



Repositorio Institucional de la Universidad Autónoma de Madrid
<https://repositorio.uam.es>

Esta es la **versión de autor** del artículo publicado en:
This is an **author produced version** of a paper published in:

IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC),
Torino, Italy, 2023

DOI: <https://doi.org/10.1109/COMPSAC57700.2023.00176>

Copyright: © 2023 IEEE

El acceso a la versión del editor puede requerir la suscripción del recurso
Access to the published version may require subscription

Measuring Bias in AI Models: An Statistical Approach Introducing N-Sigma

Daniel DeAlcala, Ignacio Serna, Aythami Morales, Julian Fierrez, Javier Ortega-Garcia
Biometrics and Data Pattern Analytics Lab, Universidad Autonoma de Madrid, Spain

Abstract—The new regulatory framework proposal on Artificial Intelligence (AI) published by the European Commission establishes a new risk-based legal approach. The proposal highlights the need to develop adequate risk assessments for the different uses of AI. This risk assessment should address, among others, the detection and mitigation of bias in AI. In this work we analyze statistical approaches to measure biases in automatic decision-making systems. We focus our experiments in face recognition technologies. We propose a novel way to measure the biases in machine learning models using a statistical approach based on the N-Sigma method. N-Sigma is a popular statistical approach used to validate hypotheses in general science such as physics and social areas and its application to machine learning is yet unexplored. In this work we study how to apply this methodology to develop new risk assessment frameworks based on bias analysis and we discuss the main advantages and drawbacks with respect to other popular statistical tests.

Index Terms—Artificial Intelligence, AI, Bias, Explainable, Risk Assessment, Trustworthiness, 5-Sigma

I. INTRODUCTION

Artificial Intelligence (AI) can play an important role to achieve the Sustainable Development Goals (SDGs) by 2030 [1]. AI brings enormous benefits in several critical areas for our society (e.g., health, security, sustainability), but it can also significantly compromise the safety of citizens worldwide. The development of a Responsible AI technology needs an international multidisciplinary effort to ensure the trustworthiness, sustainability, and safety. This effort involves a multi-stakeholder work including academia, industry, civil society, and public agencies, among others.

The absence of international standards for the development of Responsible Artificial Intelligence has motivated a wide variety of approaches [2], [3]. The regulation is moving from a technology-based framework to a risk-based framework [4]. The new regulatory framework proposed by the European Union defines 4 levels of risk in AI: *i*) Unacceptable, *ii*) High, *iii*) Limited, and *iv*) Minimal. As an example, high-risk technologies will be subject to strict obligations including adequate risk assessment and mitigation systems. This risk-based framework requires protocols and technologies capable of assessing and explaining the results of AI systems based on parameters beyond the traditional performance metrics (e.g., overall accuracy).

How to measure or assess the fairness of an automatic-decision algorithm is not a trivial task. Fairness is a human concept that can be mathematically defined in different ways [5]. Traditionally, fairness is measured as a difference in performance between population subgroups (e.g., performance

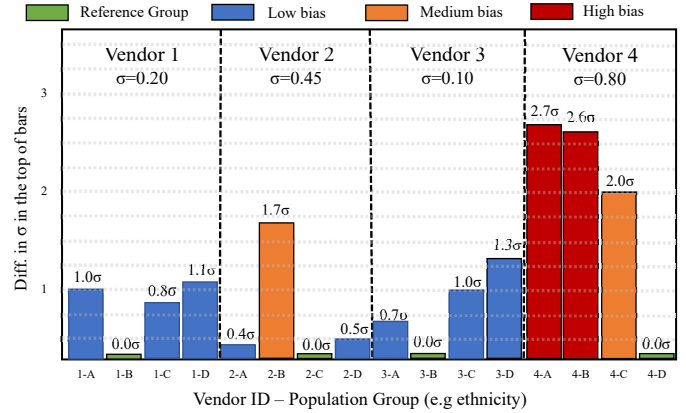


Fig. 1: Performance differences of 4 AI models (e.g., Face Recognition technology) evaluated over 4 different demographic groups (A,B,C,D). The difference is measured in number of sigmas (σ) with respect to a reference group. The bias level is represented with different colors.

for different demographic groups). The literature has proposed several approaches to measure such a difference with traditional statistical methods [6], [7] or machine learning approaches [5], [8], but none of these methods has been yet adopted as a widely recognized standard.

On the other hand, the 5-sigma approach is widely used for statistical analyses in many fields including natural [9] and human sciences [10]. In this work we extend this 5-Sigma approach and apply it to bias analysis of data-driven learning models (see Fig. 1). The contributions of this work are:

- The proposal of a common experimental protocol to achieve a fairer and more standardized evaluation of AI models.
- We analyze two pointwise metrics and a traditional distribution metric for bias analysis in machine learning models. (More specifically, in discrimination-aware face recognition models.)
- We extend 5-Sigma into *N*-Sigma for bias assessment of machine learning. This proposed extension is compatible with a risk-based evaluation framework where a variable (*N*) can be associated to each risk level (see Fig. 1).

II. RELATED WORKS

Bias in ML systems is an increasingly studied topic for which various notions of fairness have been applied [8],

[11]. The most common way to measure bias is through performance in demographic groups, but it is not the only way. Researchers have also looked at how models respond at the level of activation and how this is different across different groups [12], [13].

Among other AI application fields, face biometrics is perhaps the most popular and evolved one regarding bias analysis [14]. From the fairness criteria proposed in the literature, the statistical parity criterion is inadequate for Face Recognition (FR) models. The work [15] shows that a perfect model does not imply demographic parity with entangled variables, which is the case of FR, where sensitive demographic characteristics are linked to identity. It is unreasonable to think that the outcome of a face recognition system is independent of a person's ethnicity when a white user is trying to impersonate a black user. Equalized odds are often used; for example, the NIST report uses false negatives and false positives for each demographic group to measure the fairness [16].

Except a few exceptions [17], most of the literature studying bias in facial algorithms does not clearly define what bias is and merely shows that the performance varies between population groups. Recent research is attempting to mitigate biases after quantifying them. These recent papers typically use a form of standard deviation of the algorithm performance across individuals of different populations as a measure for bias, both implicitly and explicitly [17]–[20].

A. Bias in face recognition

The number of academic studies analyzing the fairness of face recognition algorithms has grown significantly in recent years, and the number of published works pointing out the biases in the results of face detection [21] and recognition algorithms is large [14], [22]–[24]. Facial recognition systems can suffer from a variety of biases, ranging from those arising from unconstrained environmental variables such as illumination, pose, expression, and face resolution, from systematic errors such as image quality [25], [26], and from demographic factors [27] like age, sex, and race. Among these different covariates, the skin color is repetitively remarked as a factor with high impact in the performance [28].

III. DATABASE AND MODELS

In our experiments we used Racial Faces in the Wild (RFW) [29]. This database is divided into four demographic classes: Caucasian, Indian, Asian, and African. Each class has about 10K images of 3K individuals. There are no major differences in pose, age, and sex distribution between Caucasian, Asian, and Indian groups. The African group has a smaller age difference than the others, and while females account for approximately 35% in the other groups, they account for less than 10% in the African group.

The model used is a ResNet-100 network [30], trained on the MS1Mv3¹ database [31] (93K identities and 5.2M

images) with ArcFace [32] loss function. A model with 101 convolutional layers and 44 million parameters.

When using facial recognition systems in verification mode, two faces are assigned the same identity if their similarity distance is smaller than a threshold τ . The similarity is computed between the two face descriptors \mathbf{x}_r and \mathbf{x}_s obtained from a face model. A similarity score is known as a genuine score or authentic score if it is the result of matching two samples of the same biometric trait of a user. It is known as an impostor score if it involves comparing two biometric samples originating from different users [33]. Several metrics can be used to compute similarity, the two most frequent are euclidean distance and cosine similarity.

From the similarity results, EER (Equal Error Rate) and TPR (True Positive Rate) are computed for a specific threshold. The decision threshold τ for each model is different and is set using genuine and impostor comparisons. The EER is the error at a given threshold at which FMR (False Match Rate) and FNMR (False Non-match Rate) are equal. The TPR is the probability of correctly identifying two user samples as being from the same user.

IV. METHODS: EXPERIMENTAL PROTOCOL

A. Training protocol

To have a reference of the performance of the metrics, different biased and unbiased Face Recognition models have been trained using the base model explained in the Section III.

Basically, a finetuning was carried out: a dense layer was added at the end of the model and trained with different data depending on the aimed bias. For example, to positively bias Asian ethnicity, the dense layer is trained only with faces of Asian people. If non-bias is intended, this dense layer is trained with data from all ethnicities. Triplet Loss function [34], [35] is used, whose objective is to bring the feature embeddings of the same user closer together and to pull apart those of different users, in the feature space.

For this training, 50% of the RFW database users were used. To avoid that the results depend on the training of a single model, M models are trained ($M = 20$ in this work) for each aimed bias. In this way, we obtain average results not affected by the stochastics associated with the training process. The M models are trained by bootstrapping the 75% of the users within the 50% belonging to the training set. Bootstrapping [36] is a method of inferring results of a population from the results found in a collection of smaller random samples from that population, using replacement during the sampling process.

B. Evaluation protocol

The other 50% of the RFW users serve for the evaluation stage. In this stage, the values of the 4 metrics explained in Section V (SP, EOP, T-Test, N-Sigma) were obtained for each model. These metric values were obtained by ethnic group to compare performance differences between ethnic groups and to establish the existence or not of bias.

¹https://github.com/deepinsight/insightface/tree/master/recognition/_datasets_ms1m-retinaface

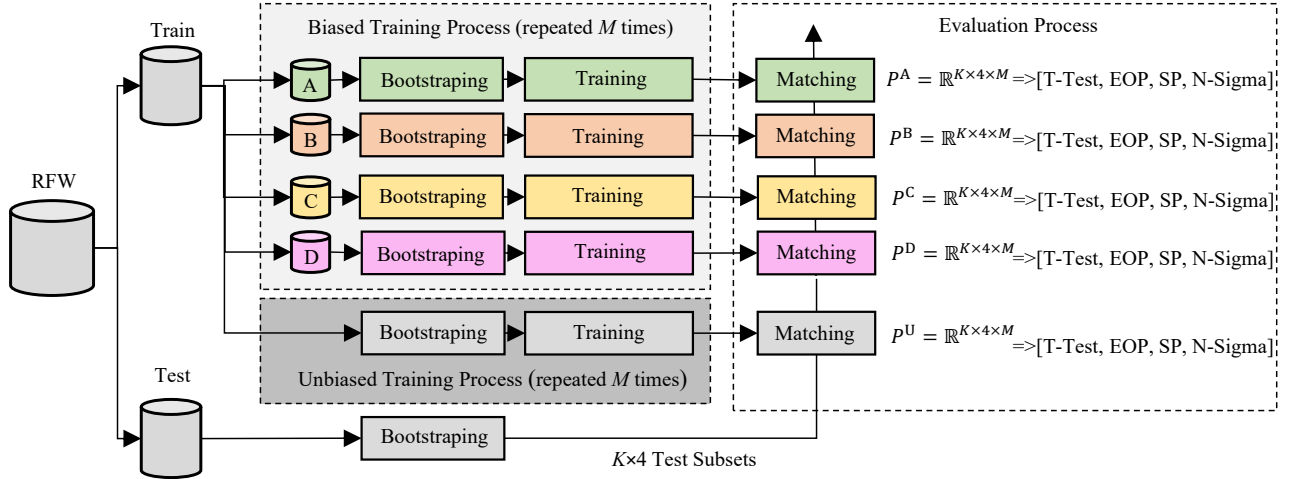


Fig. 2: Experimental Framework for the analysis of biased learning processes. A, B, C, and D represent different demographic groups used to introduce bias in the learning process.

For clarification, a model was trained to be positively biased for one ethnicity and then tested on all ethnicities independently in order to study the differences in performance across them.

Measuring bias in a model is not trivial, it can yield results that are difficult to interpret and compare. When a performance value (for example, EER or TPR) is obtained from a neural network, it has a margin of error, i.e., that value may be a little higher or a little lower depending on the specific data used. To try to avoid this variability and not lose valuable information, in this work the metrics for measuring bias, are obtained from a group of K EER/TPR values (each value obtained from a subset). In other words, several EER/TPR values are calculated to reflect the variability of the performance results in a model and in 1) **Pointwise metrics (SP, EOP)**: performance values are compared one by one, or in 2) **distribution metrics (T-Test, N-Sigma)**: the group of performance values is compared as a distribution.

C. Optimizing subsets

For the computation of the 4 metrics analyzed in this work, K subsets are created within each ethnic group, and each of these subsets is used to obtain one EER/TPR value. The K subsets have to be representative of the database to obtain a valid set of performance values. Bootstrapping [36] as explained before is a statistical technique that enables us to estimate the characteristics of a population by taking multiple random samples from it. The method involves creating smaller subsets from the larger population, with replacement, and using them to calculate the desired statistics:

- 1) A small number of samples in a subset will give values dependent on the samples selected.
- 2) A large number of samples may be unnecessary and may complicate the computation.

D. Experimental framework

The entire workflow is presented in Figure 2:

- 1) The RFW database is divided in training (50%) and evaluation (50%).
- 2) The training stage is used to create the Biased Models as explained in the Subsection IV-A. As a result, we have M biased models for each of the 4 ethnic groups and M unbiased models, in total: $M \times 4 + M$ models.
- 3) At the evaluation stage, K subsets of each ethnic group (A,B,C,D) are selected from the database, that is, $K \times 4$. Consequently, the models are evaluated with $K \times 4$ randomly selected subsets. $K \times 4$ new subsets are sampled for each batch of models. As M models of each type are trained (M batches of models), $(K \times 4) \times (M)$ subsets are created.
- 4) Each evaluation metric (SP, EOP, T-Test, N-Sigma) produces a single value for each ethnicity. Thus, initially, we obtain $(\text{number of ethnicities} = 4) \times M$ values on each evaluation metric for each of the model types (ethnically biased A,B,C,D and unbiased U).
- 5) M models per type were trained to avoid stochastic effects associated with training in the results. Therefore, the results presented in Section (VI-A) show the average over the M model batches. For this reason, the final results show $\text{number of ethnicities} \times 1$ values on each evaluation metric for each type of model (A,B,C,D,U).

V. MEASURING BIAS IN AI APPLICATIONS

Bias refers to the unequal behavior of an algorithm; this irregular behavior may render its decisions unfair and is therefore called biased. Thus, in AI terms, bias is measured in terms of differences in performance between different groups.

In this work we are going to focus on measuring bias using the experimental protocol previously described applied on two pointwise metrics and two distribution metrics.

A. Pointwise metrics

The experimental protocol is applied to two metrics used in the literature and particularized to our use case. Therefore,

the definitions of the metrics undergo certain changes which are shown below.

Consider a binary classifier \hat{Y} . An outcome $\hat{Y} = 0$ represents a “non-match” decision (i.e., comparison between samples of different classes), while $\hat{Y} = 1$ represents a “match” decision (i.e., comparison between samples from the same class). The literature includes specific measures proposed to detect biased results in machine learning models [37]. In this work we will use two:

- Statistical Parity (SP) or Demographic Parity [37]: $P(\hat{Y}|s = 0) = P(\hat{Y}|s = 1)$ which means that the predictions must be independent of attribute s , and the probability of obtaining an outcome must always be the same regardless of the attribute (e.g. gender, ethnicity, age). This metric is not suitable for FR as we already explained in Section 2, however in this case we are going to use a particularity of this definition which can be adequate to measure the bias in FR systems. Statistical parity between groups can be expressed in terms of both False Match Rate (FMR) and False Non-Match Rate (FNMR) at a certain decision threshold [38]:

$$SP(\tau) = 1 - (\alpha A(\tau) + (1 - \alpha)B(\tau)) \quad (1)$$

where α defines the weight of the importance of False Matches, $A(\tau)$ is an specific operational point (defined by the threshold τ) of the FMR differential across groups for a given threshold and $B(\tau)$ is the FNMR for this operational point. If taken the threshold at which FMR and FNMR are equal, then the FMR and FNMR became the EER (Equal Error Rate) and $A = B$. The equation simplifies to:

$$SP(\tau) = 1 - A(\tau_{EER}) = 1 - B(\tau_{EER}) \quad (2)$$

Defining $A()$ as the mean of the EER differences between two groups ($G1$ and $G2$), where each EER is obtained from each of the K subsets of the group, we get:

$$SP(\tau_{EER}) = 1 - \frac{1}{K} \sum_{i=1}^K |EER_i^{G1} - EER_i^{G2}| \quad (3)$$

- Equality of Opportunity (EOP) [5]: $P(\hat{Y} = 1|s = 0, Y = 1) = P(\hat{Y} = 1|s = 1, Y = 1)$. Used in biometric literature as *differential* value, this metric is a relaxed version of the equalized odds criterion. EOP considers only the True Positive Rates (TPR). This definition of EOP serves to indicate that the TPR between different groups must be equal. In this case, the formula depends on the operational point (threshold τ):

$$EOP(\tau) = 1 - \frac{1}{K} \sum_{i=1}^K |TPR_i^{G1}(\tau) - TPR_i^{G2}(\tau)| \quad (4)$$

The operational point chosen in this work is the one corresponding to an FPR of 0.01.

B. Distribution metrics

Instead of comparing the performance values (TPR/EER) independently as it is done with the previous metrics (Equations 3 and 4), in the distribution metrics the performance values are understood as a group, and the aim is to compare them as a distribution.

1) *Traditional statistical test (T-Test)*: In our case, an appropriate statistical approach to compare the distributions is the T-Test. This test is used to evaluate the statistical significance of the difference between the means μ_{G1} and μ_{G2} of two populations, in situations where the populations follow a normal distribution, the standard deviation is unknown, and the sample size is small. It uses an estimation of the standard deviation instead of the true value. The selected statistic in our case is the *Welch corrected unpaired T-Test*:

$$Z = \frac{\mu_{G1} - \mu_{G2}}{s} \quad \text{where} \quad s = \sqrt{\frac{s_{G1}^2 + s_{G2}^2}{n}} \quad (5)$$

where $n = n_{G1} = n_{G2}$ is the number of samples and s_{G1}^2 and s_{G2}^2 are the unbiased estimators of the population variance. The null hypothesis $H_0 : G1 = G2$ (the distribution of FR results for both groups is the same) is rejected if $|Z| > t_{1-\alpha/2}$, where t_γ is the γ -quantile value of the t distribution.

2) *The N-Sigma method*: The 5-Sigma method in particle physics refers to the probability in a mass spectrum of having a statistical fluctuation (a peak) in the background. The probability (p -value) of a chance peak must not exceed 5σ of a normalised Gaussian distribution. The sigma (σ) is the deviation from the mean (μ) of the distribution that includes approximately 68% (34% on each side of the mean) of the data. If we select two sigmas from the mean, we would have around the 95% of the data. If we select 5 sigma, the samples not included are only about a $3 \times 10^{-7}\%$.

When searching for discovery, the data statistic used to discriminate between background only (known as the null hypothesis H_0) and “background plus signal” (H_1) is usually the L_1/L_0 likelihood ratio for the two hypotheses; and the 5σ criterion is applied to the observed value of this ratio, as compared with its expected distribution assuming just background [9]. In this work the N-Sigma method can be expressed as:

$$N = \frac{\mu_{G1} - \mu_{G2}}{\sigma_{G1}} \quad (6)$$

where μ_{G1} and μ_{G2} are the means of the two populations being compared. σ_{G1} is the standard deviation of the population used as reference.

Here, unlike in the T-Test, we do not reject or accept the hypothesis by setting a threshold. In this case the result yields a distance N between the two distributions. This distance can be used to define risk levels.

VI. EXPERIMENTS

A. Results

First of all, in Table I we have the mean EER/TPR value (%) for each ethnicity group on the models created. Also,

Eth. Finetuned	Ethnicity Evaluated			
	African	Asian	Caucasian	Indian
All	1.83/78.8	1.58/78.1	1.69/82.1	2.07/79.1
African	1.98/78.3	1.88/74.4	2.04/78.4	2.31/74.3
Asian	2.26/73.4	1.72/76.6	1.95/76.7	2.51/72.3
Caucas.	2.15/74.6	1.82/73.7	1.95/79.3	2.33/73.3
Indian	2.19/73.6	1.88/72.7	2.07/77.1	2.22/76.7

TABLE I: EER/TPR mean values (%) for each ethnicity after unbiased (none, first row) and biased training (performing a fine tuning for an specific ethnicity, following 4 rows).

the values of the evaluation metrics are shown in Table II. It is important to understand that the metrics in our work are applied with an evaluating group with respect to a reference group, since bias is a human concept which must be measured with respect to something. The reference group (named as G_1 in Equation 3, 4, 5 and 6) has been chosen as the one with the lowest mean in Table I and a similarity value is given for all evaluation groups (G_2 in Equation 3, 4, 5 and 6) with respect to it. The reference group could very well be another one (e.g. the one with the highest mean). In the case of the EOP, SP, and the T-Test, a higher value implies more similarity, although understanding the meaning of the value is not trivial. In the case of the N-Sigma method, a lower value implies more similarity, the specific value being the distance between the distributions in sigmas. In the case that the reference group (G_1) is the same as the group being evaluated (G_2) the SP, EOP and T-Test metrics will give the maximum value which is 1, while the N-Sigma metric the minimum which is 0. This is because the groups are identical.

All the values present in Tables I and II are an average of the values obtained for the $M = 20$ trained models.

1) *Discussion:* The best results in terms of mean are achieved with the model trained with all the ethnicities (Table I) because they self-regulate each other. In this case, it can be seen that for the EER, the best results are achieved with the Asian ethnicity, while for the TPR the best results are obtained with the Caucasian ethnicity.

Having seen this, analyzing the tables it is observed that when the model is biased for a particular ethnicity, the similarity value increases with respect to that ethnicity for all the metrics (these results can be seen in subtables IIb, IIc, IId, IIe when compared to subtable IIa). The Caucasian ethnicity is the exception: training only with the Caucasian samples has not decreased the bias with respect to this ethnicity. We do not have a groundtruth that tells us what should come out and therefore we cannot say in terms of values which metric works best, but must speak in other terms such as interpretability.

Regarding the interpretability of the different metrics, different aspects can be analyzed.

- The EOP and the SP simply measure differences at the subset level and gives you an average value. So what you can see here is a result based simply on a mean difference in performance between groups, and the results should be understood as such.

TABLE II: The next five subtables (a-e) present performance metrics Mean for the models favored in different ethnicities.

(a) Finetuned for All Ethnicities (U). Sigma value: 0.196

Eth Eval	T-Test	EOP	SP	N-sig
African	1.00×10^{-09}	0.967	0.997	1.30
Asian	1.00	0.961	1.00	0.00
Caucasian	0.16	1.00	0.998	0.58
Indian	4.94×10^{-23}	0.970	0.995	2.53

(b) Finetuned for African (A). Sigma value: 0.206

Eth Eval	T-Test	EOP	SP	N-sig
African	0.12	0.998	0.998	0.52
Asian	1.00	0.959	1.00	0.00
Caucasian	1.00×10^{-03}	1.00	0.998	0.81
Indian	2.81×10^{-16}	0.959	0.995	2.10

(c) Finetuned for Asian (B). Sigma value: 0.205

Eth Eval	T-Test	EOP	SP	N-sig
African	8.96×10^{-31}	0.966	0.994	2.64
Asian	1.00	0.999	1.00	0.00
Caucasian	1.62×10^{-07}	1.00	0.997	1.16
Indian	5.41×10^{-51}	0.955	0.992	3.89

(d) Finetuned for Caucasian (C). Sigma value: 0.210

Eth Eval	T-Test	EOP	SP	N-sig
African	1.44×10^{-14}	0.952	0.996	1.58
Asian	1.00	0.943	1.00	0.00
Caucasian	6.0×10^{-3}	1.00	0.998	0.59
Indian	1.70×10^{-24}	0.938	0.994	2.32

(e) Finetuned for Indian (D). Sigma value: 0.207

Eth Eval	T-Test	EOP	SP	N-sig
African	5.16×10^{-18}	0.964	0.996	1.56
Asian	1.00	0.954	1.00	0.00
Caucasian	2.0×10^{-4}	1.00	0.998	0.94
Indian	3.60×10^{-27}	0.995	0.996	1.62

- The interpretability of the T-Test and the N-Sigma is a bit different. In this case, you take the values by subsets and make a comparison of the distributions between the ethnic groups. Therefore, the value tells how far apart the distributions of values are. Both metrics give a value of the difference between the distributions, however, the results of the T-Test are somewhat more difficult to interpret. In this case, a significance level that can be $\alpha = 0.05$ is chosen, and whenever the T-Test result is below that level, it is said that the distributions are statistically different. Subsequently, the lower the value, the more different they are, but the values themselves have no meaning. As for the N-Sigma method, the results represent the same thing but it is easier to interpret what the value itself means. What it means is, with respect to the variance of the favored distribution, how far away the rest of the distributions are, i.e., a value of 1 means that it is at 1 variance and a value of 2 would imply that it is twice as far away.

Dissecting the results we realize that with the distribution

comparison methods the difference between the ethnic groups can be appreciated more clearly and allows deeper analysis. And between these two methods, the N-Sigma method offers more easily understandable results.

As we have just explained, N-Sigma and T-Test values do not represent the same thing, one is the result of a statistical test while the other is a distance between distributions. Therefore, although the T-Test and N-Sigma results show some correlation (i.e. lower T-Test values are related to higher N-Sigma values), it is not a perfect inverse correlation (otherwise both values would mean the same thing without contributing anything new). For this reason, for a nearly same value of N-Sigma (0.58 in Table IIa Eth eval = Caucasian and 0.59 in Table IIb Eth eval = Caucasian) the T-Test value is different.

VII. CONCLUSION

In this paper we propose the use of metrics to measure bias under an experimental protocol and specifically, a metric called N-Sigma widely used in other fields but unexplored in AI. This metric is based on the well-established idea of 5-sigmas used in fields such as physics or economics. We have evaluated a model fine-tuned to be biased for different ethnicities with the distribution and pointwise metrics.

The results show that the distribution methods yield results more interpretable. Among the distribution comparison methods the N-Sigma method results are more user-friendly. The use of this metric makes it possible to adapt very easily to different use cases by varying the sigma at which bias is considered to occur, e.g., defining different risk levels associated with different values of N. In applications where the presence of bias is critical (high risk), a lower sigma value can be assigned, while if the application is more flexible in this regard, the allowed sigma value can be increased.

ACKNOWLEDGMENT

Support by project BBforTAI (PID2021-127641OB-I00 MICINN/FEDER). D. deAlcala is supported by a FPU Fellowship (FPU21/05785) from the Spanish MIU.

REFERENCES

- [1] UN, "United nations activities on artificial intelligence (ai)," in *ITU Publications*, International Telecommunication Union, 2021. 1
- [2] B. Lepri, N. Oliver, and A. Pentland, "Ethical machines: The human-centric use of artificial intelligence," *IScience*, vol. 24, no. 3, 2021. 1
- [3] I. Rahwan *et al.*, "Machine behaviour," *Nature*, vol. 568, no. 7753, pp. 477–486, 2019. 1
- [4] E. Commission, "Proposal for a Regulation Laying Down Harmonised Rules on Artificial Intelligence," in *Regulation of the European Parliament and of the Council*, 2021. 1
- [5] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," *NIPS*, vol. 29, 2016. 1, 4
- [6] A. Godbole, S. A. Grosz, K. Nandakumar, and A. K. Jain, "On demographic bias in fingerprint recognition," *arXiv preprint arXiv:2205.09318*, 2022. 1
- [7] I. Žliobaitė, "Measuring discrimination in algorithmic decision making," *Data Mining and Knowledge Discovery*, vol. 31, no. 4, 2017. 1
- [8] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," *ACM Computing Surveys (CSUR)*, vol. 54, no. 6, pp. 1–35, 2021. 1
- [9] L. Lyons, "Discovering the Significance of 5 Sigma," *arXiv preprint arXiv:1310.1284*, 2013. 1, 4
- [10] K. Leonard, "Schaum's outline of business statistics fourth edition," 2003. 1
- [11] S. Verma and J. Rubin, "Fairness Definitions Explained," in *IEEE/ACM International Workshop on Software Fairness (FairWare)*, 2018. 1
- [12] S. Naggal, M. Singh, R. Singh, M. Vatsa, *et al.*, "Deep Learning for Face Recognition: Pride or Prejudiced?," *arXiv:1904.01219*, 2019. 2
- [13] I. Serna *et al.*, "InsideBias: Measuring Bias in Deep Networks and Application to Face Gender Biometrics," in *ICPR*, 2021. 2
- [14] P. Terhorst *et al.*, "A comprehensive study on face recognition biases beyond demographics," *IEEE Trans. on Technology and Society*, vol. 3, pp. 16–30, March 2022. 2
- [15] F. Locatello *et al.*, "On the Fairness of Disentangled Representations," *Advances in Neural Information Processing Systems*, vol. 32, 2019. 2
- [16] P. J. Grother, M. L. Ngan, and K. K. Hanaoka, *Ongoing Face Recognition Vendor Test (FRVT) Part 3: Demographic Effects*. NIST Internal Report, National Institute of Standards and Technology, 2019. 2
- [17] I. Serna, A. Morales, J. Fierrez, *et al.*, "SensitiveLoss: Improving accuracy and fairness of face representations with discrimination-aware deep learning," *Artificial Intelligence*, vol. 305, 2022. 2
- [18] S. Gong, X. Liu, and A. Jain, "Mitigating Face Recognition Bias via Group Adaptive Classifier," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, (Nashville, TN), IEEE, June 2021. 2
- [19] M. Wang and W. Deng, "Mitigate Bias in Face Recognition using Skewness-Aware Reinforcement Learning," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2020. 2
- [20] Z. Yang *et al.*, "RamFace: Race Adaptive Margin Based Face Recognition for Racial Bias Mitigation," in *International Joint Conference on Biometrics (IJCB)*, IEEE, 2021. 2
- [21] J. Buolamwini and T. Gebru, "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification," in *Conference on Fairness, Accountability and Transparency*, 2018. 2
- [22] B. Klare *et al.*, "Face Recognition Performance: Role of Demographic Information," *IEEE TIFS*, vol. 7, no. 6, pp. 1789–1801, 2012. 2
- [23] I. Hupont and C. Fernandez, "DemogPairs: Quantifying the Impact of Demographic Imbalance in Deep Face Recognition," in *FG*, 2019. 2
- [24] P. Drozdzowski *et al.*, "Demographic Bias in Biometrics: A Survey on an Emerging Challenge," *IEEE Transactions on Technology and Society*, vol. 1, no. 2, pp. 89–103, 2020. 2
- [25] T. Schlett, C. Rathgeb, O. Henniger, J. Galbally, J. Fierrez, and C. Busch, "Face image quality assessment: A literature survey," *ACM Computing Surveys*, vol. 54, pp. 1–49, January 2022. 2
- [26] J. Hernandez-Ortega *et al.*, "FaceVec: Vector quality assessment for face biometrics based on ISO compliance," in *IEEE/CVF Winter Conf. on Applications of Computer Vision Workshops*, January 2022. 2
- [27] E. Gonzalez-Sosa *et al.*, "Facial soft biometrics for recognition in the wild: Recent works, annotation and COTS evaluation," *IEEE Trans. on Inf. Forensics and Security*, vol. 13, no. 8, pp. 2001–2014, 2018. 2
- [28] B. Lu, J.-C. Chen, C. D. Castillo, and R. Chellappa, "An Experimental Evaluation of Covariates Effects on Unconstrained Face Verification," *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 1, no. 1, pp. 42–55, 2019. 2
- [29] M. Wang *et al.*, "Racial Faces in the Wild: Reducing Racial Bias by Information Maximization Adaptation Network," in *ICCV*, 2019. 2
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *CVPR*, IEEE, 2016. 2
- [31] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, "Ms-celeb-1m: A dataset and benchmark for large-scale face recognition," in *European Conference on Computer Vision (ECCV)*, pp. 87–102, Springer, 2016. 2
- [32] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive Angular Margin Loss for Deep Face Recognition," in *CVPR*, 2019. 2
- [33] A. Ross and A. K. Jain, *Biometrics, Overview*, pp. 168–172. Boston, MA: Springer US, 2009. 2
- [34] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," in *CVPR*, 2015. 2
- [35] A. Morales *et al.*, "SetMargin loss applied to deep keystroke biometrics with circle packing interpretation," *Pattern Recognition*, vol. 122, p. 108283, February 2022. 2
- [36] B. Efron and R. LePage, *Introduction to Bootstrap*. Wiley & Sons, New York, 1992. 2, 3
- [37] S. Barocas, M. Hardt, and A. Narayanan, "Fairness in machine learning," *NIPS tutorial*, vol. 1, p. 2017, 2017. 4
- [38] T. Freitas Pereira and S. Marcel, "Fairness in biometrics: A figure of merit to assess biometric verification systems," *IEEE Trans. on Biometrics, Behavior, and Id. Science*, vol. 4, no. 1, pp. 19–29, 2021. 4