

Contributions in VL-XAI to Address and Democratize Software Fairness



Giordano d'Aloisio

SoBigData@Univaq - 25 March, 2025

giordano.daloisio@univaq.it

Università degli Studi dell'Aquila / Italy



COMPAS

- COMPAS is an ML algorithm used by some courts in the US to predict recidivism of condemned people
- A study showed that, given two people with the same features but different ethnicity, the system was giving higher probability of recidivism to non-white people



Bervard Parker, left, was rated high risk; Dylan Progett was rated low risk. [Derek Ritchie for ProPublica]

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

ON A SPRING AFTERNOON IN 2014, Brisha Borden was running late to pick up her god-sister from school when she spotted an unlocked kid's blue Huffy bicycle and a silver Razor scooter. Borden and a friend grabbed the bike and scooter and tried to ride them down the street in the Fort Lauderdale suburb of Coral Springs.

Just as the 18-year-old girls were realizing they were too big for the tiny conveyances — which belonged to a 6-year-old boy — a woman came running after them saying, "That's my kid's stuff!" Borden and her friend immediately dropped the bike and scooter and walked away.

But it was too late — a neighbor who witnessed the heist had already called the police. Borden and her friend were arrested and charged with burglary and petty theft for the items, which were valued at a total of \$80.

COMPAS

- ▶ COMPAS is an ML algorithm used by some courts in the US to predict recidivism of condemned people
- ▶ A study showed that, given two people with the same features but different ethnicity, the system was giving higher probability of recidivism to non-white people

The system was biased against non-white people



Bernard Parker, left, was rated high risk; Dylan Proctor was rated low risk. [Drew Ritchie for ProPublica]

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

ON A SPRING AFTERNOON IN 2014, Brisha Borden was running late to pick up her god-sister from school when she spotted an unlocked kid's blue Huffy bicycle and a silver Razor scooter. Borden and a friend grabbed the bike and scooter and tried to ride them down the street in the Fort Lauderdale suburb of Coral Springs.

Just as the 18-year-old girls were realizing they were too big for the tiny conveyances — which belonged to a 6-year-old boy — a woman came running after them saying, "That's my kid's stuff!" Borden and her friend immediately dropped the bike and scooter and walked away.

But it was too late — a neighbor who witnessed the heist had already called the police. Borden and her friend were arrested and charged with burglary and petty theft for the items, which were valued at a total of \$80.

Let's Define Bias and Fairness

- **BIAS:** systematic favoritism or discrimination in models' predictions towards individuals based on some sensitive features (like *gender, race, and others*)
- **FAIRNESS:** absence of favoritism or discrimination in models' predictions



Is the concept of bias that simple?

Is the concept of bias that simple?

7

A Survey on Bias and Fairness in Machine Learning

115:5

- (1) **Measurement Bias.** *Measurement, or reporting, bias arises from how we choose, utilize, and measure particular features [140].* An example of this type of bias was observed in the recidivism risk prediction tool COMPAS, where prior arrests and friend/family arrests were used as proxy variables to measure level of “riskiness” or “crime”—which on its own can be viewed as mismeasured proxies. This is partly due to the fact that minority communities are controlled and policed more frequently, so they have higher arrest rates. However, one should not conclude that because people coming from minority groups have higher arrest rates, therefore they are more dangerous, as there is a difference in how these groups are assessed and controlled [140].
- (2) **Omitted Variable Bias.** *Omitted variable bias⁴ occurs when one or more important variables are left out of the model [38, 110, 127].* An example for this case would be when someone designs a model to predict, with relatively high accuracy, the annual percentage rate at which customers will stop subscribing to a service, but soon observes that the majority of users are canceling their subscription without receiving any warning from the designed model. Now imagine that the reason for canceling the subscriptions is appearance of a new strong competitor in the market that offers the same solution, but for half the price. The appearance of the competitor was something that the model was not ready for; therefore, it is considered to be an omitted variable.
- (3) **Representation Bias.** *Representation bias arises from how we sample from a population during data collection process [140].* Non-representative samples lack the diversity of the population, with missing subgroups and other anomalies. Lack of geographical diversity in datasets like ImageNet (as shown in Figures 3 and 4) results in demonstrable bias towards Western cultures.

Is the concept of bias that simple?

A Survey on Bias and Fairness in Machine Learning

115:5

- (1) **Measurement Bias.** *Measurement, or reporting, bias arises from how we choose, utilize, and measure particular features [140].* An example of this type of bias was observed in the recidivism risk prediction tool COMPAS, where prior arrests and friend/family arrests were used as proxy variables to measure level of “riskiness” or “crime”—which on its own can be

3.1.2 *Algorithm to User.* Algorithms modulate user behavior. Any biases in algorithms might introduce biases in user behavior. In this section, we talk about biases that are as a result of algorithmic outcomes and affect user behavior as a consequence.

- (1) **Algorithmic Bias.** *Algorithmic bias is when the bias is not present in the input data and is added purely by the algorithm [9].* The algorithmic design choices, such as use of certain optimization functions, regularizations, choices in applying regression models on the data as a whole or considering subgroups, and the general use of statistically biased estimators in algorithms [44], can all contribute to biased algorithmic decisions that can bias the outcome of the algorithms.
- (2) **User Interaction Bias.** *User Interaction bias is a type of bias that can not only be observable on the Web but also get triggered from two sources—the user interface and through the user itself by imposing his/her self-selected biased behavior and interaction [9].* This type of bias can be influenced by other types and subtypes, such as presentation and ranking biases.
- (a) **Presentation Bias.** *Presentation bias is a result of how information is presented [9].* For example, on the Web users can only click on content that they see, so the seen content gets clicks, while everything else gets no click. And it could be the case that the user does not see all the information on the Web [9].
- (b) **Ranking Bias.** *The idea that top-ranked results are the most relevant and important will result in attraction of more clicks than others.* This bias affects search engines [9] and crowdsourcing applications [92].
- (3) **Popularity Bias.** *Items that are more popular tend to be exposed more. However, popularity metrics are subject to manipulation—for example, by fake reviews or social bots [113].* As an

Is the concept of bias that simple?

A Survey on Bias and Fairness in Machine Learning

115:5

- (1) **Measurement Bias.** *Measurement bias refers to the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.*

3.1.2 *Algorithm to User.* *Algorithms can introduce biases in user behavior by influencing outcomes and affecting user interactions.*

- (1) **Algorithmic Bias.** *Algorithmic bias is introduced purely by the algorithmic optimization functions as a whole or considering individual algorithms [44], can all influence the output of the algorithms.*

- (2) **User Interaction Bias.** *Users interact with the Web but also get influenced by other types of bias.*

- (a) **Presentation Bias.** *For example, on the Web, a person gets clicks, while another person does not see all the content.*

- (b) **Ranking Bias.** *The ranking of items based on popularity result in attraction of more clicks than others. This bias affects search engines [9] and crowdsourcing applications [92].*

- (3) **Popularity Bias.** *Items that are more popular tend to be exposed more. However, popularity metrics are subject to manipulation—for example, by fake reviews or social bots [113]. As an*

to the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.

- (2) **Population Bias.** *Population bias arises when statistics, demographics, representatives, and user characteristics are different in the user population of the platform from the original target population [116]. Population bias creates non-representative data. An example of this type of bias can arise from different user demographics on different social platforms, such as women being more likely to use Pinterest, Facebook, Instagram, while men being more active in online forums like Reddit or Twitter. More such examples and statistics related to social media use among young adults according to gender, race, ethnicity, and parental educational background can be found in Reference [64].*

- (3) **Self-selection Bias.** *Self-selection bias⁴ is a subtype of the selection or sampling bias in which subjects of the research select themselves. An example of this type of bias can be observed in an opinion poll to measure enthusiasm for a political candidate, where the most enthusiastic supporters are more likely to complete the poll.*

- (4) **Social Bias.** *Social bias happens when others' actions affect our judgment [9]. An example of this type of bias can be a case where we want to rate or review an item with a low score, but when influenced by other high ratings, we change our scoring thinking that perhaps we are being too harsh [9, 147].*

- (5) **Behavioral Bias.** *Behavioral bias arises from different user behavior across platforms, contexts, or different datasets [116]. An example of this type of bias can be observed in Reference [104], where authors show how differences in emoji representations among platforms can result in*

result in attraction of more clicks than others. This bias affects search engines [9] and crowdsourcing applications [92].

- (3) Popularity Bias. Items that are more popular tend to be exposed more. However, popularity metrics are subject to manipulation—for example, by fake reviews or social bots [113]. As an*

Is the concept of bias that simple?

A Survey on Bias and Fairness in Machine Learning

115:5

- (1) **Measurement Bias.** *Measurement bias refers to the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.*

3.1.2 *Algorithm to User.* *Algorithms can introduce biases in user behavior by influencing the outcome of the algorithm.*

- (1) **Algorithmic Bias.** *Algorithms can add bias to the system by introducing optimization heuristics as a whole. This can lead to algorithmic bias.*
- (2) **User Interaction Bias.** *User interaction bias arises from the way users interact with the system by impacting their behavior and influence on the system.*
- (a) **Preference Bias.** *Preference bias arises from the way users prefer certain types of content over others. For example, if a user prefers one type of news over another, the algorithm will show them more of that type of news.*
- (b) **Reward Bias.** *Reward bias arises from the way users receive rewards for interacting with the system. For example, if a user receives a reward for sharing a post, they are more likely to share it again.*

- (3) **Population Bias.** *Population bias arises from the way different groups of people interact with the system. For example, if a certain demographic group is underrepresented in the user base, the algorithm may show them less relevant content.*

to the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.

- (2) **Population Bias.** *Population bias arises when statistics, demographics, representatives, and user characteristics are different in the user population of the platform from the original target population [116]. Population bias creates non-representative data. An example of this type of bias is when a platform's user base is predominantly male, but the target population is predominantly female.*

A Survey on Bias and Fairness in Machine Learning

115:9

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

- (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.

such as
re active
to social
cational

in which
erved in
nusastic

ample of
core, but
s we are

contexts,
ce [104],
result in

Is the concept of bias that simple?

A Survey on Bias and Fairness in Machine Learning

115:5

- (1) **Measurement Bias.** *Measurement bias refers to the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.*

3.1.2 *Algorithm to User.* *Algorithms can introduce biases in user behavior by introducing biases in user behavioral outcomes.*

- (1) **Algorithmic Bias.** *Algorithms can introduce biases in user behavior by adding heuristics to the optimization process, such as when the algorithm prioritizes certain users over others.*

- (2) **User Interaction Bias.** *User interaction bias arises when the way users interact with the system influences the way the system works.*

- (a) **Preference Bias.** *Preference bias arises when users get different recommendations based on their past interactions with the system.*

- (b) **Reward Bias.** *Reward bias arises when the system rewards certain types of user interactions over others.*

- (3) **Population Bias.** *Population bias arises when the user population of the platform from the original target population [116]. Population bias creates non-representative data. An example of this type of bias is the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.*

to the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.

- (2) **Population Bias.** *Population bias arises when statistics, demographics, representatives, and user characteristics are different in the user population of the platform from the original target population [116]. Population bias creates non-representative data. An example of this type of bias is the fact that only 5% of Fortune 500 CEOs were women—which would cause the search results to be biased towards male CEOs [140]. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering.*

A Survey on Bias and Fairness in Machine Learning

115:9

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

- (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.

such as
re active
to social
cational

in which
erved in
nusastic

ample of
core, but
s we are

contexts,
ce [104],
result in

**At least 23 different definitions
of bias in the literature**

From many definitions come many metrics...

From many definitions come many metrics...

Generic metrics

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

From many definitions come many metrics...

Generic metrics

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

Individual fairness metrics

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

From many definitions come many metrics...

Generic metrics

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

Individual fairness metrics

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

Group fairness metrics

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

From many definitions come many metrics...

Generic metrics

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

Individual fairness metrics

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

Group fairness metrics

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

At least 29 different metrics available in the AIF360 library

Mitigating Bias

Mitigating Bias

aif360.algorithms.preprocessing

```
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.Reweighting (...)
```

Reweighting is a preprocessing technique that Weights the examples in each (group, label) combination differently to ensure fairness before classification [4].

Mitigating Bias

aif360.algorithms.preprocessing

<code>algorithms.preprocessing.DisparateImpactRemover ([...])</code>	Disparate impact remover is feature values increase group ordering within groups [1]_.
<code>algorithms.preprocessing.LFR ([..., k, Ax, ...])</code>	Learning fair representations finds a latent representation obfuscates information about
<code>algorithms.preprocessing.OptimPreproc ([..., ...])</code>	Optimized preprocessing is a probabilistic transformation of the data with group fairness, fidelity constraints and objective functions.
<code>algorithms.preprocessing.Reweighting (...)</code>	Reweighting is a preprocessing examples in each (group, label) pair to achieve fairness before classification.

aif360.algorithms.inprocessing

<code>algorithms.inprocessing.AdversarialDebiasing (...)</code>	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]_.
<code>algorithms.inprocessing.ARTClassifier (...)</code>	Wraps an instance of an <code>art.classifiers.Classifier</code> to extend <code>Transformer</code> .
<code>algorithms.inprocessing.GerryFairClassifier ([...])</code>	Model is an algorithm for learning classifiers that are fair with respect to rich subgroups.
<code>algorithms.inprocessing.MetaFairClassifier ([...])</code>	The meta algorithm here takes the fairness metric as part of the input and returns a classifier optimized w.r.t.
<code>algorithms.inprocessing.PrejudiceRemover ([...])</code>	Prejudice remover is an in-processing technique that adds a discrimination-aware regularization term to the learning objective [6]_.
<code>algorithms.inprocessing.ExponentiatedGradientReduction (...)</code>	Exponentiated gradient reduction for fair classification.
<code>algorithms.inprocessing.GridSearchReduction (...)</code>	Grid search reduction for fair classification or regression.

Mitigating Bias

aif360.algorithms.preprocessing

<code>algorithms.preprocessing.DisparateImpactRemover ([...])</code>	Disparate impact remover is a feature selector that increases group ordering within groups [1].
<code>algorithms.preprocessing.LFR ([..., k, Ax, ...])</code>	Learning fair representations finds a latent representation that obfuscates information about protected attributes.
<code>algorithms.preprocessing.OptimPreproc ([..., ...])</code>	Optimized preprocessing is a probabilistic transformation of the data with group fairness, fidelity constraints and objective functions.

aif360.algorithms.inprocessing

<code>algorithms.inprocessing.AdversarialDebiasing (...)</code>	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].
<code>algorithms.inprocessing.ARTClassifier (...)</code>	Wraps an instance of an <code>art.classifiers.Classifier</code> to extend <code>Transformer</code> .
<code>algorithms.inprocessing.GerryFairClassifier ([...])</code>	Model is an algorithm for learning classifiers that are fair with respect to rich subgroups.

aif360.algorithms.postprocessing

<code>algorithms.postprocessing.CalibratedEqOddsPostprocessing (...)</code>	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].
<code>algorithms.postprocessing.EqOddsPostprocessing (...)</code>	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].
<code>algorithms.postprocessing.RejectOptionClassification (...)</code>	Reject option classification is a postprocessing technique that gives favorable outcomes to unprivileged groups and unfavorable outcomes to privileged groups in a confidence band around the decision boundary with the highest uncertainty [10].

<code>Optimizer ([...])</code>	The meta algorithm here takes the fairness metric as part of the input and returns a classifier optimized w.r.t.
<code>PrejudiceRemover ([...])</code>	Prejudice remover is an in-processing technique that adds a discrimination-aware regularization term to the learning objective [6].
<code>GradientReduction (...)</code>	Exponentiated gradient reduction for fair classification.
<code>GridSearch (...)</code>	Grid search reduction for fair classification or regression.

Mitigating Bias

aif360.algorithms.preprocessing

```
algorithms.preprocessing.DisparateImpactRemover (...)
```

Disparate impact remover is a feature that increases group ordering within groups [1].

```
algorithms.preprocessing.LFR (...[, k, Ax, ...])
```

Learning fair representations finds a latent representation that obfuscates information about

```
algorithms.preprocessing.OptimPreproc (...[, ...])
```

Optimized preprocessing is a probabilistic transformation

aif360.algorithms.postprocessing

```
algorithms.postprocessing.CalibratedEqOddsPostprocessing (...)
```

processing technique that optimizes over calibrated classifier score outputs to find probabilities with which to change output labels with an equalized odds objective [7].

```
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].

```
algorithms.postprocessing.RejectOptionClassification (...)
```

Reject option classification is a postprocessing technique that gives favorable outcomes to unprivileged groups and unfavorable outcomes to privileged 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 `art.classifiers.Classifier` to extend `Transformer`.

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

```
ver (...)
```

discrimination aware regularization term to the learning objective [6].

```
GradientReduction (...)
```

Exponentiated gradient reduction for fair classification.

```
duction (...)
```

Grid search reduction for fair classification or regression.

Our Contributions

Our Contributions

Challenge 1 (CH1)

Developing approaches for bias mitigation both in binary and multi-class classification settings.

Challenge 2 (CH2)

Democratizing the development of fair learning-based systems to actors with different expertise.

Challenge 1: Bias in Multi-Class Classification

24

- Most of the bias mitigation approaches focus on binary classification
- However, many multi-class classification approaches have been proposed in sensitive domains

Computing, Artificial Intelligence and Information Technology

A data-driven software tool for enabling cooperative information sharing among police departments

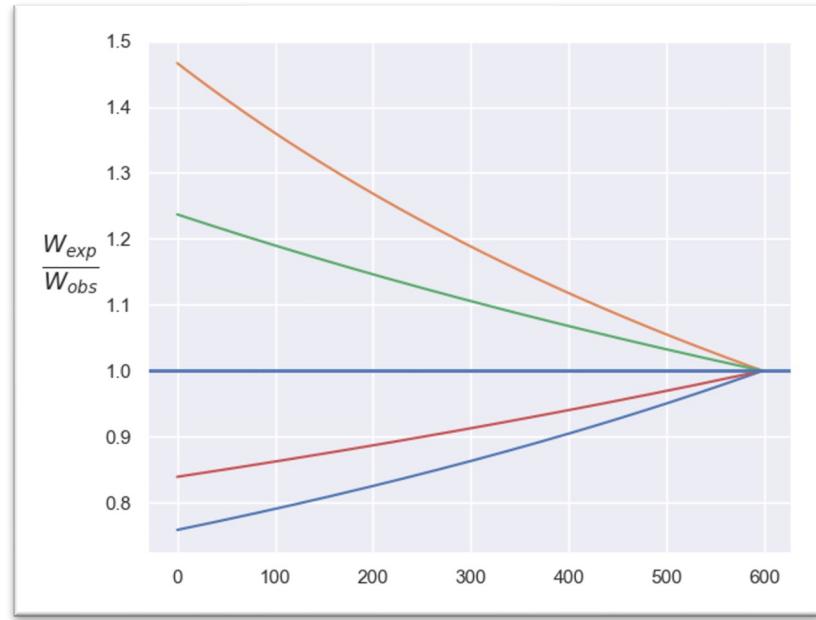
Will I Pass the Bar Exam: Predicting Student Success Using LSAT Scores and Law School Performance

Nuclear feature extraction for breast tumor diagnosis

Contribution 1: Debiaser for Multiple Variables

25

- ▶ DEMV is a pre-processing approach to improve fairness in binary and multi-class classification tasks
- ▶ Overcomes all the other state-of-the-art multi-class bias mitigation algorithms in the literature
- ▶ Algorithm available on SoBigData RI and PIP



Challenge 2: Democratising Software Fairness

Challenge 2: Democratising Software Fairness

27

23 Definitions of
Bias

Challenge 2: Democratising Software Fairness

23 Definitions of
Bias



29 Different Metric

Challenge 2: Democratising Software Fairness

23 Definitions of
Bias



29 Different Metric



14 Different
Methods

Challenge 2: Democratising Software Fairness

30

23 Definitions of
Bias



29 Different Metric



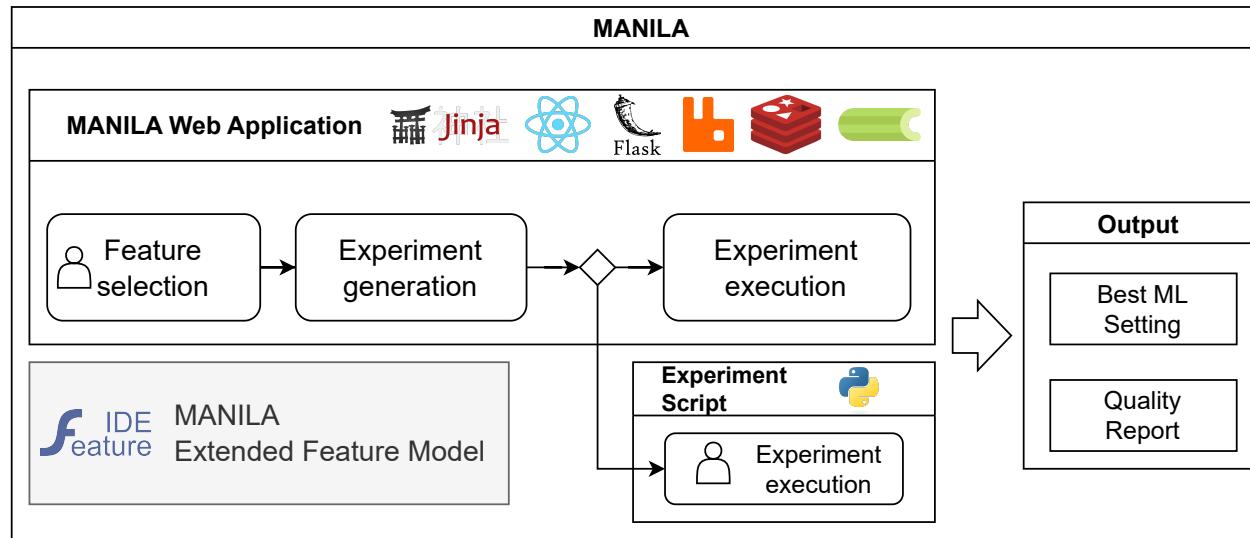
14 Different
Methods



Data Scientist less-expert
on fairness

Contribution 2: MANILA

- We propose MANILA, a web-based application to design, implement and execute fairness evaluations
- Uses the Extended Feature Model (ExtFM) formalism to model the evaluation workflow as a Software Product Line



Available in
SoBigData RI





Thank you for your attention!

UNIVERSITÀ
DEGLI STUDI
DELL'AQUILA



DISIM
Dipartimento di Ingegneria
e Scienze dell'Informazione
e Matematica