Democratizing the Development of Fair and Effective Machine Learning Systems



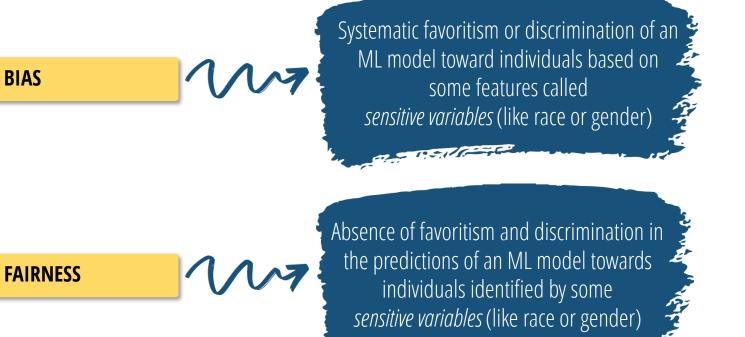
SOBIG**DATA***

Giordano d'Aloisio

Università degli Studi dell'Aquila / Italy

Al-Gap 20th November 2023, L'Aquila

Let's formally define bias and fairness



Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A Survey on Bias and Fairness in Machine Learning. ACM Computing Surveys, 54(6), 1–35. https://doi.org/10.1145/3457607

Is the concept of bias and fairness that simple?

Actually not...

A Survey on Bias and Fairness in A Survey on Bias and Fairnes

- (1) Measurement Bias. Mea measure particular feature cidivism risk prediction to used as proxy variables to viewed as mismeasured p are controlled and policed should not conclude that rates, therefore they are r assessed and controlled [1 (2) Omitted Variable Bias, (are left out of the model [: designs a model to predict. customers will stop subsc are canceling their subscr Now imagine that the reas competitor in the market t of the competitor was som to be an omitted variable. (3) Representation Bias. Re
- (5) Representation Blas. Re ing data collection process | tion, with missing subgrou like ImageNet (as shown i cultures.
 (4) Aggregation Blas. Aggre
- drawn about individuals fr can be seen in clinical aid differences across ethnicit to diagnose and monitor Therefore, a model that is ethnic and gender groups sented equally in the train population can result in as (a) Simpson's Paradox. analysis of heterogene in aggregated data dis underlying subgroups paradox arose during t ley [16]. After analyzin toward women, a sma compared to their male analyzed over the depa a small advantage ove partments with lower served in a variety of d and computational soc
- (b) Modifiable Areal Un when modeling data a different trends learne

three subgroups, which study was conducted, c ure 2(b)). The positive pletely disappears (soli green lines) are unaffed (6) Longitudinal Data F analysis to track cohor modeled using cross-s point. The heterogeneo clusions than longitudi that comment length (cross-sectional snapsho joined Reddit in differ length within each coh (7) Linking Bias. Linking tivities, or interactions a [101] authors show how considering the links in in the network. Referen

consider the example in

(1) from social link pattern The differences and bia sampling, as shown in cause different types of 3.1.2 Algorithm to User, j

 introduce biases in user beha rithmic outcomes and affect
Algorithmic Bias. Al added purely by the al optimization functions as a whole or consideri algorithms [44], can all of the algorithms.
User Interaction Bias the Web but also get tri by imposing his/her sel influenced by other ty;
(a) Presentation Bias

 (a) I reservation plass example, on the We gets clicks, while e does not see all the (b) Ranking Bias. The result in attraction crowdsourcing app

eling data a ends learne (3) **Popularity Bias**. Item metrics are subject to n instance, this type of bias where popular objects w be a result of good qualii (4) Emergent Bias. Emerge bias arises as a result of some time after the comp in user interfaces, since i prospective users by des as discussed in detail in 1 (5) Evaluation Bias. Evalua use of inappropriate and as Adience and IJB-A be recognition systems thai examples for this type of

115:8

3.1.3 User to Data. Many c inherent biases in users migh behavior is affected/modulatec troduce bias in the data genera

(1) Historical Bias. Histori world and can seep into feature selection [140]. A result where searching fc to the fact that only 5% results to be biased towa the reality, but whether worth considering.

- (2) Population Bias. Popul du ser characteristics are di ser duracteristics are di population [116]. Popula of bias can arise from di women being more likely in online forums like Remedia use among young background can be foum-ri (3) Self-selection Bias. Self subjects of the research si an opinion poll to measu su popreters are more like si supporters are more like the (4) Social Bias. Social bias 1
 - this type of bias can be a when influenced by othe being too harsh [9, 147].
 (5) Behavioral Bias. Behav or different datasets [116]

or different datasets [116] where authors show how A Survey on Bias and Fairness in Machine Learning

different reactions and behavior from people and sometimes even leading to communication

- (6) Temporal Bias. Temporal bias arises from differences in populations and behaviors over time [116]. An example can be observed in Twitter where people talking about a particular topic start using a hashtag at some point to capture attention, then continue the discussion about the event without using the hashtag [116, 142].
- (7) Content Production Bias. Content Production bias arises from structural, lexical, semantic, and syntactic differences in the contents generated by users [116]. An example of this type of bias can be seen in Reference [114] where the differences in use of language across different gender and age groups is discussed. The differences in use of language can also be seen across and within countries and populations.

Existing work tries to categorize these bias definitions into groups, such as definitions falling solely under data or user interaction. However, due to the existence of the feedback loop phenomenon [36], these definitions are intertwined, and we need a categorization that closely models this situation. This feedback loop is not only existent between the data and the algorithm, but also between the algorithms and user interaction [29]. Inspired by these papers, we modeled categorization of bias definitions, as shown in Figure 1, and grouped these definitions on the arrows of the loop where we thought they were most effective. We emphasize the fact again that these definitions are intertwined, and one should consider how they affect each other in this cycle and address them accordingly.

3.2 Data Bias Examples

There are multiple ways that discriminatory bias can seep into data. For instance, using unbalanced data can create biases against underrepresented groups. Reference [166] analyzes some examples of the biases that can exist in the data and algorithms and offers some recommendations and suggestions toward mitigating these issues.

3.2.1 Examples of Bias in Machine Learning Data. In Reference [24], the authors show that datasets such as UB-A and 8.6.2% in Adience are imbalanced and contain maindy light-skinned subjects—75.6% in IJB-A and 8.6.2% in Adience. This can bias the analysis towards dark-skinned groups who are underrepresented in the data. In another instance, the way we use and analyze our data can create bias when we do not consider different subgroups in the data. In Reference [24], the authors also show that considering only male-female groups into light-skinned females, light-skinned males, and ark-skinned females. It is only in this case that we can clearly observe the bias towards dark-skinned females, as previously dark-skinned finales would compromise for dark-skinned females, as previously dark-skinned shi subgroups. Popular machine-learning datasets that serve as a base for most of the developed algorithms and tools can also be biased—which can be harmful to the downstream applications that are based in machine learning. In Reference [138], researchers showed that these datasets suffer from representation bias and datasets.

3.2.2 Examples of Data Bias in Medical Applications. These data biases can be more dangerous in other sensitive applications. For example, in medical domains there are many instances in which the data studied and used are skewed toward certain populations—which can have dangerous consequences for the underrepresented communities. Reference [97] showed how exclusion of African-Americans resulted in their misclassification in clinical studies, so they became advocates

At least 23 different definitions of bias and fairness are available from the literature

115:9

From many definitions come many metrics

Generic metrics

<pre>metrics.num_samples (y_true[, y_pred,])</pre>	Compute the number of samples.
<pre>metrics.num_pos_neg (y_true[, y_pred,])</pre>	Compute the number of positive and negative samples.
${\tt metrics.specificity_score} \ ({\tt y_true}, {\tt y_pred}, {\tt *})$	Compute the specificity or true negative rate.
<pre>metrics.sensitivity_score (y_true, y_pred[,])</pre>	Alias of sklearn.metrics.recall_score() for binary classes only.
<pre>metrics.base_rate (y_true[, y_pred,])</pre>	Compute the base rate, $Pr(Y = \text{pos_label}) = rac{p}{P+N}.$
<pre>metrics.selection_rate (y_true, y_pred, *[,])</pre>	Compute the selection rate, $Pr(\hat{Y} = ext{pos_label}) = rac{TP+FP}{P+N}.$
<pre>metrics.smoothed_base_rate (y_true[, y_pred,])</pre>	Compute the smoothed base rate, $\frac{P+lpha}{P+N+ R_Y lpha}.$
metrics.smoothed_selection_rate (y_true,)	Compute the smoothed selection rate, $\frac{TP+FP+lpha}{P+N+ R_Y lpha}.$
metrics.generalized_fpr (y_true, probas_pred, *)	Return the ratio of generalized false positives to negative examples in the dataset, $GFPR=\frac{GFP}{N}.$
$\tt metrics.generalized_fnr \ (y_true, probas_pred, *)$	Return the ratio of generalized false negatives to positive examples in the dataset, $GFNR=\frac{GFN}{P}.$

Individual fairness metrics

<pre>metrics.generalized_entropy_index (b[, alpha])</pre>	Generalized entropy index measures inequality over a population.
$\tt metrics.generalized_entropy_error (y_true, y_pred)$	Compute the generalized entropy.
<pre>metrics.theil_index (b)</pre>	The Theil index is the <code>generalized_entropy_index()</code> with $lpha=1.$
${\tt metrics.coefficient_of_variation} \ (b)$	The coefficient of variation is the square root of two times the generalized_entropy_index() with $\alpha=2.$
${\tt metrics.consistency_score} \; (X, y[, n_neighbors])$	Compute the consistency score.

Group fairness metrics

metrics.statistical_parity_difference (y_true)	Difference in selection rates.
<pre>metrics.mean_difference (y_true[, y_pred,])</pre>	Alias of statistical_parity_difference() .
<pre>metrics.disparate_impact_ratio (y_true[,])</pre>	Ratio of selection rates.
$\tt metrics.equal_opportunity_difference~(y_true,)$	A relaxed version of equality of opportunity.
metrics.average_odds_difference (y_true,)	A relaxed version of equality of odds.
<pre>metrics.average_odds_error (y_true, y_pred, *)</pre>	A relaxed version of equality of odds.
<pre>metrics.class_imbalance (y_true[, y_pred,])</pre>	Compute the class imbalance, $\frac{N_u-N_p}{N_u+N_p}.$
metrics.kl_divergence (y_true[, y_pred,])	Compute the Kullback-Leibler divergence, $KL(P_p P_u) = \sum_y P_p(y) \log \Big(rac{P_p(y)}{P_u(y)} \Big)$
metrics.conditional_demographic_disparity (y_true)	Conditional demographic disparity, $CDD = rac{1}{\sum_i N_i} \sum_i N_i \cdot DD_i$
<pre>metrics.smoothed_edf (y_true[, y_pred,])</pre>	Smoothed empirical differential fairness (EDF).
metrics.df_bias_amplification (y_true, y_pred, *)	Differential fairness bias amplification.
metrics.between_group_generalized_entropy_error ()	Compute the between-group generalized entropy.
<pre>metrics.mdss_bias_scan (y_true, probas_pred)</pre>	DEPRECATED: Change to new interface - aif360.sklearn.detectors.mdss_detector.bias_scan by version 0.5.0.
<pre>metrics.mdss_bias_score (y_true, probas_pred)</pre>	Compute the bias score for a prespecified group of records using a given scoring function.

At least 29 different bias and fairness metrics are available in the AIF360 repository

Bias mitigation methods

aif360.algorithms.preprocessing

$\verb+algorithms.preprocessing.DisparateImpactRemover ([])$	Disparate impact remover is a preprocessing technique that edits feature values increase group fairness while preserving rank- ordering within groups [1]
algorithms.preprocessing.LFR $([, k, Ax,])$	Learning fair representations is a pre-processing technique that finds a latent representation which encodes the data well but obfuscates information about protected attributes [2]
algorithms.preprocessing.OptimPreproc ([,])	Optimized preprocessing is a preprocessing technique that learns a probabilistic transformation that edits the features and labels in the data with group fairness, individual distortion, and data fidelity constraints and objectives [3]
algorithms.preprocessing.Reweighing (\dots)	Reweighing is a preprocessing technique that Weights the examples in each (group, label) combination differently to ensure fairness before classification [4]

aif360.algorithms.postprocessing

${\tt algorithms.postprocessing.CalibratedEqOddsPostprocessing}\;()$	Calibrated equalized odds postprocessing is a post- processing technique that optimizes over calibrated classifier score outputs to find probabilities with which to change output labels with an equalized odds objective [7]
$\verb"algorithms.postprocessing.EqOddsPostprocessing" ()$	Equalized odds postprocessing is a post-processing technique that solves a linear program to find probabilities with which to change output labels to optimize equalized odds [8]_ [9]
$\verb"algorithms.postprocessing.RejectOptionClassification"()$	Reject option classification is a postprocessing technique that gives favorable outcomes to unpriviliged groups and unfavorable outcomes to priviliged groups in a confidence band around the decision boundary with the highest uncertainty [10]

aif360.algorithms.inprocessing

algorithms.inprocessing.AdversarialDebiasing ()	Adversarial debiasing is an in-processing technique that learns a classifier to maximize prediction accuracy and simultaneously reduce an adversary's ability to determine the protected attribute from the predictions [5]
algorithms.inprocessing.ARTClassifier ()	Wraps an instance of an $\ensuremath{\operatorname{art.classifiers.classifier}}$ to extend $\ensuremath{\operatorname{Transformer}}$.
algorithms.inprocessing.GerryFairClassifier $([])$	Model is an algorithm for learning classifiers that are fair with respect to rich subgroups.
$\verb+algorithms.inprocessing.MetaFairClassifier ([])$	The meta algorithm here takes the fairness metric as part of the input and returns a classifier optimized w.r.t.
algorithms.inprocessing.PrejudiceRemover $([])$	Prejudice remover is an in-processing technique that adds a discrimination-aware regularization term to the learning objective [6]
${\tt algorithms.inprocessing.ExponentiatedGradientReduction} \ ()$	Exponentiated gradient reduction for fair classification.
algorithms.inprocessing.GridSearchReduction ()	Grid search reduction for fair classification or regression.

14 bias mitigation methods are available in the AIF360 repository... but many more are available from the literature!

What does it mean?



At least 23 different definitions of bias and fairness



At least 29 different bias and fairness metrics



At least 14 different bias mitigation algorithms



How can we solve this issue?

- Software engineering approaches can help us to formalise and standardise the development of fair ML systems
- Having a more formal and standard workflow will ease the development of fair ML systems and make it accessible also to nonexpert users

To this aim we propose MANILA, a web-based application to *democratize* the development of fair and effective (i.e., correct) ML systems

MANILA

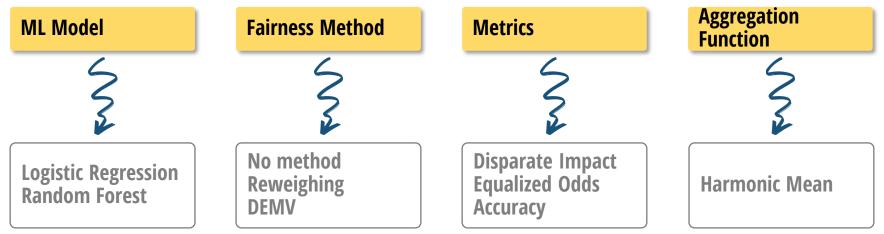
- MANILA is a tool that guides users in defining and performing fairness and effectiveness evaluations of different ML models and fairness enhancing methods
- Automatically disables methods and metrics that are not compatible with other selected features
- Eventually selects and returns the setting having the best fairness and effectiveness trade-off, based on the selected metrics
- Freely available in the SoBigData RI: <u>https://sobigdata.d4science.org/group/sobigdata.it/manila-univaq</u>



MANILA in action

- We train a Logistic Regression and a Random Forest classifier to predict the recidivism of condemned people using the COMPAS dataset
- We evaluate the fairness and effectiveness of different settings against *non-white* people







Thank you for your attention!

UNIVERSITÀ DEGLI STUDI DELL'AQUILA



DISIM Dipartimento di Ingegneria e Scienze dell'Informazione e Matematica