

Similar to the behavior of IR^{-1} , some of the metrics, such as SEI and NSD, are directly ill-defined in the the trivial cases where a single demographic group is present in the dataset. Metrics related to R, such as ENS, 1/D or R itself, are instead robust to this cases, simply confirming the intuition that only one group is represented, and giving a result always lower than for datasets with more groups represented, no matter their evenness.

Interpretability: To analyze the interpretability of these metrics one of the key factors is the range of the metric, as the bounds contextualize the values of a metric. Some of the metrics, such as R, ENS, and 1/D, are bounded in a $[1, R]$ range, immediately interpretable as a number of represented groups. From these, ENS and 1/D complement the pure richness information with the evenness of the population, lowering the value when some of the groups are underrepresented. In the case of metrics focused on a single characteristic of representational bias, such as evenness (SEI, NSD) and dominance (BP, IR^{-1}), the most common range is $(0,1)$ (including or excluding the bounds). This unitary range is easy to interpret for these characteristics, and allows for a quick conversion between diversity and bias when one of the meanings is preferred, as we have done with D , NSD and BP . Among the metrics associated with this unitary range, we can observe that IR^{-1} tends to exaggerate the biases compared to the other metrics in components with more than two groups (all but gender). This results in IR^{-1} values close to 0 in these components for most datasets. Finally, the H value lies in the range $[0, \ln(R(X))]$, with no natural interpretation. As the information conveyed by H is the same as the one in the ENS, while being less interpretable, ENS is generally preferable.

Statistical Agreement Assessment: To evaluate the agreement between the metrics, we employ the Spearman’s ρ , a measure of the strength and direction of the monotonic relationship of two variables interpreted as ordinal, that is, where only the ranking produced is considered. To compute the pairwise agreement between metrics, we employ Spearman’s ρ to compare each pair of metrics on each component (using the results in Fig. 2) and then average across the four components. This procedure highlights the similarity across the metrics. Pairs of metrics that order the bias in the same way (indicating redundancy among the metrics) produce values of ρ close to 1 or -1 , depending on the direction of the relationship, and less related metrics produce values close to 0. Fig. 3 shows these agreement results between the different metrics of representational bias.

Except for R and IR^{-1} , the agreement value is greater than 0.71 for most metrics, with an average of 0.88. In the case of R, we can observe very low correlations with the rest of the metrics. The low robustness of this metric, where adding a single example of a missing demographic group to a whole dataset will produce a change in its value, accounts for these low correlations. In the same way, IR^{-1} is especially unreliable on demographic axes with many potential groups, such as race and gender, becoming overly sensitive to the representation of the least represented group.

Outside of these two exceptions, the rest of the metrics mostly conform to the taxonomy presented in Section III-A. SEI and

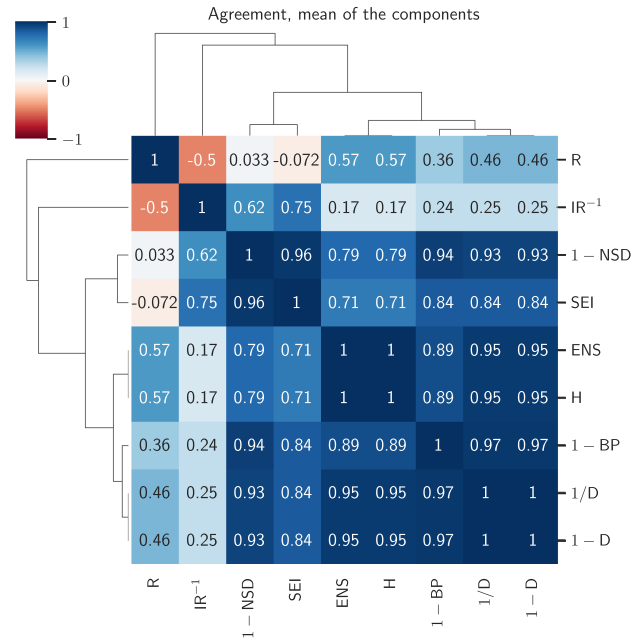


Fig. 3. Spearman’s ρ agreement between the representational bias metrics, measured independently for each component and then averaged for each pair of metrics. Higher ρ values indicate high coherence between the rankings generated by the metrics.

NSD are strongly correlated ($\rho = 0.96$), as both measure evenness, and both of them are, as expected, independent of R, with $|\rho| < 0.1$ in the two cases. The other cluster of metrics with high correlation is mainly made up of combined metrics, namely, ENS, H, 1/D and 1-D. Unexpectedly, the other dominance metric, BP, is also included in this cluster, with high correlation with the others (between 0.85 and 0.96). This metric, different from IR^{-1} in that it does not directly consider the least represented group, is more reliable and appears to experimentally capture information similar to the rest of the metrics. Considering the simplicity and high interpretability of BP, this makes it an interesting metric for general representational bias characterization.

It is also noteworthy the high correlation between the combined and evenness metrics clusters, showing how in this particular case of study both components of representational bias are highly related.

Recommended Metrics: In general, characterizing the bias of a dataset for a given component appears to be appropriately summarized by one of the combined metrics bounded by R, namely, ENS and 1/D, as these combined metrics are highly interpretable variants of the pure richness, and one of the evenness metrics to highlight this independent component, such as SEI or NSD. For more succinct analysis where evenness and richness are not expected to substantially disagree, Berger-Parker appears to maintain a high correlation with other metrics while having an intuitive meaning and simple implementation.

B. Agreement Between Stereotypical Bias Metrics

Overview: Fig. 4 shows the stereotypical bias metric results for the three main demographic components against the target

	WSEFEP LAB	ADFS LAB	KDEF LAB	JAFFE LAB	Outliu-CASIA LAB	MUG LAB	CAERS ITW-M	EXPW ITW-I	CK+ LAB	MMAFEDB ITW-I	CK LAB	AFFECTNET ITW-I	FER+ ITW-I	FER2013 ITW-I	iSAFE LAB	RAF-DB ITW-I	GEMEP LAB	NHFIER ITW-I	SFEW ITW-M	LIRIS-CSE LAB
Label - Age	ϕ_C 0	0	.002	-	.021	.053	.052°	.072°	.125°	.07°	.137°	.095°	.105°	.104°	.117°	.144 Δ	.165 Δ	.138 Δ	.115°	.134°
	T 0	0	< .001	-	.019	.034	.048	.067	.105	.066	.104	.088	.097	.097	.089	.134	.156	.128	.111	.089
	C 0	0	.002	0	.037	.053	.126	.174	.212	.17	.19	.226	.248	.247	.163	.334	.314	.32	.271	.132
	U < .001	< .001	< .001	0	< .001	< .001	.004	.009	.011	.009	.009	.017	.021	.018	.008	.036	.037	.031	.019	.007
	U ^R 0	< .001	< .001	-	< .001	.002	.009	.008	.032	.008	.026	.014	.019	.019	.023	.034	.042	.035	.024	.042
	NMI < .001	< .001	< .001	0	< .001	< .001	.003	.004	.008	.004	.007	.008	.01	.009	.006	.018	.02	.017	.011	.006
Label - Race	ϕ_C -	0	-	.039	.038	.035	.04	.041	.073°	.058°	.096°	.041	.068°	.086°	.197°	.058°	-	.087°	.076°	.256 Δ
	T -	0	-	.025	.03	.026	.04	.041	.073	.058	.092	.041	.068	.086	.149	.058	-	.087	.076	.204
	C 0	0	0	.039	.054	.049	.096	.1	.175	.141	.21	.099	.164	.206	.268	.14	0	.208	.184	.34
	U 0	< .001	0	< .001	< .001	< .001	.003	.013	.006	.019	.004	.009	.012	.021	.006	0	0	.011	.01	.051
	U ^R -	0	-	.002	.002	.002	.007	.004	.019	.008	.031	.004	.013	.02	.107	.008	-	.02	.03	.118
	NMI 0	< .001	0	< .001	< .001	< .001	.002	.002	.008	.004	.012	.002	.006	.007	.018	.004	0	.007	.008	.037
Label - Gender	ϕ_C 0	0	.002	-	.026	.054	.138°	.157°	.023	.167°	.053	.195°	.171°	.178°	.108°	.131°	.093	.172°	.261°	.313 Δ
	T 0	0	.001	-	.018	.035	.088	.1	.014	.107	.034	.125	.109	.114	.069	.084	.062	.11	.167	.21
	C 0	0	.002	0	.026	.054	.136	.155	.023	.164	.053	.192	.169	.176	.107	.13	.093	.17	.253	.299
	U < .001	< .001	< .001	0	< .001	< .001	.005	.008	< .001	.009	< .001	.013	.009	.009	.003	.005	.002	.008	.018	.034
	U ^R < .001	0	< .001	-	< .001	.002	.015	.018	< .001	.02	.002	.028	.021	.023	.009	.012	.006	.022	.05	.074
	NMI 0	< .001	< .001	0	< .001	< .001	.004	.006	< .001	.006	< .001	.009	.007	.006	.002	.004	.002	.006	.014	.024

Less biased datasets ← Datasets sorted by mean normalized values → More biased datasets

Fig. 4. Stereotypical bias metrics for the three demographic components against the target label. Higher values correspond to higher amounts of stereotypical bias. The graphical representation at each row, corresponding to a single metric and demographic component, is normalized to the maximum value of the row. In the ϕ_C row a $^\circ$ mark indicates a statistically weak association and a Δ mark a statistically medium association. The datasets are sorted by the average of the normalized metrics, from lower values (less stereotypical bias) in the left to higher values (more stereotypical bias) in the right.

label. These metrics target bias directly, thus, higher values, shown in the right, relate to more stereotypically biased datasets. The graphical representation of the values is normalized by the maximum value of each line (metric in each component). Following the thresholds for ϕ_C presented in Section III-D, we indicate the bias strength in the row ϕ_C . $^\circ$ marks a weak bias, and Δ marks a medium bias. No strong bias are found according to these thresholds.

The ranking produced by these stereotypical bias metrics is highly coherent, especially in the gender component, where a binary classification is used. According to the magnitude of the values reported, two different groups are observed. The three metrics based on the χ^2 metric, namely, ϕ_C , T, and C, all report values in a similar range for all datasets. The NMI and the U metric, applied both in the forward and reverse direction, noted here as U^R, tend to report coherent but lower values than the other three metrics. U and U^R result in similar values and rankings, with a few exceptions, namely, the CK+, CK, iSAFE, and LIRIS-CSE datasets in the age component, and the iSAFE dataset in the race component.

Interpretability: All of the stereotypical bias metrics lie in the range [0,1] and share the same interpretation of these bounds, namely, 0 for no bias and 1 for maximal bias, which allow for a simple interpretation. The distribution of the values, as mentioned earlier, is different, with U, U^R and NMI resulting in very low values for all of the datasets, which is indicative of a low sensitivity. The three other metrics, ϕ_C , T and C, all show higher values, related to a higher sensitivity. Additionally, for ϕ_C both the availability of predefined thresholds and the fact that it can be maximized even when the number of demographic

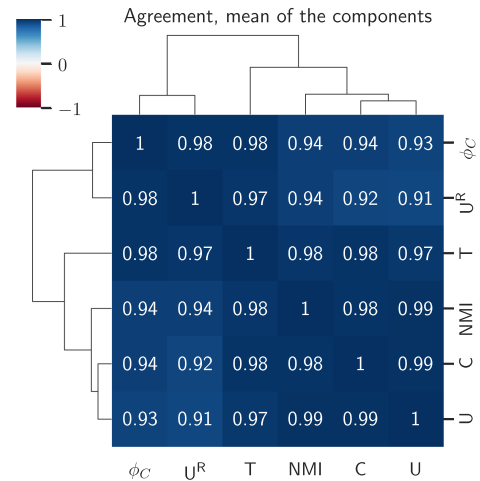


Fig. 5. Spearman's ρ agreement between the stereotypical bias metrics, measured independently for each component and then averaged for each pair of metrics. Higher ρ values indicate high coherence between the rankings generated by the metrics.

groups and target labels is not equal (see Section III-D) improve its interpretability.

Statistical Agreement Assessment: Fig. 5 uses Spearman's ρ to compare the agreement between the different stereotypical bias metrics, measured between the three main demographic components and the output label. As intuitively observed in Fig. 4, the agreement values are high in all cases, with a minimum of 0.91, between U and U^R. It can be concluded that for this case

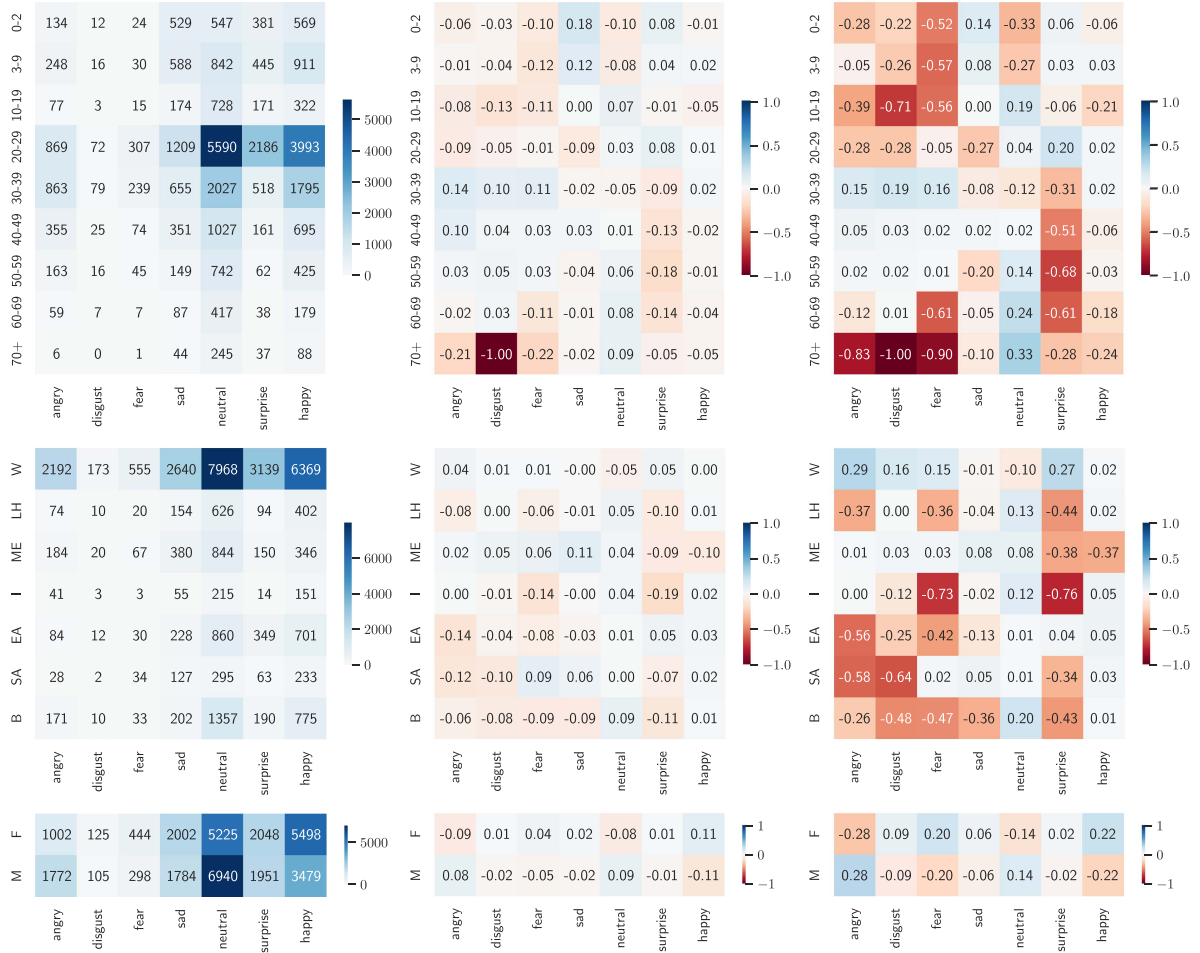


Fig. 6. Local stereotypical bias metrics for the FER+^{ITW} dataset across the three demographic components. The first column corresponds to the raw number of samples in each combination of demographic group and target label, the second column corresponds to the NPMI metric and the third column to the Z metric. For the two metrics, higher absolute values indicate higher local stereotypical bias. Negative values indicate underrepresentation and positive values overrepresentation.

study all metrics provide similar information about the amount of information shared between the three demographic components and the output classes.

Recommended Metrics: We observe no notable differences or clusters in the stereotypical bias metrics. As the ranking of biases is experimentally similar, any of the metrics is sufficient for comparing the bias between datasets. In spite of this, we suggest using ϕ_C for stereotypical bias characterization, with a more natural range of values and predefined thresholds for determining the intensity of the bias, providing a more interpretable result.

C. Agreement Between Local Stereotypical Bias Metrics

Overview: In the case of local stereotypical bias the metrics report not a single value for a dataset and component combination, but a matrix, making impractical to include the full results here. Due to this, Fig. 6 shows an example of the application of local stereotypical bias metrics, in this case for the FER+ dataset, while the results for the rest of the datasets are included in the Supplementary Material. To illustrate this case, we show

the support of each subgroup (first column), in addition to the two stereotypical bias metrics, namely, NPMI (second column) and Z (third column). The support of the subgroups contains the full information needed to identify the local biases, but are hard to interpret manually, especially in cases where the target label, the demographic groups, or both of them are unbalanced or representationally biased, such as in this example.

Statistical Agreement Assessment: Overall, we can intuitively observe a high correlation in the rank sorting of the biases detected by the two metrics, in this case when comparing the results within a single dataset and demographic component. We can again confirm this correlation by using the Spearman's ρ . In particular, for each dataset and demographic component we compute a NPMI and a Z matrix and compute the ρ agreement score between them. The final averaged value is $\rho = 0.96 \pm 0.02$, indicating a strong correlation between the two metrics.

Interpretability: Despite the similarity in the produced rankings, the range of values of both metrics is remarkably different. In particular, the chosen example dataset highlights the main shortcoming of the NPMI metric, an oversensitivity to missing

	Lab											ITW-I										ITW-M		
	MUG	GEMEP	ADFS	Oulu-CASIA	KDEF	CK+	CK	WSEFEP	iSAFE	LIRIS-CSE	JAFFE	Average	EXPW	MMAFEDB	RAF-DB	AFFECTNET	FER+	FER2013	NHFER	Average	SFEW	CAER-S	Average	
Age	9 – ENS	7.164	4.395	7.644	6.376	7.351	7.275	7.298	7.404	7.100	7.732	8.000	7.067 ± 0.932	3.017	3.144	3.258	3.149	3.407	3.414	3.954	3.334 ± 0.286	4.561	6.541	5.551 ± 0.990
	1 – SEI	0.124	0.051	0.561	0.304	0.279	0.607	0.516	0.574	0.416	0.657	–	0.409 ± 0.200	0.186	0.196	0.205	0.196	0.216	0.217	0.263	0.211 ± 0.024	0.283	0.590	0.437 ± 0.154
	ϕ_C	0.053	0.165 ^Δ	0.000	0.021	0.002	0.125 ^Δ	0.137 ^Δ	0.000	0.117 ^Δ	0.134 ^Δ	–	0.075 ± 0.063	0.072 ^Δ	0.070 ^Δ	0.144 ^Δ	0.095 ^Δ	0.105 ^Δ	0.104 ^Δ	0.138 ^Δ	0.104 ± 0.027	0.115 ^Δ	0.052 ^Δ	0.084 ± 0.031
Race	7 – ENS	4.889	6.000	4.522	4.806	6.000	4.141	4.488	6.000	5.581	5.096	5.320	5.168 ± 0.634	3.334	3.445	3.317	3.848	4.011	4.014	4.100	3.724 ± 0.321	5.151	4.990	5.070 ± 0.080
	1 – SEI	0.320	–	0.174	0.285	–	0.460	0.486	–	0.682	0.414	0.252	0.384 ± 0.151	0.332	0.348	0.330	0.410	0.437	0.438	0.453	0.393 ± 0.050	0.684	0.641	0.663 ± 0.021
	ϕ_C	0.035	–	0.000	0.038	–	0.073 ^Δ	0.096 ^Δ	–	0.197 ^Δ	0.256 ^Δ	0.039	0.092 ± 0.083	0.041	0.058 ^Δ	0.058 ^Δ	0.041	0.068 ^Δ	0.086 ^Δ	0.087 ^Δ	0.063 ± 0.018	0.076 ^Δ	0.040	0.058 ± 0.018
Gender	2 – ENS	0.014	0.001	0.074	0.233	0.098	0.055	0.039	0.000	0.012	0.001	1.000	0.139 ± 0.280	0.013	0.010	0.008	0.000	0.000	0.000	0.006	0.005 ± 0.005	0.019	0.129	0.074 ± 0.055
	1 – SEI	0.010	0.000	0.054	0.179	0.073	0.040	0.029	0.000	0.009	0.001	–	0.039 ± 0.052	0.009	0.007	0.006	0.000	0.000	0.004	0.004 ± 0.004	0.013	0.096	0.055 ± 0.041	
	ϕ_C	0.054	0.093	0.000	0.026	0.002	0.023	0.053	0.000	0.108 ^Δ	0.313 ^Δ	–	0.067 ± 0.090	0.157 ^Δ	0.167 ^Δ	0.131 ^Δ	0.195 ^Δ	0.171 ^Δ	0.178 ^Δ	0.172 ^Δ	0.167 ± 0.018	0.261 ^Δ	0.138 ^Δ	0.199 ± 0.062

Fig. 7. Summary of demographic bias in FER datasets, using the selected metrics. Higher values in any metric indicate a higher amount of bias. The representational bias metrics are shown in their bias formulation, in the case of ENS by subtracting it from the number of groups in each demographic component, and in the case of SEI from 1, its theoretical maximum. The datasets in each group are sorted by the average of the normalized metrics.

subgroups, such as the 70+ age group in the disgust label. These underrepresentation biases obtain a NPMI of -1 , but arguably similar biases, such as the small representation of the 70+ age group in the fear label, obtain only a NPMI of -0.22 , scaling in an unintuitive way. In the Z metric, we can observe values of -1 and -0.9 for these same groups, closer to the natural intuition. Generally speaking, the observed values for Z are better distributed in the range, making them more interpretable than those of the NPMI.

Recommended Metrics: From these results, we recommend the usage of Z to evaluate local stereotypical bias, as the biases detected are similar to those of the NPMI, but having a better interpretability.

D. Demographic Bias in FER Datasets

In this section, we focus on analyzing the FER case study, employing the information provided by the metrics applied in the previous sections to identify which biases are the most prevalent across the datasets, and which datasets are the least affected by them. For the analysis in this section we group the datasets according to the source of the data in each dataset (according to the classification presented in Section IV-A), as this is one of the key factors determining the type of biases exhibited.

Fig. 7 summarizes the findings in Figs. 2 and 4, according to the suggested metric selection of the previous sections. In particular, for representational bias we show the ENS and SEI to characterize representational bias, and ϕ_C to characterize stereotypical bias. In the case of ENS and SEI, we convert them to their bias variant, to facilitate their interpretation. In the case of ENS we use $|G| - \text{ENS}$, where $|G|$ is the number of groups in the demographic component, with the meaning of equivalent number of *unrepresented* groups. In the case of SEI, we use $1 - \text{SEI}$, with the meaning of *unevenness* between represented groups. The results are split in three sections, according to the data source of each dataset (LAB, ITW-I and ITW-M).

Representational Bias: Overall, we observe the lowest representational bias values in the ITW-I datasets, followed by the ITW-M and finally the LAB datasets. Gender is the least biased component across the three groups, followed by age,

and race exhibits the highest bias. We observe almost no representational bias in the gender component for all datasets, with most having a $|G| - \text{ENS}$ close to 0. The exception to this is JAFFE, a laboratory-gathered dataset taken from a small sample of Japanese women that was never intended to be used as a general dataset for ML training [47]. In the age and race components, the biases are more generalized, with the LAB datasets exhibiting the highest representational bias, while the ITW-M and ITW-I datasets seem to be less biased. Despite this, the ITW-M appear to be closer to the representation profiles of the LAB datasets than to those of the ITW-I datasets. Globally, EXPW, MMAFEDB and RAF-DB are the least representationally biased datasets, with low $|G| - \text{ENS}$ in the age (≤ 3.258) and the race (≤ 3.445) components. For the evenness aspect of representational bias, we can observe the results of $1 - \text{SEI}$, which is relatively low and homogeneous across the ITW-I datasets. The two ITW-M datasets are less even across the represented groups than either the ITW-I and LAB datasets. Finally, the LAB datasets have a large and less coherent range of evenness values, showing a larger variety of demographic profiles between them.

Global Stereotypical Bias: Overall, stereotypical bias seems to be present at a lower rate than representational bias for the three groups, with only weak and medium bias found in some cases according to the ϕ_C metric thresholds. The least biased group is the LAB group in two of the components (gender and age), while the ITW-M is the least biased in the remaining one (race). As expected, this type of bias is almost absent from most LAB datasets, as they usually take samples for all the target classes for each subject. Despite this, some LAB datasets, such as LIRIS-CSE, GEMEP, iSAFE, CK, and CK+, do not follow this rule. These datasets only include certain classes for each subject, and in these cases they exhibit as much stereotypical bias as the ITW-I datasets. This effect is present mostly in the age and race components, while in the gender category these LAB datasets still exhibit lower bias scores than the ITW-I ones. ITW-I datasets have an overall higher stereotypical bias, with most of them showing weak or medium bias in the three demographic components, especially for the gender component. From the ITW-I datasets, the least stereotypically biased are EXPW and

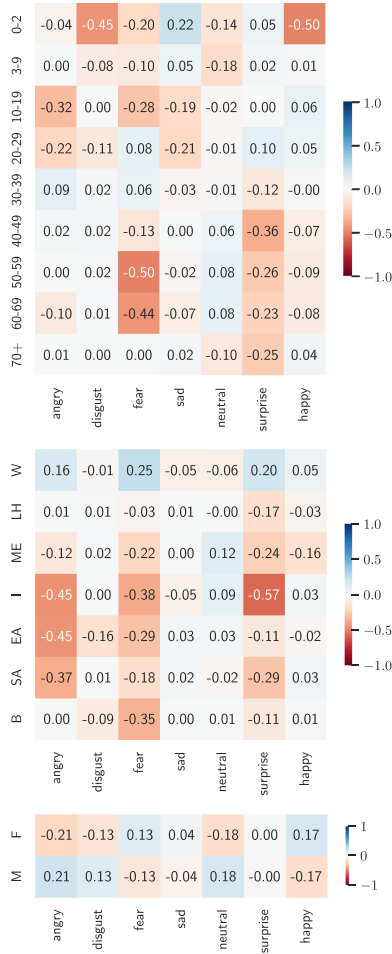


Fig. 8. Local stereotypical bias according to Z for the EXPW^{ITW} dataset in each of the demographic components.

MMAFEDB. The two ITW-M, namely, SFEW and CAER-S, show very different stereotypical bias profiles, with SFEW having higher stereotypical bias scores in the three demographic components, while CAER-S is unbiased in the race component and less biased in the two other components compared to SFEW.

Local Stereotypical Bias: According to these results, we can highlight EXPW as one of the least biased datasets overall. For this reason, in Fig. 8 we measure its local stereotypical bias using Z. The highest biases are found in the age and race components, where six subgroups have large underrepresentations ($Z \leq -0.45$), namely the 0–2 age group in the disgust and happy labels, the 50–59 group in the fear label, the east asian and indian groups in the angry label, and the indian group in the surprise label. The angry, fear, and surprise labels exhibit the highest overall local stereotypical biases, especially in the race and age components.

FER Dataset Selection: With this information, currently no completely unbiased dataset is found for FER. Although ITW datasets such as RAF-DB, MMAFEDB, and EXP have relatively low representational bias and high evenness, their stereotypical bias results demand caution when using them as a sole source of training data. From these, EXPW can be highlighted as the least

stereotypically biased dataset of the ITW category, and having a relatively large number of samples at 91,793 images, it can be generally recommended. When using this dataset, precautions should be taken to evaluate the potential impact of stereotypical bias, especially in the angry, fear, and surprise labels in the race component.

VI. CONCLUSION AND FUTURE WORK

In this work, we have proposed a taxonomy of metrics applicable to some of the main types of demographic biases found in datasets, namely, representational and stereotypical bias. We have incorporated into our review both metrics previously employed for this purpose and new proposals adapted from other fields, such as information theory and ecology. In particular, we have shown how metrics intended for species diversity in ecosystems can be also used to measure representational bias in datasets. After presenting these metrics, we have employed FER as a case of study, comparing the biases present in 20 datasets for this task to evaluate the behavior of the metrics. With this information, we are able to highlight the most interpretable metrics for each subtype of bias, while avoiding redundant metrics that do not offer additional information.

Regarding representational bias, defined as an unequal representation of different demographic groups, we have found a relatively high experimental agreement between the different metrics available. Despite the diverse theoretical basis and implementation details, we conclude that the joint usage of a combined metric, where we suggest the Effective Number of Species (ENS) for its interpretability, and an evenness metric, such as the Shannon Evenness Index (SEI), selected for the same reason, are generally sufficient to characterize representational bias.

Regarding stereotypical bias, an undesired association between a demographic component and the target label, we find no significant difference between the metrics available. For interpretability reasons, we recommend the usage of Cramer’s V (ϕ_C) when measuring stereotypical bias in a dataset. Additionally, for the metrics used to measure stereotypical bias in a local way, that is, for a specific subgroup defined by both a demographic group and a target label, we suggest using Ducher’s Z, as an interpretable and informative metric.

As a case study, we have applied these metrics to a collection of twenty FER datasets. We find evidence of representational bias in most of the datasets, especially those taken in laboratory conditions, as the low subject number and collection conditions lead to constrained demographic profiles. The datasets taken in the wild, especially those from internet searches, exhibit lower representational bias, but at the cost of higher stereotypical bias. Overall, we find that the EXPW dataset exhibits the lowest representational bias, with relatively low stereotypical bias. Furthermore, we apply a local stereotypical bias metric to identify the specific stereotypical biases that could be of concern, and find that special considerations should be taken when analyzing the angry, fear and surprise labels, as they are racially biased.

For future work, we can note that although this work has focused on representational and stereotypical bias, dataset

demographic bias can manifest in many other ways [96], such as image quality, image context and label quality, and further research is still needed in the way these other manifestations can be measured. Additionally, we have only considered both demographic components and target variables as nominal variables, but our research could be extended to continuous and ordinal demographic components and regression problems. Furthermore, we have focused on stereotypical bias as an association between demographic components and target classes, although the association between several demographic components also poses potential risks and can be measured with the same metrics. Currently, there has been no research on the potential impact of biases in this regard.

Another limitation that could be improved in future work is the usage of demographic labels derived from a demographic relabeling model, FairFace. As these labels can be inaccurate, new datasets that include demographic information, or new models capable of more accurate demographic predictions, could support more robust bias analysis.

The dataset bias analysis found here is supported by previous work [8], [94] that has shown bias transference to the final trained model in the specific context of FER. Nonetheless, further work is still required to comprehend the implications and reach of this bias transference in different problems and contexts.

In this work, we have focused on how to measure demographic bias, but there is still work to be done on the usage of the reviewed metrics to study the transference of bias from the dataset to the model, as a way to improve bias mitigation strategies. To this end, we suggest the generation of intentionally biased synthetic datasets derived from real datasets as a general application-agnostic framework to evaluate both the limits of these demographic bias metrics and potential mitigation strategies. In this sense, and to inspire future works, in the Supplementary Material, we provide the results of a series of experiments showcasing this methodology and showing how different types of dataset bias, measured with our proposed metrics, propagate in different ways to the final model predictions.

Finally, it is crucial to understand any bias results not only as a statistical discovery but as a potential harm to real people. In this sense, more work is needed on the psychological and sociological impact of potentially biased systems in the final applications.

REFERENCES

- [1] A. Jobin, M. Ienca, and E. Vayena, "The global landscape of AI ethics guidelines," *Nat. Mach. Intell.*, vol. 1, no. 9, pp. 389–399, Sep. 2019.
- [2] A. Winfield, "Ethical standards in robotics and AI," *Nat. Electron.*, vol. 2, no. 2, pp. 46–48, Feb. 2019.
- [3] B. C. Stahl and D. Wright, "Ethics and privacy in AI and Big Data: Implementing responsible research and innovation," *IEEE Secur. Privacy*, vol. 16, no. 3, pp. 26–33, May/June 2018.
- [4] S. Mitchell, E. Potash, S. Barocas, A. D'Amour, and K. Lum, "Algorithmic Fairness: Choices, assumptions, and definitions," *Annu. Rev. Statist. Appl.*, vol. 8, no. 1, pp. 141–163, Mar. 2021.
- [5] E. Ntoutsis et al., "Bias in data-driven artificial intelligence systems—An introductory survey," *WIREs Data Mining Knowl. Discov.*, vol. 10, no. 3, May 2020, Art. no. e1356.
- [6] S. Verma and J. Rubin, "Fairness definitions explained," in *Proc. ACM Int. Workshop Softw. Fairness*, Gothenburg Sweden, 2018, pp. 1–7.
- [7] J. Zhao, T. Wang, M. Yatskar, V. Ordonez, and K.-W. Chang, "Men also like shopping: Reducing gender bias amplification using Corpus-level Constraints," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, Copenhagen, Denmark, 2017, pp. 2979–2989.
- [8] I. Dominguez-Catena, D. Paternain, and M. Galar, "Assessing demographic bias transfer from dataset to model: A case study in facial expression recognition," in *Proc. Workshop Artif. Intell. Saf.*, Vienna, Austria, 2022. [Online]. Available: <https://ceur-ws.org/Vol-3215/12.pdf>
- [9] M. Mitchell et al., "Model cards for model reporting," in *Proc. Conf. Fairness Accountability Transparency*, 2019, pp. 220–229.
- [10] C. Dulhanty and A. Wong, "Auditing ImageNet: Towards a model-driven framework for annotating demographic attributes of large-scale image datasets," Jun. 2019, *arxiv.1905.01347*.
- [11] J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in *Proc. 1st Conf. Fairness Accountability Transparency*, ser. Proceedings of Machine Learning Research, S. A. Friedler and C. Wilson, Eds., PMLR, 2018, pp. 77–91. [Online]. Available: <https://proceedings.mlr.press/v81/buolamwini18a.html>
- [12] B. Kim, H. Kim, K. Kim, S. Kim, and J. Kim, "Learning not to learn: Training deep neural networks with biased data," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Long Beach, CA, USA, 2019, pp. 9004–9012.
- [13] P. Werner, D. Lopez-Martinez, S. Walter, A. Al-Hamadi, S. Gruss, and R. W. Picard, "Automatic recognition methods supporting pain assessment: A survey," *IEEE Trans. Affective Comput.*, vol. 13, no. 1, pp. 530–552, First Quarter, 2022.
- [14] R. Nimmagadda, K. Arora, and M. V. Martin, "Emotion recognition models for companion robots," *J. Supercomput.*, vol. 78, pp. 13710–13727, Mar. 2022.
- [15] R. Berk, H. Heidari, S. Jabbari, M. Kearns, and A. Roth, "Fairness in criminal justice risk assessments: The state of the art," *Sociol. Methods Res.*, vol. 50, no. 1, pp. 3–44, Jul. 2018.
- [16] R. E. Jack, C. Blais, C. Scheepers, P. G. Schyns, and R. Caldara, "Cultural confusions show that facial expressions are not universal," *Curr. Biol.*, vol. 19, no. 18, pp. 1543–1548, Sep. 2009.
- [17] H. A. Elfenbein and N. Ambady, "On the universality and cultural specificity of emotion recognition: A meta-analysis," *Psychol. Bull.*, vol. 128, no. 2, pp. 203–235, 2002.
- [18] D. Danks and A. J. London, "Algorithmic bias in autonomous systems," in *Proc. 26th Int. Joint Conf. Artif. Intell.*, Melbourne, Australia, 2017, pp. 4691–4697.
- [19] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning*. Cambridge, MA, USA: MIT Press, 2019.
- [20] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," *ACM Comput. Surv.*, vol. 54, no. 6, pp. 1–35, Jul. 2021.
- [21] R. Schwartz, A. Vassilev, K. Greene, L. Perine, A. Burt, and P. Hall, "Towards a standard for identifying and managing bias in artificial intelligence," National Institute of Standards and Technology, Tech. Rep., Mar. 2022.
- [22] A. Torralba and A. A. Efros, "Unbiased look at dataset bias," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.* 2011, 2011, pp. 1521–1528.
- [23] S. Das et al., "Fairness measures for machine learning in finance," *J. Financial Data Sci.*, vol. 3, no. 4, pp. 33–64, Oct. 2021.
- [24] R. Thomas and D. Uminsky, "The problem with metrics is a fundamental problem for AI," Feb. 2020, *arxiv.2002.08512*.
- [25] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," in *Proc. 30th Int. Conf. Neural Inf. Process. Syst.*, Red Hook, NY, USA, 2016, pp. 3323–3331.
- [26] M. B. Zafar, I. Valera, M. G. Rogriguez, and K. P. Gummadi, "Fairness constraints: Mechanisms for fair classification," in *Proc. 20th Int. Conf. Artif. Intell. Statist.*, PMLR, 2017, pp. 962–970.
- [27] H. Suresh and J. Gutttag, "A framework for understanding sources of harm throughout the machine learning life cycle," in *Proc. Conf. Equity Access Algorithms Mechanisms Optim.*, 2021, pp. 1–9.
- [28] P. Grother, M. Ngan, and K. Hanaoka, "Face recognition vendor test part 3: Demographic effects," National Institute of Standards and Technology, Gaithersburg, MD, Tech. Rep. NIST IR 8280, Dec. 2019.
- [29] S. Dooley et al., "Comparing human and machine bias in face recognition," Oct. 2021, *arxiv.2110.08396*
- [30] O. Keyes, "The misgendering machines: Trans/HCI implications of automatic gender recognition," in *Proc. ACM Hum.-Comput. Interaction*, vol. 2, no. CSCW, pp. 1–22, Nov. 2018.

- [31] A. Birhane and V. U. Prabhu, "Large image datasets: A pyrrhic win for computer vision?," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, Waikoloa, HI, USA, 2021, pp. 1536–1546.
- [32] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Miami, FL, 2009, pp. 248–255.
- [33] E. Denton, A. Hanna, R. Amironesi, A. Smart, and H. Nicole, "On the genealogy of machine learning datasets: A critical history of ImageNet," *Big Data Soc.*, vol. 8, no. 2, Jul. 2021, Art. no. 205395172110359.
- [34] D. Zhao, A. Wang, and O. Russakovsky, "Understanding and evaluating racial biases in image captioning," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 14 810–14 820.
- [35] N. Garcia, Y. Hirota, Y. Wu, and Y. Nakashima, "Uncurated image-text datasets: Shedding light on demographic bias," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 6957–6966.
- [36] C. Hazirbas, J. Bitton, B. Dolhansky, J. Pan, A. Gordo, and C. C. Ferrer, "Towards measuring fairness in AI: The casual conversations dataset," *IEEE Trans. Biometrics, Behav., Identity Sci.*, vol. 4, no. 3, pp. 324–332, Jul. 2022.
- [37] B. Porgali, V. Albiero, J. Ryda, C. C. Ferrer, and C. Hazirbas, "The Casual Conversations V2 Dataset : A diverse, large benchmark for measuring fairness and robustness in audio/vision/speech models," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, 2023, pp. 10–17.
- [38] A. Wang et al., "REVISE: A tool for measuring and mitigating bias in visual datasets," *Int. J. Comput. Vis.*, vol. 130, no. 7, pp. 1790–1810, Jul. 2022.
- [39] J. Chen, Q. Ou, Z. Chi, and H. Fu, "Smile detection in the wild with deep convolutional neural networks," *Mach. Vis. Appl.*, vol. 28, no. 1, pp. 173–183, Feb. 2017.
- [40] P. Ekman and W. V. Friesen, "Constants across cultures in the face and emotion," *J. Pers. Social Psychol.*, vol. 17, no. 2, pp. 124–129, 1971.
- [41] P. Ekman, "What Scientists who study emotion agree about," *Perspectives Psychol. Sci.*, vol. 11, no. 1, pp. 31–34, Jan. 2016.
- [42] C. Chen, C. Crivelli, O. G. B. Garrod, P. G. Schyns, J.-M. Fernández-Dols, and R. E. Jack, "Distinct facial expressions represent pain and pleasure across cultures," in *Proc. Nat. Acad. Sci. USA*, vol. 115, no. 43, pp. E10 013–E10 021, Oct. 2018.
- [43] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
- [44] S. Li and W. Deng, "Deep facial expression recognition: A survey," *IEEE Trans. Affective Comput.*, vol. 13, no. 3, pp. 1195–1215, Third Quarter, 2022.
- [45] H. Guerdelli, C. Ferrari, W. Barhoumi, H. Ghazouani, and S. Berretti, "Macro- and micro-expressions facial datasets: A survey," *Sensors*, vol. 22, no. 4, Feb. 2022, Art. no. 1524.
- [46] M. Lyons, S. Akamatsu, M. Kamachi, and J. Gyoba, "Coding facial expressions with Gabor wavelets," in *Proc. IEEE 3rd Int. Conf. Autom. Face Gesture Recognit.*, 1998, pp. 200–205.
- [47] M. J. Lyons, "'Excavating AI' Re-excavated: Debunking a fallacious account of the JAFFE dataset," Jul. 2021, *arxiv.2107.13998*.
- [48] S. Singh and S. Benedict, "Indian semi-acted facial expression (iSAFE) dataset for human emotions recognition," in *Advances in Signal Processing and Intelligent Recognition Systems*, S. M. Thampi, R. M. Hegde, S. Krishnan, J. Mukhopadhyay, V. Chaudhary, and O. Marques, Eds., Singapore: Springer, 2020, pp. 150–162.
- [49] E. Kim, D. Bryant, D. Srikanth, and A. Howard, "Age bias in emotion detection: An analysis of facial emotion recognition performance on young, middle-aged, and older adults," in *Proc. AAAI/ACM Conf. AI Ethics Soc.*, 2021, pp. 638–644. [Online]. Available: <https://doi.org/10.1145/3461702.3462609>
- [50] K. Ahmad, S. Wang, C. Vogel, P. Jain, O. O'Neill, and B. H. Sufi, "Comparing the performance of facial emotion recognition systems on real-life videos: Gender, ethnicity and age," in *Proc. Future Technol. Conf.*, Springer International Publishing, 2022, pp. 193–210.
- [51] J. J. Greene et al., "The spectrum of facial palsy: The MEEI facial palsy photo and video standard set," *The Laryngoscope*, vol. 130, no. 1, pp. 32–37, 2020.
- [52] T. Xu, J. White, S. Kalkan, and H. Gunes, "Investigating bias and fairness in facial expression recognition," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 506–523.
- [53] A. Domnich and G. Anbarjafari, "Responsible AI: Gender bias assessment in emotion recognition," Mar. 2021, *arxiv.2103.11436*.
- [54] S. R. Jannat and S. Canavan, "Expression recognition across age," in *Proc. IEEE 16th Int. Conf. Autom. Face Gesture Recognit.*, 2021, pp. 1–5.
- [55] J. Deuschel, B. Finzel, and I. Rieger, "Uncovering the Bias in Facial Expressions," Nov. 2021, *arXiv: 2011.11311*.
- [56] R. Poyiadzi, J. Shen, S. Petridis, Y. Wang, and M. Pantic, "Domain generalisation for apparent emotional facial expression recognition across age-groups," Oct. 2021, *arxiv.2110.09168*.
- [57] J. Hernandez et al., "Guidelines for assessing and minimizing risks of emotion recognition applications," in *Proc. Int. Conf. Affect. Comput. Intell. Interaction*, 2021, Art. no. 8.
- [58] D. Pessach and E. Shmueli, "Algorithmic fairness," Jan. 2020, *arxiv.2001.09784*.
- [59] E. Pielou, "The measurement of diversity in different types of biological collections," *J. Theor. Biol.*, vol. 13, pp. 131–144, Dec. 1966.
- [60] W. H. Berger and F. L. Parker, "Diversity of planktonic foraminifera in deep-sea sediments," *Science*, vol. 168, no. 3937, pp. 1345–1347, Jun. 1970.
- [61] L. Jost, "Entropy and diversity," *Oikos*, vol. 113, no. 2, pp. 363–375, May 2006.
- [62] E. H. Simpson, "Measurement of diversity," *Nature*, vol. 163, no. 4148, pp. 688–688, Apr. 1949.
- [63] C. E. Shannon and W. Weaver, *The Mathematical Theory of Communication*, ser. The Mathematical Theory of Communication. Champaign, IL, USA: Univ. Illinois Press, 1949.
- [64] H. Cramér, "Chapter 21. The two-dimensional case," in *Mathematical Methods of Statistics*, ser. Princeton Mathematical Series. Princeton, NJ, USA: Princeton Univ. Press, 1991, no. 9, Art. no. 282.
- [65] A. A. Tschuprow, *Principles of the Mathematical Theory of Correlation*. New York, NY, USA: W. Hodge and Co., 1939.
- [66] J. M. Sakoda, "Measures of association for multivariate contingency tables," in *Proc. Social Statist. Sect. Amer. Statist. Assoc.*, 1977, pp. 777–780.
- [67] H. Theil, "On the estimation of relationships involving qualitative variables," *Amer. J. Sociol.*, vol. 76, no. 1, pp. 103–154, 1970. [Online]. Available: <http://www.jstor.org/stable/2775440>
- [68] M. Ducher, C. Cerutti, M. P. Gustin, and C. Z. Paultre, "Statistical relationships between systolic blood pressure and heart rate and their functional significance in conscious rats," *Med. Biol. Eng. Comput.*, vol. 32, no. 6, pp. 649–655, Nov. 1994.
- [69] R. H. Whittaker, "Vegetation of the Siskiyou mountains, Oregon and California," *Ecolog. Monographs*, vol. 30, no. 3, pp. 279–338, 1960.
- [70] A. Kammoun, R. Slama, H. Tabia, T. Ouni, and M. Abid, "Generative adversarial networks for face generation: A survey," *ACM Comput. Surv.*, vol. 55, pp. 1–37, Mar. 2022.
- [71] J. Cohen, *Statistical Power Analysis for the Behavioral Sciences*, 2nd ed. Evanston, IL, USA: Routledge, Jul. 1988.
- [72] G. Bouma, "Normalized (pointwise) mutual information in collocation extraction," in *Proc. Int. Conf. GSKL*, 2009, pp. 31–40.
- [73] D. Lundquist, A. Flykt, and A. Ohman, "Karolinska directed emotional faces," *PsychTESTS Dataset*, vol. 91, 1998, Art. no. 630.
- [74] T. Kanade, J. Cohn, and Y. Tian, "Comprehensive database for facial expression analysis," in *Proc. IEEE 4th Int. Conf. Autom. Face Gesture Recognit.*, 2000, pp. 46–53.
- [75] G. Zhao, X. Huang, M. Taini, S. Z. Li, and M. Pietikäinen, "Facial expression recognition from near-infrared videos," *Image Vis. Comput.*, vol. 29, no. 9, pp. 607–619, Aug. 2011.
- [76] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The extended Cohn-Kanade dataset (CK+): A complete dataset for action unit and emotion-specified expression," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. - Workshops*, 2010, pp. 94–101.
- [77] T. Bänziger, M. Mortillaro, and K. R. Scherer, "Introducing the Geneva multimodal expression corpus for experimental research on emotion perception," *Emotion*, vol. 12, no. 5, pp. 1161–1179, 2012.
- [78] N. Aifanti, C. Papachristou, and A. Delopoulos, "The MUG facial expression database," in *Proc. 11th Int. Workshop Image Anal. Multimedia Interactive Serv.*, 2010, pp. 1–4. [Online]. Available: <https://www.semanticscholar.org/paper/The-MUG-facial-expression-database-Aifanti-Papachristou/f1af714b92372c8e606485a3982eab2f16772ad8>
- [79] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon, "Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*, 2011, pp. 2106–2112.
- [80] I. J. Goodfellow et al., "Challenges in representation learning: A report on three machine learning contests," *Neural Netw.*, vol. 64, pp. 59–63, Apr. 2015.
- [81] M. Olszanowski, G. Pochwatko, K. Kuklinski, M. Scibor-Rylski, P. Lewinski, and R. K. Ohme, "Warsaw set of emotional facial expression pictures: A validation study of facial display photographs," *Front. Psychol.*, vol. 5, Jan. 2015, Art. no. 1516.

- [82] J. van der Schalk, S. T. Hawk, A. H. Fischer, and B. Doosje, "Moving faces, looking places: Validation of the Amsterdam Dynamic Facial Expression Set (ADFES)," *Emotion*, vol. 11, no. 4, pp. 907–920, 2011.
- [83] E. Barsoum, C. Zhang, C. C. Ferrer, and Z. Zhang, "Training deep networks for facial expression recognition with crowd-sourced label distribution," in *Proc. 18th ACM Int. Conf. Multimodal Interaction*, 2016, pp. 279–283.
- [84] A. Mollahosseini, B. Hasani, and M. H. Mahoor, "AffectNet: A database for facial expression, valence, and arousal computing in the wild," *IEEE Trans. Affective Comput.*, vol. 10, no. 1, pp. 18–31, First Quarter, 2019.
- [85] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, "From facial expression recognition to interpersonal relation prediction," *Int. J. Comput. Vis.*, vol. 126, no. 5, pp. 550–569, May 2018.
- [86] S. Li, W. Deng, and J. Du, "Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2584–2593.
- [87] S. Li and W. Deng, "Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition," *IEEE Trans. Image Process.*, vol. 28, no. 1, pp. 356–370, Jan. 2019.
- [88] J. Lee, S. Kim, S. Kim, J. Park, and K. Sohn, "Context-aware emotion recognition networks," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 10 142–10 151.
- [89] R. A. Khan, A. Crenn, A. Meyer, and S. Bouakaz, "A novel database of children's spontaneous facial expressions (LIRIS-CSE)," *Image Vis. Comput.*, vol. 83–84, pp. 61–69, Mar. 2019.
- [90] MMA Facial Expression, Jun. 2020. [Online]. Available: <https://www.kaggle.com/mahmoudima/mma-facial-expression>
- [91] Natural Human Face Images for Emotion Recognition, Dec. 2020. [Online]. Available: <https://www.kaggle.com/sudarshanvaidya/random-images-for-face-emotion-recognition>
- [92] K. Karkkainen and J. Joo, "FairFace: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2021, pp. 1547–1557.
- [93] D. E. King, "Max-margin object detection," Jan. 2015, *arxiv.1502.00046*.
- [94] I. Dominguez-Catena, D. Paternain, and M. Galar, "Gender stereotyping impact in facial expression recognition," in *Machine Learning and Principles and Practice of Knowledge Discovery in Databases*, vol. 1752. Berlin, Germany: Springer, 2023, pp. 9–22.
- [95] The GenIUSS Group, "Best practices for asking questions to identify transgender and other gender minority respondents on population-based surveys," The Williams Institute, Los Angeles, CA, Tech. Rep., 2014. [Online]. Available: <https://williamsinstitute.law.ucla.edu/publications/geniuss-trans-pop-based-survey/>
- [96] S. Fabbri, S. Papadopoulos, E. Ntoutsis, and I. Kompatsiaris, "A survey on bias in visual datasets," *Comput. Vis. Image Understanding*, vol. 223, Oct. 2022, Art. no. 103552.



Iris Dominguez-Catena (Student Member, IEEE) received the BSc and MSc degrees in computer science from the Public University of Navarra, in 2015 and 2020, respectively. She worked for several private software companies from 2015 to 2018. She is currently working toward the PhD degree with the Public University of Navarra, working on demographic bias issues in Artificial Intelligence. Her research interests focus on AI fairness, bias detection and mitigation, and other ethical problems of AI deployment in society.



Daniel Paternain (Member, IEEE) received the MSc and PhD degrees in computer science from the Public University of Navarra, Pamplona, Spain, in 2008 and 2013, respectively. He is currently associate professor with the Department of Statistics, Computer Science and Mathematics. He is also the author or coauthor of almost 40 articles in journals from JCR and more than 50 international conference communications. His research interests include both theoretical and applied aspects of information fusion, computer vision, and machine learning.



Mikel Galar (Member, IEEE) received the MSc and PhD degrees in computer science from the Public University of Navarra, Pamplona, Spain, in 2009 and 2012, respectively. He is currently an associate professor with the Public University of Navarra. He is the author of 50 published original articles in international journals and more than 80 contributions to conferences. He is a coauthor of a book on imbalanced datasets and a book on large-scale data analytics. His research interests are machine learning, deep learning, ensemble learning, and Big Data. He received the extraordinary prize for his PhD thesis from the Public University of Navarra and the 2013 IEEE Transactions on Fuzzy System Outstanding Paper Award (bestowed in 2016).