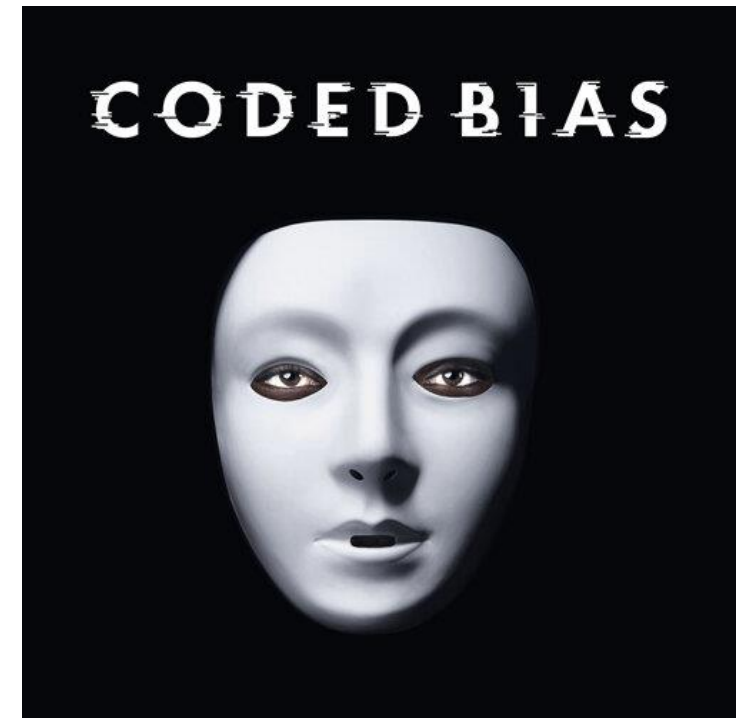# Approaches to Measure and Mitigate *Algorithmic Bias*

Giordano d'Aloisio

Università degli Studi dell'Aquila e SoBigData RI

**PARTNER**

# Coded Bias

- In 2018 an MIT researcher was studying Amazon's face recognition systems

- Those systems were not able to recognize her face

- At first, she thought it was an error in the system

- But then, she noticed that wearing a white mask the system was able to recognize her
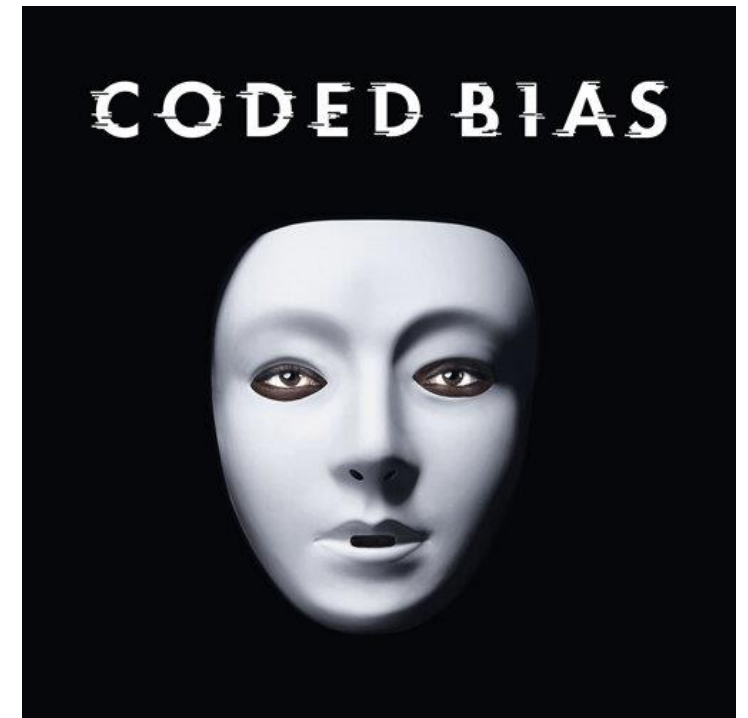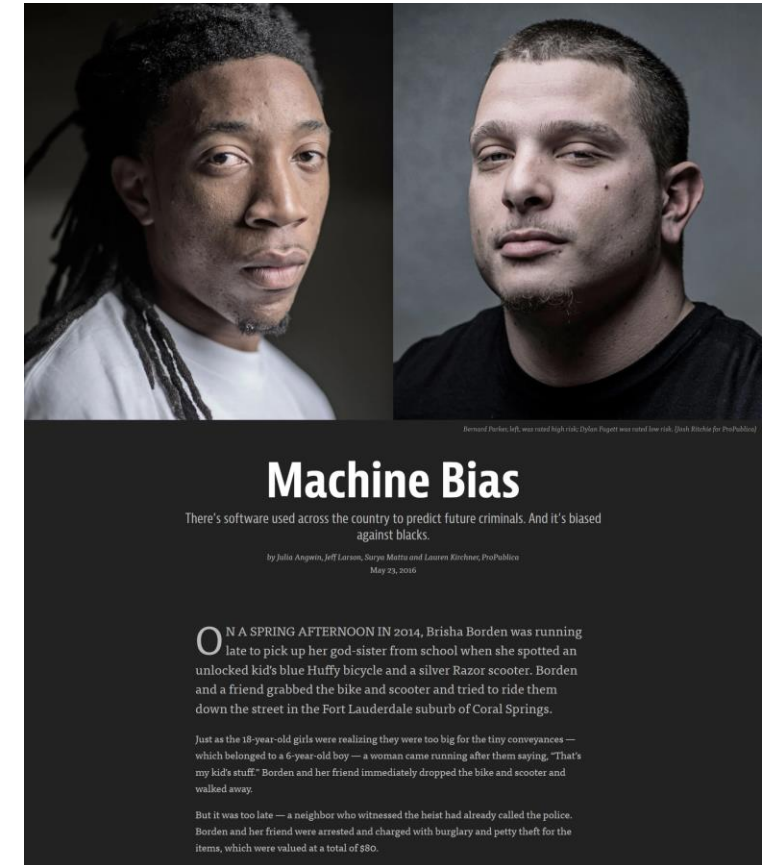
# Coded Bias

- In 2018 an MIT researcher was studying Amazon's face recognition systems

- Those systems were not able to recognize her face

- At first, she thought it was an error in the system

- But then, she noticed that wearing a white mask the system was able to recognize her

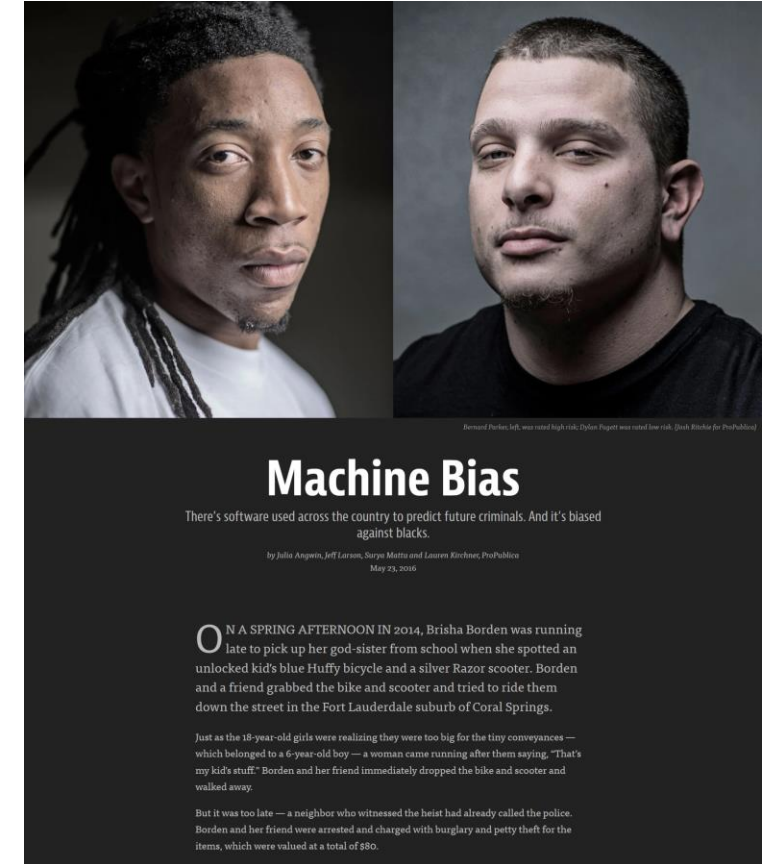The system was biased against non-white women

# Another Example

- 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 race, the system was giving higher probability of recidivism to non-white people



**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.

# Another Example

- 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 race, 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 Fugett was rated low risk. (Josh 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 better Bias and Fairness

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

- **FAIRNESS:** absence of favouritism or discrimination in models' predictions

*N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, 'A Survey on Bias and Fairness in Machine Learning', ACM Comput. Surv., vol. 54, no. 6, pp. 1–35, Jul. 2021, doi: 10.1145/3457607.*

# Is the concept of bias that simple?

- Actually not…

# Is the concept of bias that simple?

- Actually not…

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 bias[4] 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?

- Actually not…



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

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 observant 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?

- Actually not…

A Survey on Bias and Fairness in Machine Learning

(1) **Measurement Bias.** *Measurement, or reporting, bias arises from how* *measure particular features* [140]. An example of this type of bias cidivism risk prediction tool COMPAS, where prior arrests and frie used as proxy variables to measure level of "riskiness" or "crime"—w viewed as mismeasured proxies. This is partly due to the fact that are controlled and policed more frequently, so they have higher arre

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

(1) **Algorithmic Bias.** *Algorithmic bias is when the bias is not present in th* *added purely by the algorithm* [9]. The algorithmic design choices, suc optimization functions, regularizations, choices in applying regression as a whole or considering subgroups, and the general use of statistically b algorithms [44], can all contribute to biased algorithmic decisions that ca of the algorithms.

(2) **User Interaction Bias.** *User Interaction bias is a type of bias that can not* *the Web but also get triggered from two sources—the user interface and th* *by imposing his/her self-selected biased behavior and interaction* [9]. This 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

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 bias[4] 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

# Is the concept of bias that simple?

- Actually not…

A Survey on Bias and Fairness in Machine Learning

(1) **Measurement Bias.** *Measurement, or reporting, bias arises from hov measure particular features* [140]. An example of this type of bias v cidivism risk prediction tool COMPAS, where prior arrests and frie used as proxy variables to measure level of "riskiness" or "crime"—w viewed as mismeasured proxies. This is partly due to the fact that are controlled and policed more frequently, so they have higher arr

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

(1) **Algorithmic Bias.** *Algorithmic bias is when the bias is not present in th added purely by the algorithm* [9]. The algorithmic design choices, suc optimization functions, regularizations, choices in applying regression as a whole or considering subgroups, and the general use of statistically b algorithms [44], can all contribute to biased algorithmic decisions that ca of the algorithms.

(2) **User Interaction Bias.** *User Interaction bias is a type of bias that can not the Web but also get triggered from two sources—the user interface and th by imposing his/her self-selected biased behavior and interaction* [9]. This influenced by other types and subtypes, such as presentation and ranking

  (a) **Presentation Bias.** *Presentation bias is a result of how information is example, on the Web users can only click on content that they see, so gets clicks, while everything else gets no click. And it could be the c does not see all the information on the Web* [9].

  (b) **Ranking Bias.** *The idea that top-ranked results are the most relevant a result in attraction of more clicks than others. This bias affects search crowdsourcing applications* [92].

(3) **Popularity Bias.** *Items that are more popular tend to be exposed more. Ho metrics are subject to manipulation—for example, by fake reviews or social

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

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.

# Is the concept of bias that simple?

- Actually not…

A Survey on Bias and Fairness in Machine Learning

(1) **Measurement Bias.** *Measurement, or reporting, bias arises from how* *measure particular features* [140]. An example of this type of bias v cidivism risk prediction tool COMPAS, where prior arrests and frie used as proxy variables to measure level of "riskiness" or "crime"—w viewed as mismeasured proxies. This is partly due to the fact that are controlled and policed more frequently, so they have higher arre

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

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

(1) **Algorithmic Bias.** *Algorithmic bias is when the bias is not present in th* *added purely by the algorithm* [9]. The algorithmic design choices, suc optimization functions, regularizations, choices in applying regression as a whole or considering subgroups, and the general use of statistically b algorithms [44], can all contribute to biased algorithmic decisions that ca of the algorithms.

(2) **User Interaction Bias.** *User Interaction bias is a type of bias that can not* *the Web but also get triggered from two sources—the user interface and thi* *by imposing his/her self-selected biased behavior and interaction* [9]. This influenced by other types and subtypes, such as presentation and ranking

   (a) **Presentation Bias.** *Presentation bias is a result of how information is* example, on the Web users can only click on content that they see, so gets clicks, while everything else gets no click. And it could be the c does not see all the information on the Web [9].

   (b) **Ranking Bias.** *The idea that top-ranked results are the most relevant a* *result in attraction of more clicks than others.* This bias affects search crowdsourcing applications [92].

(3) **Popularity Bias.** *Items that are more popular tend to be exposed more. Hc* *metrics are subject to manipulation—for example, by fake reviews or social*

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.

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

# From many definitions come many metrics

# From many definitions come many metrics

## Generic metrics

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

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

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

# From many definitions come many metrics



## Generic metrics

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

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

## Group fairness metrics

| | |
|---|---|
| metrics.statistical_parity_difference (y_true) | Difference in selection rates. |
| metrics.mean_difference (y_true[, y_pred, ...]) | Alias of `statistical_parity_difference()`. |
| metrics.disparate_impact_ratio (y_true[, ...]) | Ratio of selection rates. |
| 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. |
| metrics.average_odds_error (y_true, y_pred, *) | A relaxed version of equality of odds. |
| metrics.class_imbalance (y_true[, y_pred, ...]) | 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, $KL(P_p||P_u) = \sum_y P_p(y) \log\left(\frac{P_p(y)}{P_u(y)}\right)$ |
| metrics.conditional_demographic_disparity (y_true) | Conditional demographic disparity, $CDD = \frac{1}{\sum_i N_i} \sum_i N_i \cdot DD_i$ |
| metrics.smoothed_edf (y_true[, y_pred, ...]) | 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. |
| metrics.mdss_bias_scan (y_true, probas_pred) | DEPRECATED: Change to new interface - aif360.sklearn.detectors.mdss_detector.bias_scan by version 0.5.0. |
| metrics.mdss_bias_score (y_true, probas_pred) | Compute the bias score for a prespecified group of records using a given scoring function. |

# From many definitions come many metrics

## Generic metrics

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

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

## Group fairness metrics

| | |
|---|---|
| metrics.statistical_parity_difference (y_true) | Difference in selection rates. |
| metrics.mean_difference (y_true[, y_pred, ...]) | Alias of statistical_parity_difference() . |
| metrics.disparate_impact_ratio (y_true[, ...]) | Ratio of selection rates. |
| 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. |
| metrics.average_odds_error (y_true, y_pred, *) | A relaxed version of equality of odds. |
| metrics.class_imbalance (y_true[, y_pred, ...]) | 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, $KL(P_p\|\|P_u) = \sum_y P_p(y) \log\left(\frac{P_p(y)}{P_u(y)}\right)$ |
| metrics.conditional_demographic_disparity (y_true) | Conditional demographic disparity, $CDD = \frac{1}{\sum_i N_i} \sum_i N_i \cdot DD_i$ |
| metrics.smoothed_edf (y_true[, y_pred, ...]) | 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. |
| metrics.mdss_bias_scan (y_true, probas_pred) | DEPRECATED: Change to new interface - aif360.sklearn.detectors.mdss_detector.bias_scan by version 0.5.0. |
| metrics.mdss_bias_score (y_true, probas_pred) | Compute the bias score for a prespecified group of records using a given scoring function. |

**At least 29 different metrics are available in the AIF360 library**

# Addressing the bias issue

# Addressing the bias issue

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

# Addressing the bias issue



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

**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`. |
| algorithms.inprocessing.GerryFairClassifier ([...]) | Model is an algorithm for learning classifiers that are fair with respect to rich subgroups. |
| 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]_. |
| algorithms.inprocessing.ExponentiatedGradientReduction (...) | Exponentiated gradient reduction for fair classification. |
| algorithms.inprocessing.GridSearchReduction (...) | Grid search reduction for fair classification or regression. |

# Addressing the bias issue

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

## 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`. |
| algorithms.inprocessing.GerryFairClassifier ([...]) | Model is an algorithm for learning classifiers that are fair with respect to rich subgroups. |
| 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]_. |
| algorithms.inprocessing.ExponentiatedGradientReduction (...) | Exponentiated gradient reduction for fair classification. |
| algorithms.inprocessing.GridSearchReduction (...) | Grid search reduction for fair classification or regression. |

## aif360.algorithms.postprocessing

| | |
|---|---|
| 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]_. |
| 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 unpriviliged groups and unfavorable outcomes to priviliged groups in a confidence band around the decision boundary with the highest uncertainty [10]_. |

# Addressing the bias issue

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

## 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`. |
| algorithms.inprocessing.GerryFairClassifier ([...]) | Model is an algorithm for learning classifiers that are fair with respect to rich subgroups. |
| 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]_. |
| algorithms.inprocessing.ExponentiatedGradientReduction (...) | Exponentiated gradient reduction for fair classification. |
| algorithms.inprocessing.GridSearchReduction (...) | Grid search reduction for fair classification or regression. |

## aif360.algorithms.postprocessing

| | |
|---|---|
| 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]_. |
| 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 unpriviliged groups and unfavorable outcomes to priviliged groups in a confidence band around the decision boundary with the highest uncertainty [10]_. |

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

# What is missing?

# What is missing?

- **Challenge 1:** Most bias mitigation methods address binary classification. What about multi-class classification?

# What is missing?

- **Challenge 1:** Most bias mitigation methods address binary classification. What about multi-class classification?

- **Challenge 2:** Plenty of bias definitions, metrics and methods are available. How can we guide non-expert users in developing fair ML systems?

# What is missing?

- **Challenge 1:** Most bias mitigation methods address binary classification. What about multi-class classification?

- **Challenge 2:** Plenty of bias definitions, metrics and methods are available. How can we guide non-expert users in developing fair ML systems?

- **Challenge 3:** Most fairness assessment approaches are domain and definition specific. How can we address non-traditional use cases?

# Challenge 1

- Most of the bias mitigation approaches focus on binary classification
- However, bias is a relevant issue in many multi-class problems

Computing, Artificial Intelligence and Information Technology

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

Michael Redmond [a], Alok Baveja [b]

## Will I Pass the Bar Exam: Predicting Student Success Using LSAT Scores and Law School Performance
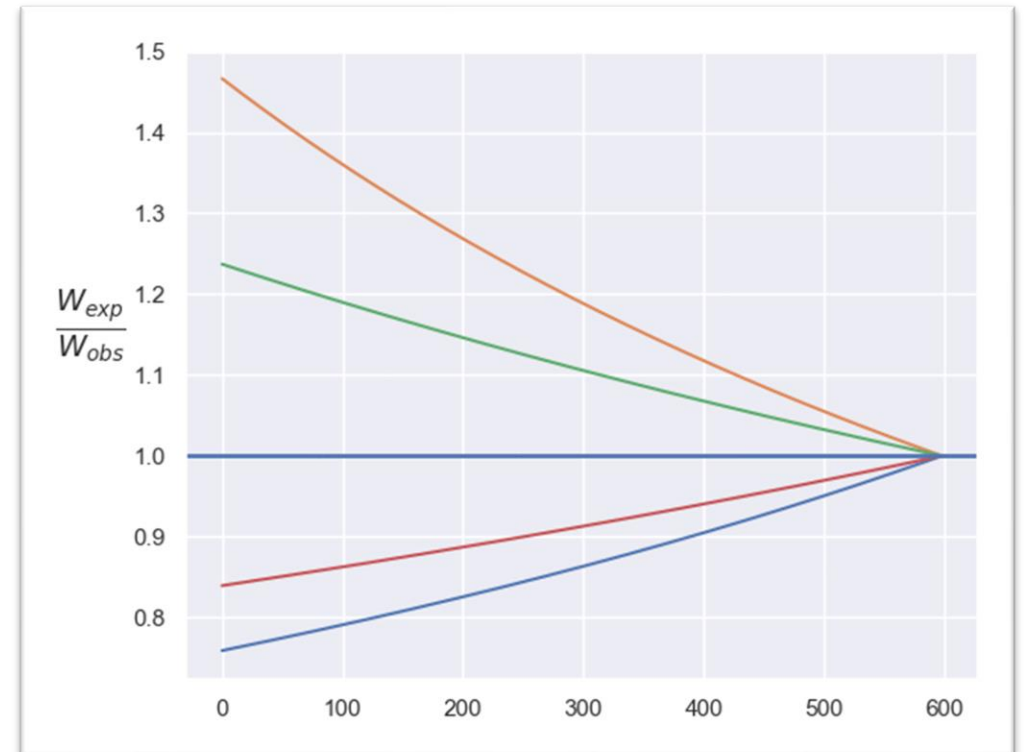
Katherine A. Austin

Catherine Martin Christopher

Darby Dickerson

## Nuclear feature extraction for breast tumor diagnosis

W. Nick Street, W. H. Wolberg, O. L. Mangasarian

# Addressing challenge 1

- To address this challenge, we propose the Debiaser for Multiple Variables (DEMV) [1-2]

- Novel algorithm to improve fairness in binary and multi-class classification problems

- Works by perfectly rebalance the dataset's sensitive groups

- Overcomes all the other multi-class debiaser algorithms in the literature

# Challenge 2

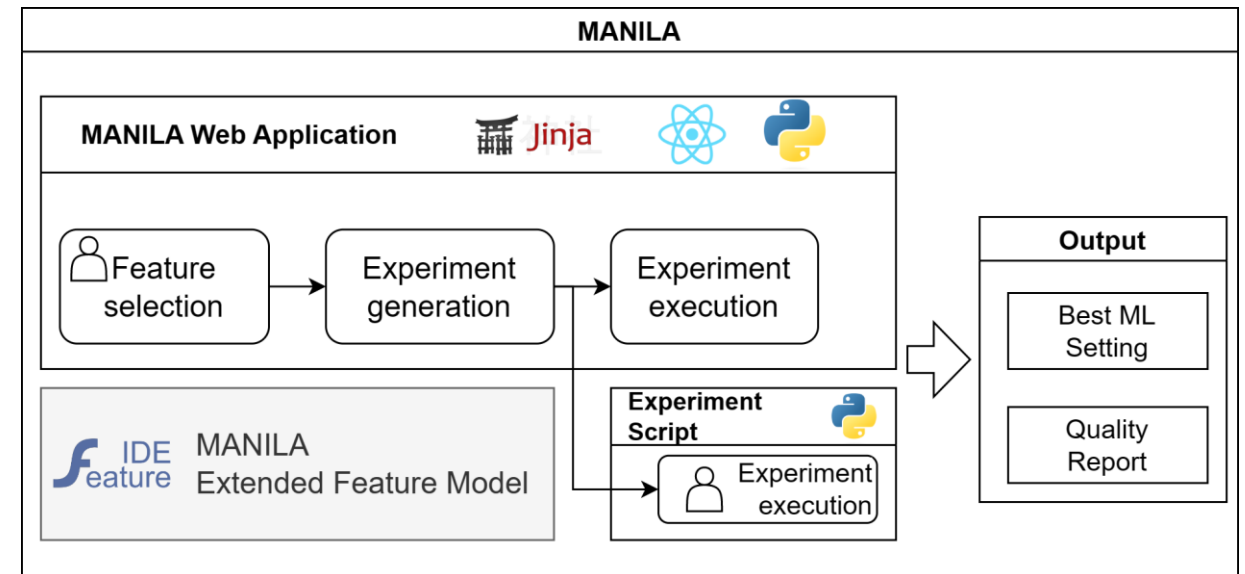| 23 different definitions of bias | **+** | 29 different metrics | **+** | 14 different methods |
|---|---|---|---|---|

- This can be a challenge for users that are non-expert on fairness

- Software engineering approaches can help us to formalize and standardize the development of fair ML systems

- Hence, make the development easier even for non-expert users

# Addressing challenge 2

- To this aim we propose MANILA [3-4]

- A web application that guides users in defining and performing fairness and effectiveness evaluations

- Available as an application in the SoBigData RI
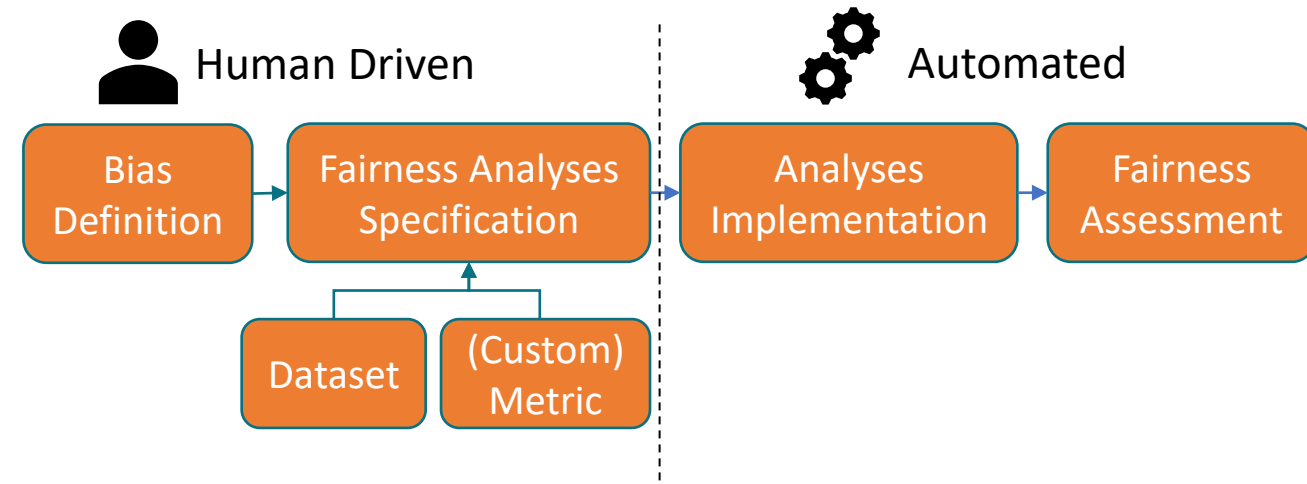
# Challenge 3

- Most of the fairness assessment tools available focus on specific definitions of fairness or cover traditional use cases

- What about non-traditional use cases (e.g., IoT?)

Co-zyBench: Using Co-Simulation and Digital Twins to Benchmark Thermal Comfort Provision in Smart Buildings

ResyDuo: Combining Data Models and CF-Based Recommender Systems to Develop Arduino Projects

# Addressing challenge 3

- We propose MODNESS [5], a model-driven framework to conceptualize, design, implement, and execute fairness assessment analyses

- Allows to define different concrete fairness analyses starting from a single high-level bias definition

- Allows to model custom metrics for fairness assessment

# Many challenges are still open

- Addressing the trade-off between fairness and other quality properties (e.g., privacy, computational complexity,…)

- Early identify features leading to bias in a dataset

- Suggest to the user the best bias definition and metric starting for a specific requirement

- Formally model a fairness specification

- And many more…

# References

[1] G. d'Aloisio, G. Stilo, A. Di Marco, e A. D'Angelo, «Enhancing Fairness in Classification Tasks with Multiple Variables: A Data-and Model-Agnostic Approach», in International Workshop on Algorithmic Bias in Search and Recommendation, Springer, 2022, pp. 117–129.

[2] G. d'Aloisio, A. D'Angelo, A. Di Marco, e G. Stilo, «Debiaser for Multiple Variables to enhance fairness in classification tasks», Information Processing & Management, vol. 60, fasc. 2, p. 103226, 2023.

[3] G. d'Aloisio, A. Di Marco, G. Stilo, «A Framework to Democratize the Quality-Based Machine Learning Development Through Extended Feature Models: Strengths and Limitations», Science of Computer Programming (Under review)

[4] G. d'Aloisio, A. Di Marco, e G. Stilo, «Democratizing Quality-Based Machine Learning Development through Extended Feature Models», in International Conference on Fundamental Approaches to Software Engineering, Springer Nature Switzerland Cham, 2023, pp. 88–110.

[5] G. d'Aloisio, C. Di Sipio, A. Di Marco, D. Di Ruscio, «How fair are we? From conceptualization to automated assessment of fairness definitions», Software and Systems Modeling (Under review)

# Thank you for your attention!