MEASURE TWICE, CUT ONCE: QUANTIFYING BIAS AND FAIRNESS IN DEEP NEURAL NETWORKS

Anonymous authors Paper under double-blind review

ABSTRACT

Algorithmic bias is of increasing concern, both to the research community, and society at large. Bias in AI is more abstract and unintuitive than traditional forms of discrimination and can be more difficult to detect and mitigate. A clear gap exists in the current literature on evaluating the relative bias in the performance of multi-class classifiers. In this work, we propose two simple yet effective metrics, Combined Error Variance (CEV) and Symmetric Distance Error (SDE), to quantitatively evaluate the class-wise bias of two models in comparison to one another. By evaluating the performance of these new metrics and by demonstrating their practical application, we show that they can be used to measure fairness as well as bias. These demonstrations show that our metrics can address specific needs for measuring bias in multi-class classification.

1 INTRODUCTION

Broad acceptance of the large-scale deployment of AI and neural networks depends on the models' perceived trustworthiness and fairness. However, research on evaluating and mitigating bias for neural networks in general and compressed neural networks in particular is still in its infancy. Because deep neural networks (DNNs) are "black box" learners, it can be difficult to understand what correlations they have learned from their training data, and how that affects the downstream decisions that are made in the real world. Two models may appear to have very similar performance when only measured in terms of accuracy, precision, etc. but deeper analysis can show uneven performance across many classes. Moreover, when the number of tasks grows beyond one or two, the difficulty in reasoning and quantifying trade-offs when selecting or validating a model also grows.



Figure 1: Google NGram Data Michel et al. (2011) showing relative usage of machine learning related terms over time. Deep learning has quickly passed up the use of statistical terms like logistic regression. Computer Vision tasks like Object Detection and Image Recognition are growing at faster rates than Binary Classification which fairness metrics can address.

Widely accepted and effective metrics for measuring the bias of several neural networks against one another are still missing. Issues of both fairness and bias, which will be discussed as distinct but related phenomena in this paper, can seriously degrade the trustworthiness of a machine learning model in real-world conditions. It is important to quantify the performance of models in terms both of bias and fairness. While there exists extensive work on AI fairness regarding binary classification tasks Borkan et al. (2019); Hinnefeld et al. (2018); Dixon et al. (2018); Maughan & Near

(2020), there is a shortage of metrics extending these ideas to other machine learning tasks. Many researchers have taken an interest in wanting to ensure their models are fair, currently, we simply do not have the tools to measure for many domains.

In fact, recent trends shown in Figure 1, are exacerbating the divide with the majority of new research in neural networks exploring topics outside of the reach of existing fairness and bias metrics. While difficult to quantify exactly we see from Google NGram data Michel et al. (2011) since 2010 the focus of machine learning is increasingly on multi-class classification and other difficult to quantify tasks rather than binary classification, revealing a need for metrics that can accommodate multi-class tasks. Worse still, we appear to be at the edge of another inflection point in AI where Large Language Models (LLMs) and Foundational models Bommasani et al. (2021) like GPT-3 Brown et al. (2020) are ingesting the entire corpus of human thought with limited supervision.

In this paper, we introduce two new metrics based on simple principles, whose purpose is to quantify a change in per-class bias between two or more models. These metrics provide singular data points that are easier to consider than the laborious checking of distributions of class-wise error rates. We will discuss their intuition and their application for comparing the relative performance of deep learning models, and classifiers in general, in terms of bias and fairness. To the best of our knowledge, these new metrics are distinct from all existing methods in that they expose per-group, per-class bias not neatly captured by other metrics, enabling the examination of issues of fairness and bias in great depth. While our proposed work is not a panacea for all emerging trends in AI we believe it represents a starting point to address the current gap.

The remaining sections are organized as follows. In Section 2 we will contextualize the field of fairness metrics and their shortcomings as they relate to our considered problem domain. In Section 3 we will define the intuition for our metrics, and provide a mathematical definition. In Section 4 we will provide specific use cases as experiments we envision the metrics will be used in, and how to reason about their differences. Section 5 will discuss some limitations of our metrics and how we might improve or extend them to other domains.

2 BACKGROUND

Bias and fairness in machine learning have received increasing attention in recent years. The advantages of algorithmic decision-making can be very attractive to large organizations, but there is a risk that the output of these algorithms can be unfair Mehrabi et al. (2019). Unfairness can have serious perceptual and legal consequences for organizations who choose to rely on machines to make important decisions Caton & Haas (2020). This makes it imperative that quantitative measures for bias and fairness in machine learning be defined.

Bias, discrimination, and unfairness are terms that are often used interchangeably but we would like to make a distinction to better dissect the problem. We will refer to bias as meaning the behavior of a machine learning model giving preference to one characterization over another Mooney (1996). Or put simply, having a lower error rate on one class than another. When discussing fairness in this paper, we will be referring to group fairness, which in general is concerned with outcomes for privileged and unprivileged groups Maughan & Near (2020), where a group is a protected feature of an instance from the training data characterizing the instance in some way. Typically we do not want membership in a group to affect the outcome of a prediction, e.g. considering race or gender for ranking resumes or home loan applications.

There are other accepted definitions of fairness as well Mehrabi et al. (2019). Individual fairness, requiring a model given similar predictions for similar individuals. Subgroup fairness, which uses notions of both individual and group fairness by holding some constraint over large collections of a subgroup. However, group fairness is the most commonly measured by metrics of fairness Mehrabi et al. (2019).

Both bias and unfairness can degrade the performance of a model in ways that are not well captured by accuracy, precision, and other measures of ML performance. Biased and unfair models can perform very well on biased or unfair data. A nuanced metric can reveal conditions under which a model's performance might be degraded by the bias of the model or its training data. Many good metrics exist for measuring the group or individual unfairness of a model, but the focus has overwhelmingly been on tasks of supervised, binary classification Caton & Haas (2020). Substantial literature has emerged concerning algorithmic bias, discrimination, and fairness. Mehrabi et al. conducted a survey on bias and fairness in machine learning Mehrabi et al. (2019). Mitchell et al. explored how model cards can be used to provide details on model performance across cultural, racial, and inter-sectional groups and to inform when their usage is not well-suited Mitchell et al. (2019). Gebru et al. proposed using datasheets as a standard process for documenting datasets Gebru et al. (2018). Amini, et al. proposed to mitigate algorithmic bias through re-sampling datasets by learning latent features of images Amini et al. (2019). Wang et al. (2020).

Other metrics of fairness have been described in recent works Borkan et al. (2019); Hinnefeld et al. (2018); Dixon et al. (2018); Maughan & Near (2020) whose purpose is to measure unfairness in machine learning models. Measurements of fairness based on the area under the receiver operating characteristic curve (AUC-ROC) are described in Dixon et al. (2018) and expanded in Borkan et al. (2019). These metrics measure group-wise accuracy using AUC. Prediction Sensitivity, described in Maughan & Near (2020) fills the need for a reliable measure of individual fairness, as opposed to group fairness. A common shortcoming of these metrics is that they focus exclusively on binary classification Caton & Haas (2020) and are not meaningful in tasks of multi-class classification. Our proposed metrics are usable with any number of classes. Additionally, few works have studied how bias can present as unfairness and vice versa. To the best of our knowledge, our work is the first to propose a single metric shown to express both bias and unfairness when comparing two models.

3 PROPOSED METRICS FOR CLASS-WISE BIAS

We propose two new metrics, Combined Error Variance (CEV) and Symmetric Distance Error (SDE)s. Both measure changes in the class-wise false positive and false negative rates of two models, and each has its own advantages which will be explored in Section 4. When calculating both CEV and SDE one model is used as the base and another model as the alternative.

3.1 COMBINED ERROR VARIANCE

The concept of the Combined Error Variance (CEV) metric is to measure the tendency of DNNs to sacrifice one class for the benefit of another class. CEV approximates the variance of the change in False Negative Rate (FNR) and change in False Positive Rate (FPR). It summarizes changes in FNR/FPR away from the model's average. Mathematically, CEV is defined as follows.

$$\delta X_{ie} = \frac{X_{ie} - \hat{X}_{ie}}{\hat{X}_{ie}}$$
(1) $\delta X_{\mu e} = \frac{1}{n} \sum_{i=0}^{n} (\delta X_{ie})$ (2)

$$cev = \frac{1}{n} \sum_{i=1}^{n} (dist((\delta X_{\mu pos}, \delta X_{\mu neg}), (\delta X_{i pos}, \delta X_{i neg})))^2$$
(3)

Let X_{ie} be a pair of values for the FPR and FNR for class *i* of the comparison model and \hat{X}_{ie} be the original models FPR/FNR pair, with e indicating either the false-positive or false-negative rate. We first find the normalized change in FPR/FNR δX_{ie} by subtracting the error rates for the two models from each other and dividing by the original. The mean change $\delta X_{\mu e}$ is found by averaging the values of δX_{ie} , keeping in mind that every δX_{ie} is the change in two values FPR and FNR. The CEV is calculated by treating each δX_{ie} as a point in a 2-dimensional space of FNR and FPR. The square of the euclidean distances between each δX_{ie} and the mean change represented by $\delta X_{\mu e}$ are summed and divided by the total number of classes *n*.

3.2 Symmetric Distance Error

The principle of the Symmetric Distance Error (SDE) metric is to measure another undesirable bias behavior that presents in simple models. That is, a class with more training examples or that has similar features to another class is more frequently to be chosen by the model with limited capacity. To reflect this biased behavior, SDE calculates how "far away" from balanced is the change in FPR/FNR for each single class error. Intuitively, if we make a scatter plot with changes in FPR and FNR as X and Y values, the diagonal line in that plot would be a perfectly balanced change in FPR/FNR. Therefore, the SDE can be calculated as the symmetric distance of each change to that balance line.

$$d = \frac{|a(x_0) + b(y_0)|}{\sqrt{a^2 + b^2}}$$
(4)
$$d = \frac{|(1)(x_0) + (-1)(y_0)|}{\sqrt{(1)^2 + (-1)^2}} = \frac{|x_0 - y_0|}{\sqrt{2}}$$
(5)

For a line in the Cartesian plane described by the equation ax + by + c = 0, the distance d from any point (x_0, y_0) can be derived from the equation in 4. In our specific context, the diagonal of the Cartesian plane (i.e. the balance line) is x = y or x - y = 0 will represent an equal difference in FNR and FPR between two models. Given any change of FNR and FPR the symmetric distance of that change to the balance line can be calculated as:

$$sde = \frac{1}{n} \sum_{i=0}^{n} |\delta F N R_i - \delta F P R_i|$$
(6)

Once the symmetric distance of each change is calculated, the SDE of a model can be calculated as the mean absolute change of normalized FP/FN rate, with the change being calculated as described in Equation 1. The $\sqrt{2}$ has been omitted from the final equation as a constant that has no effect on the meaning of the metric. This metric will therefore reveal that one model or the other is more biased toward false positives or false negatives in a class-wise fashion.

3.3 NORMALIZATION

It is frequently true that a metric is meaningless without some numerical context. Both CEV and SDE may produce a large range of values depending on the specific dataset, number of classes, and performance of the models trained on that data. While not strictly necessary, in order to make the outputs of our metrics more interpretable, we follow a procedure for normalizing their values based on a hypothetical "worst performing" model to give us a reference. To do this a set of predictions for all test instances is produced at random with all classes being equally likely, and the FPR/FNR of these random predictions is calculated. The CEV and SDE of the random predictions is generated relative to the original model. These are then used as a divisor to normalize the other CEV and SDE values of a group of models. Following this process, our metrics now indicate a change in algorithmic bias relative to a random predictor. Thus, a CEV value of 0.5 shows that the class-wise bias of model 2 relative to model 1 has increased by 50% of the change between model 1 and a random predictor.

4 EXAMPLE APPLICATIONS

We have explored several applications of CEV and SDE for comparing the performance of two models. While we don't believe this list is exhaustive, in this section we illustrate several ways our proposed metrics can be used. We group these applications into two categories:

- 1. Comparing models w.r.t each other for the purpose of evaluating change in bias. We demonstrate using CEV/SDE in the context of model compression to detect compression-induced bias. We then further generalize this concept by informing and selecting from any number of trained low resource models to replace a higher capacity model.
- 2. Evaluating group fairness. We demonstrate how CEV/SDE can be used to measure relative bias w.r.t protected groups. We also compare our results to existing binary classification fairness metrics and demonstrate the use of our metrics on multi-class data.

4.1 MODEL SELECTION

4.1.1 MODEL COMPRESSION

Compressed neural networks can offer significant reductions in the computing power required for model inference. However, it has been shown that compressed models are often more biased than the

original when the per-class error rates are examined Hooker et al. (2019; 2020). CEV and SDE allow one to reason about two values instead of auditing all FP/FN rates of the classes. They can quickly expose compressed models that cannibalize a subset of their target classes to preserve top-1/top5 accuracy.

COMPRESSION IDENTIFIED EXEMPLARS

To the best of our knowledge, only one other work proposes a method for measuring desperate impact and bias for multi-class classification. Compression Identified Exemplars (CIEs) Hooker et al. (2019; 2020) are proposed specifically to categorize changes in model behavior and attribute the extent to which model compression techniques are responsible. While CIEs have many advantages beyond measuring bias (i.e. human-in-the-loop data auditing) the means by which they are found make counting CIEs impractical for many problems and impossible for many others. CIEs are defined as the following equation, which represents images in a data set for which compression explicitly changes the behavior of the answer. In Hooker et al. (2019), a population of 30 models were trained for each compression method and sparsity level. They then defined an image i as an exemplar if modal label $Y_{i,t}^M$, or class most predicted by the *t*-compressed model population disagrees with the label produced from the original networks.

$$CIE_{i,t} = \begin{cases} 1, & \text{if } y_{i,0}^M \neq y_{i,t}^M \\ 0, & \text{otherwise} \end{cases}$$

While CIEs help reveal the bias issues in pruned models, not every image reported as a CIE represents a problem. A good portion of CIEs represent images that are equally hard for a human to classify and may simply be a case where the uncompressed networks overfit to learn the example. Our proposed metrics, CEV SDE, address some of the weaknesses of CIEs and give us true metrics. Pruning and quantization have been observed to sacrifice accuracy on a subset of classes in classification tasks in order to retain overall top-k accuracy Hooker et al. (2019). To catch and reflect this bias, our metrics are designed to quantify both the *spread* of the change in classification error as well as changes in *how* the model is making mistakes. As a result, both of our proposed metrics consider the distribution of change in false positive and false negative rates (FPR, FNR) for all classes.

To demonstrate the efficacy of CEV and SDE for evaluating bias in compressed models, we evaluate the following well-known compression algorithms and compare CEV and SDE metrics with the CIEs counts. The following example is conducted to illustrate the intuitive application of CEV/SDE. In this experiment, we measure the change that structured pruning has on a convolutional neural network as well as how various distillation methods mitigate bias. We train a CNN ResNet He et al. (2016) (ResNet32x4) model on CIFAR100 Krizhevsky (2009) and prune it to various sparsities using Filter-wise Structured Pruning.

Hooker et al. reported that the negative effects of compression are most observed on underrepresented groups Hooker et al. (2020). Therefore, we resample the CIFAR100 dataset to purposely underrepresent certain classes in the training set. We select 15 classes and reduce the remaining training samples to 50%, 20%, and 10% of their original numbers. Our hypothesis is that these biased classes with fewer data samples will more likely be picked as "victims" by the pruned models.

We use pruning as the baseline model and the rest of the models are pruned jointly with one of the state-of-the-art KD methods shown in Table 1. We also test whether combining our feature map based distillation methods with KD (e.g. AT + KD or PKT + KD) can achieve additional benefits. We utilize Tian et al's Tian et al. (2019) implementation for all distillation methods tested. For all experiments, the models are pruned initially to 10% sparsity then gradually pruned every 5 epochs until the desired sparsity is reached according to the AGP Zhu & Gupta (2017) schedule until they reach 45% sparsity. All pruning is completed at the halfway point of training and allowed to continue to fine-tune for another 120 epochs. We do not perform any layer sensitivity analysis or prune layers at different ratios. Although that may have resulted in higher accuracy, our goal is not to reach state-of-the-art compression ratios but to demonstrate how CEV/SDE capture the effect of pruning on bias and any methods that might mitigate it.

Two observations can be made from the results presented in Table 1 and Figure 2a: (1) CEV and SDE generally agree with the CIE count. They achieve this while being easier to calculate and not



Figure 2: Normalized FP/FN Rate Change for distillation methods on biased CIFAR100 dataset. Ellipses represent 95% confidence intervals for data

Table 1: CIE count, CEV, SDE, and Accuracy for pruning with and without KD on biased CIFAR100

Method	# of CIEs	CEV	SDE	Accuracy
AT + KD	742	0.00187	0.13173	77.100
PKT + KD	748	0.00199	0.13098	77.335
SP + KD	768	0.00331	0.16162	76.927
FSP + KD	742	0.00333	0.16002	75.285
KD (Hinton et al., 2015)	770	0.00338	0.16065	78.142
AT (Zagoruyko & Komodakis, 2017)	909	0.00430	0.19306	78.097
PKT (Passalis & Tefas, 2018)	881	0.00481	0.19891	78.963
SP (Tung & Mori, 2019)	838	0.00583	0.21591	78.520
FSP (Yim et al., 2017)	877	0.00638	0.22525	78.413
Struct Pruning	887	0.00931	0.26687	77.242

requiring multiple models to be trained. (2) Accuracy alone is a poor indicator of model quality. In Table 1 Structured pruning has accuracy comparable to AT + KD and PKT + KD but in Figure 2 we see that the change in accuracy in structured pruning is not distributed equitably with some classes having a 300% change in FNR. SDE also neatly captures that the skew in FP/FN resulting for most models. These experiments show the great potential of our proposed CEV/SDE metrics in distinguishing desirable models from biased models that appear to be equal at a surface level.





Figure 3: Change-in-Top1, normalized CEV, and normalized SDE adjacency matrices of models listed in Table 2. Each entry displays the given metric calculated for the columns model w.r.t the rows model. Reading the table row-wise you see trade offs going from the row model to another. Reading the column you see the trade offs for other models going to the column model. Higher values CEV/SDE (shown by darker cells) indicate moving towards a more biased model.

There are many algorithms and model architectures for low resource inference in image recognition alone. We have discussed examples of using CEV/SDE to analyze compression effects. Now we consider this problem more broadly. Bias is an important consideration when selecting a pre-trained

index	model	top1	top5	img size	params $x10^6$
0	efficientnet_b2 Tan & Le (2019)	80.608	95.310	288	9.11
1	efficientnet_b1 Tan & Le (2019)	78.792	94.342	256	7.79
2	efficientnet_b1_pruned Aflalo et al. (2020)	78.242	93.832	240	6.33
3	mobilenetv3_large_100_miil Howard et al. (2019)	77.914	92.914	224	5.48
4	mobilenetv2_120d Sandler et al. (2018)	77.294	93.502	224	5.83
5	mobilenetv3_large_100 Howard et al. (2019)	75.768	92.540	224	5.48
6	mobilenetv3_rw	75.628	92.708	224	5.48
7	mobilenetv2_110d Sandler et al. (2018)	75.052	92.180	224	4.52
8	pit_ti_distilled_224 Heo et al. (2021)	74.536	92.096	224	5.10
9	deit_tiny_distilled_patch16_224 Touvron et al. (2021)	74.504	91.890	224	5.91
10	mobilenetv2_100	72.978	91.016	224	3.50
11	resnet18 He et al. (2016)	69.758	89.078	224	11.69

Table 2: Low resource ImageNet models from TIMM github Wightman (2021). Top1/Top5, input image sizes, and parameter count listed. Index corresponds to axis labels in Figure 3

model from one of the dozens which are available in many problem spaces. Unfortunately, CIE count is not applicable in the case where you have models already trained and simply want to understand the trade-offs you will be making. Here we see how one might use CEV/SDE to detect and avoid a model more biased than models of similar accuracy. For this example, we have selected a set of models from the TIMM model repository (Wightman, 2021) that have between 3.5×10^6 and 11.7×10^6 parameters. Each model has been pre-trained on the Imagenet dataset (Russakovsky et al., 2015). Table 2 lists the specific models, their top1/top5 accuracy, image input size, and number of parameters. We have constructed heat maps of the CEV/SDE values by calculating the interaction between each model and building an adjacency matrix. In both Table 2 and Figure 3 we sort the models by Top-1 accuracy. With our constructed matrices we can quickly glance and observe that mobilenetv3_large100 on row 5 column 5 stands out clearly in the CEV/SDE matrices. We see that although the model has comparable accuracy and parameters to mobilenetv3_rw and mobilenetv2_110d, it is actually measured to have worse trade-offs of FPR/FNR w.r.t to the tables best model in terms of accuracy efficientnet_b2, and is no better or worse than several of the next several models on our accuracy sorted list. CEV and SDE have prevented us from making a poor selection with relative ease. Again, we find accuracy alone is a poor indicator of model quality.

4.2 FAIRNESS

Fairness has rightly been enjoying increased attention over the last several years when measured by the total number of papers addressing it Caton & Haas (2020). Fairness is often defined as the ability of a model to classify all groups within the testing data equally well. For example, a model trained to recognize human faces should be equally good at recognizing the faces regardless of demographic traits (e.g race, gender, age). Unfortunately, unintentionally biased data collected in real-world datasets and even train methodologies can cause undesired performance in models. Our metrics were developed specifically to measure the bias of classifiers, but we will demonstrate they may also be used for measuring fairness as well. Importantly, this methodology allows the metrics to measure fairness in multi-class examples.

To measure bias with CEV or SDE, one model is compared to another. This process can be adapted to measure fairness by comparing a models performance on its test data to its performance on a subset of its test instances. For this purpose, we will select from specific protected attributes and calculate bias with respect to the groups. A large value in CEV or SDE will indicate that per-class bias is increased for one group in the data, and that the model's performance is lower for that group.

4.2.1 BINARY CLASSIFICATION

To demonstrate measuring fairness in binary classification, we trained several common machine learning models on the Titanic dataset Frank E. Harrell Jr. (2017): a shallow neural network(NN), a support vector machined(SVM), and a gradient tree boosting classifier(GTB). This dataset offers information about Titanic passengers with the labels Survived and Did Not Survived. The sex of each passenger is included as a feature of each instance. Sex was excluded in the model training and used later for group-wise fairness testing. These metrics are presented along with the False Positive Equality Difference(FPED), False Negative Equality Difference(FNED)Dixon et al. (2018),

Table 3: Comparison of Error Rate Equality Difference(ERED) Dixon et al. (2018) metrics, and Difference in Expected Value(DEV) Hinnefeld et al. (2018) Metrics with our proposed CEV and SDE on measuring bias. We make our comparison to measure fairness traits of models classifying the Titanic dataset Frank E. Harrell Jr. (2017). CEV/SDE are calculated w.r.t the whole dataset errors, and given protected group. All values are averaged over 5 runs of train/test.

Model	Our Metrics				Existing Metrics			
-	CEV		SDE		ERED		DEV	
-	$All{\rightarrow}Men$	All→Women	$All {\rightarrow} Men$	$All {\rightarrow} Women$	FPED	FNED	DIMS	DIAMR
NN	0.013557	0.012737	0.115002	0.093218	0.548443	0.458016	-0.269742	0.288790
SVM	0.012089	0.000736	0.109744	0.027081	0.412500	0.593508	-0.067460	0.491071
GTB	0.000107	0.000941	0.010341	0.030619	0.458462	0.513932	-0.193700	0.364831

Difference in Mean Scores(DIMS), and Difference in Average Model ResidualsHinnefeld et al. (2018) in Table 3.

The four metrics presented for comparison are all zero for perfectly fair predictions. The relatively small value generated for each of the eight metrics is an effect of the small size of the dataset. The fact the FPED, FNED, DIMS, and DIAMR are not 0, shows that some unfairness has been learned by our neural network. The differences in the CEV and SDE scores moving from all data to men only, and all data to women only also indicates biased and unfair performance by the classifier. So we can confirm that in tasks of binary classification, our new metrics conform to the established work in the field of fairness. But as will be shown in Section 4.2.2, CEV and SDE are not limited to the analysis of binary labels.

4.2.2 MULTI-CLASS CLASSIFICATION



Figure 4: Change in FP/FN rate for protected subgroups of ResNet model trained on CelebA dataset. Change is calculated w.r.t the to complete validation set

We have asserted that CEV and SDE can be used to measure fairness in multi-class classification. We will now demonstrate that process using the CelebA dataset Liu et al. (2015). This dataset contains several thousand images of celebrates and public figures with 40 binary attributes. We have selected from the provided attributes a subset representing hair color to serve as training labels. We then trained a ResNet34 image recognition model to identify the hair color of the image subjects. From the remaining provided attributes, we have identified several to serve as protected groups ("Attractive", "Male", "Pale Skin", "Young"). As these labels come from what might be described as "privileged", we also consider subsets formed from the conjugate of these labels. We should note that the conjugate does not imply the opposite. The absence of a Pale Skin label for example does not explicitly mean dark skin but would contain all of those examples.

The results are contained in Figure 4 and Table 4. We find that groups "Male" and "Pale Skin" have the highest Top-1 Accuracy. However, we also find they have high levels of class unfairness. Specifically for Male, our model is far less likely to correctly identify Male as having Blond Hair,

Protected Attribute	Top-1	CEV	SDE	Change in FPR	Change in FNR
Full Test Set	0.9212				
Attractive	0.9222	0.0015	0.0331	-31.0809	80.4380
Male	0.9225	0.1413	0.2205	12.8440	77.8003
Pale Skin	0.9224	0.0035	0.0465	-43.8572	-33.9335
Young	0.9215	0.0002	0.0082	-27.6765	150.5065
Not Attractive	0.9208	0.0034	0.0493	45.6423	6.8297
Not Male	0.9207	0.0053	0.0562	1.2762	47.2981
Not Pale_Skin	0.9207	0.0000	0.0021	1.9565	1.4648
Not Young	0.9213	0.0035	0.0313	146.3057	0.2381

Table 4: Top-1, CEV, SDE, and change in FPR/FNR for selected protected class from ResNet model trained on CelebA dataset

and more likely to incorrectly guess they have Gray Hair. Meanwhile, "Not Pale Skin" has lower accuracy, but the accuracy and FPR/FNRs are much closer to the average of the model as a whole. This is easily visible in Figure 4. This unevenness is neatly captured by the corresponding CEV and SDE values or the groups in our data.

5 DISCUSSION AND LIMITATIONS

As with any metric, it is also important to remember that CEV and SDE are only meaningful in context. A higher value for CEV indicates that the second model has a higher class-wise bias. A higher value for SDE indicates that the second model is skewing towards false positives or false negatives. Either behavior represents a degraded real-world performance for a model in a way that may not be captured by accuracy or precision as demonstrated in Section 4.

The importance of measuring fairness and bias in multi-class data should not be underestimated. Data that meaningfully describes the real world is often multi-class. While it is true to that multiclass classification can be re-framed as many binary classification problems, re-framing a problem as 100 or 1,000 one-vs-each problems would only serve to make reasoning about the implications much more difficult. We believe CEV and SDE are applicable to many real-world problems completely ignored by their binary cousins.

CEV and SDE can be used to measure the fairness of a machine learning model, but only group fairness. Individual fairness, which is defined as the degree to which similar individuals are classified similarly, is not measured in any of the use cases presented in Section 4.

We have not found any consistent threshold that indicates by itself that a model is or is not biased. It may be that such a threshold exists. Also important to remember that biased performance may be the result of algorithmic bias, or it may be a reflection of biased data and CEV/SDE alone cannot determine its source. But with these limitations in mind, CEV and SDE reliably indicate that one model is more or less biased than another. As concluded in Hinnefeld et al. (2018), "...fairness metrics in machine learning must be interpreted with a healthy dose of human judgment."

CEV and SDE are calculated w.r.t to some other classifier and only classifiers. As such they are not suitable for every situation. However, we believe they provide a good starting point for the community to begin to address measuring more sophisticated machine learning tasks. Additionally, we endeavor to extend the concepts of CEV/SDE to other tasks like image segmentation which are harder still to quantify. We also believe our insights from CEV/SDE can be used to create standalone metrics to measure bias and fairness without making direct model comparisons.

6 CONCLUSION

Unfairness is a persistent and difficult problem in machine learning. Bias is more quantifiable but just as dangerous to the reliable performance of machine learning models in the real world. In this paper, we have introduced two new metrics: CEV and SDE. These metrics can reliably reveal that a model is more or less biased compared to another model. We have also demonstrated that these new metrics can be used to measure the fairness of a model used for classification. Importantly, these metrics are meaningful when used with multi-class data, even with a very large number of classes.

REFERENCES

- Yonathan Aflalo, Asaf Noy, Ming Lin, Itamar Friedman, and Lihi Zelnik. Knapsack pruning with inner distillation, 2020.
- Alexander Amini, Ava P Soleimany, Wilko Schwarting, Sangeeta N Bhatia, and Daniela Rus. Uncovering and mitigating algorithmic bias through learned latent structure. In *Proceedings of the* 2019 AAAI/ACM Conference on AI, Ethics, and Society, pp. 289–295, 2019.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258, 2021.
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pp. 491–500, 2019.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
- Simon Caton and Christian Haas. Fairness in machine learning: A survey. arXiv preprint arXiv:2010.04053, 2020.
- Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference* on AI, Ethics, and Society, pp. 67–73, 2018.
- Thomas Cason Frank E. Harrell Jr. Titanic dataset, oct 2017. URL https://www.openml.org/d/40945.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778, 2016.
- Byeongho Heo, Sangdoo Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh. Rethinking spatial dimensions of vision transformers. *arXiv preprint arXiv:2103.16302*, 2021.
- J Henry Hinnefeld, Peter Cooman, Nat Mammo, and Rupert Deese. Evaluating fairness metrics in the presence of dataset bias. *arXiv preprint arXiv:1809.09245*, 2018.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2015.
- Sara Hooker, Aaron Courville, Gregory Clark, Yann Dauphin, and Andrea Frome. What do compressed deep neural networks forget? *arXiv preprint arXiv:1911.05248*, 2019.
- Sara Hooker, Nyalleng Moorosi, Gregory Clark, Samy Bengio, and Emily Denton. Characterising bias in compressed models. *arXiv preprint arXiv:2010.03058*, 2020.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1314–1324, 2019.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- Krystal Maughan and Joseph P Near. Towards a measure of individual fairness for deep learning. arXiv preprint arXiv:2009.13650, 2020.

- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635*, 2019.
- Jean-Baptiste Michel, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K Gray, Google Books Team, Joseph P Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, et al. Quantitative analysis of culture using millions of digitized books. *science*, 331(6014):176–182, 2011.
- Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In Proceedings of the conference on fairness, accountability, and transparency, pp. 220–229, 2019.
- Raymond J Mooney. Comparative experiments on disambiguating word senses: An illustration of the role of bias in machine learning. *arXiv preprint cmp-lg/9612001*, 1996.
- Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pp. 6105–6114. PMLR, 2019.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv* preprint arXiv:1910.10699, 2019.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pp. 10347–10357. PMLR, 2021.
- Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1365–1374, 2019.
- Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. Towards fairness in visual recognition: Effective strategies for bias mitigation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8919–8928, 2020.
- Ross Wightman. Pytorch image models). https://github.com/rwightman/ pytorch-image-models, August 2021.
- Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, pp. 4133–4141, 2017.
- Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR*, 2017. URL https://arxiv.org/abs/1612.03928.
- Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878*, 2017.