

# Algorithmic fairness in computational medicine

Jie Xu,<sup>a,b</sup> Yunyu Xiao,<sup>b</sup> Wendy Hui Wang,<sup>c</sup> Yue Ning,<sup>c</sup> Elizabeth A. Shenkman,<sup>a</sup>  
Jiang Bian,<sup>a</sup> and Fei Wang<sup>b,\*</sup>

<sup>a</sup>Department of Health Outcomes and Biomedical Informatics, University of Florida, Gainesville, FL, USA

<sup>b</sup>Department of Population Health Sciences, Weill Cornell Medicine, New York, NY, USA

<sup>c</sup>Department of Computer Science, Stevens Institute of Technology, Hoboken, NJ, USA



## Summary

Machine learning models are increasingly adopted for facilitating clinical decision-making. However, recent research has shown that machine learning techniques may result in potential biases when making decisions for people in different subgroups, which can lead to detrimental effects on the health and well-being of specific demographic groups such as vulnerable ethnic minorities. This problem, termed algorithmic bias, has been extensively studied in theoretical machine learning recently. However, the impact of algorithmic bias on medicine and methods to mitigate this bias remain topics of active discussion. This paper presents a comprehensive review of algorithmic fairness in the context of computational medicine, which aims at improving medicine with computational approaches. Specifically, we overview the different types of algorithmic bias, fairness quantification metrics, and bias mitigation methods, and summarize popular software libraries and tools for bias evaluation and mitigation, with the goal of providing reference and insights to researchers and practitioners in computational medicine.

**eBioMedicine 2022;84:104250**

Published online 6 September 2022

<https://doi.org/10.1016/j.ebiom.2022.104250>

**Copyright** © 2022 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>)

**Keywords:** Algorithmic fairness; Computational medicine

## Introduction

Recent years have witnessed a surge of interests in the development and deployment of machine learning algorithms in healthcare. These algorithms were trained on massive health data and have demonstrated promising performance in a diverse set of problems such as skin cancer detection from lesion images,<sup>1</sup> prediction of the risk of acute kidney injury based on electronic health records (EHR),<sup>2</sup> adaptive learning of the optimal treatment regimens for sepsis patients in intensive care<sup>3</sup> and others.<sup>4</sup>

Despite the promise, however, there is growing concern that machine learning algorithms may lead to unintended bias when making decisions involving ethnic minorities, both through the algorithms themselves and the data used to learn them. For example, associations between Framingham risk factors and cardiovascular events have been shown to be significantly different across different ethnic groups.<sup>5</sup> Video stream analysis algorithms for measuring the body's spontaneous blink rate have been found to be particularly challenging for Asian individuals.<sup>6,7</sup> Undiagnosed silent hypoxemia, detected from pulse oximetry, occurred approximately three times more frequently in Black people due to the fact that dark skin responds differently

to those light wavelengths.<sup>8</sup> In these cases, the software system may introduce or exacerbate health equity issues.<sup>7</sup>

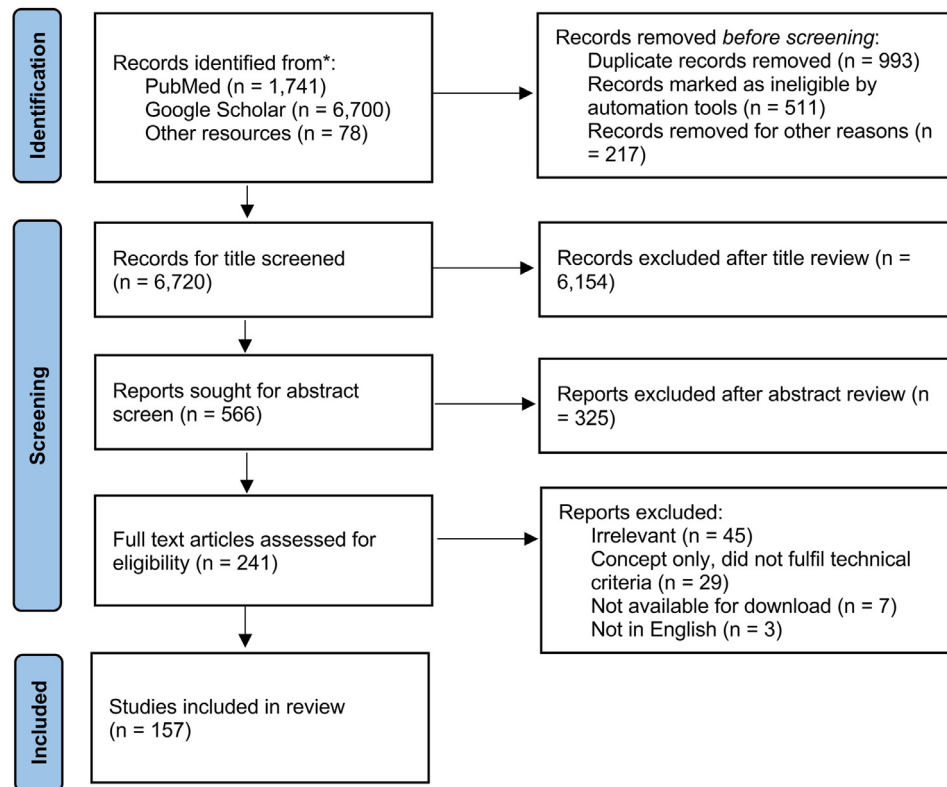
With machine learning models gaining increasing attention in medicine, it is crucial to be aware of the potential related bias and disparities, understand their causes, and methods to mitigate them. This review will help achieve this goal by providing an overview of the existing literature studying the sources of bias and disparities in computational medicine, their quantification metrics, and mitigation strategies. We will also summarize outstanding questions and point out future directions. The PRISMA diagram of the literature reviewed in this paper is shown in [Figure 1](#).

## Distinguishing from existing reviews

It is worth noting other reviews of AI Fairness in the literature and how they differ. Mehrabi et al.<sup>9</sup> built a taxonomy of machine learning related fairness in different real world application contexts. Rajkomar et al.<sup>10</sup> introduced the principles of distributive justice and provided guidance to clinicians on how to prioritize each principle when facing with potential bias in model development and deployment. Gianfrancescogian et al.<sup>11</sup> summarized the potential bias sources for electronic health records (EHRs) and provided recommendations on appropriately mitigating them. Fletcher et al.<sup>12</sup> described three basic criteria (i.e., Appropriateness, Fairness, and Bias) for evaluating machine learning and

\*Corresponding author.

E-mail address: [few2001@med.cornell.edu](mailto:few2001@med.cornell.edu) (F. Wang).



**Figure 1.** PRISMA flow diagram: disparity and fairness in computational medicine.

AI systems in the context of global health. Mhasawade et al.<sup>13</sup> focused on the interactions among different cultural, social, and environmental factors, their impact on the fairness of machine learning algorithms, and how machine learning, public and population health can work together to achieve health equity. Unlike these existing works, this review summarizes sources and quantification methods for bias in computational medicine and how they will impact downstream machine learning models, as well as potential strategies to mitigate them through computational algorithms.

### Computational bias

We categorize computational biases into three different types according to the source of bias: data bias, measurement bias, and algorithm bias. We will introduce them in this section and provide examples in medical context.

#### Data bias

In supervised learning, machine learning algorithms are trained from data sets.<sup>14</sup> For example, classification models try to accurately map the sample input features to a set of pre-specified classes based on the observations from a set of training data. Clustering models aim at identifying grouping structures of a given data set. In

this case, if the data from a specific demographic group is not properly represented, the machine learning models trained from the data will be biased.

As a simple practical example, studies have found that patients of low socioeconomic status may have limited access to health care.<sup>15,16</sup> Consequently, compared to patients with higher socioeconomic status, these patients may generate proportionately less data in their electronic health records which will lead to underrepresentation if a machine learning model is trained using this data. This will lead to poorer model performance on this particular patient group. Below we list potential sources of data bias in medicine.

**Sampling bias.** Sampling bias, also known as selection bias, occurs when the selected data does not represent the real environment in which a model will be deployed.<sup>17</sup> For example, melanoma detection algorithms based on classification of skin lesion images<sup>1</sup> may perform poorly on dark-pigmented skin if the training images contain predominantly lighter skin.<sup>18</sup> For the same reason, Face2Gene, a machine learning algorithm to recognize Down syndrome based on facial images, performed much better in Caucasian (accuracy 80%) than in African (accuracy 36.8%).<sup>19</sup>

**Allocation bias.** Allocation bias pertains to clinical trials that contain interventions and arises if there are systematic differences in how participants are allocated to the treatment and control groups.<sup>20</sup> If researchers knew which participants would benefit from an intervention, it could bias how they recruit participants and how they assign them to different groups so that they can select subjects with a good prognosis for trials. Recently there were studies trying to emulate clinical trials with real world data such as EHRs.<sup>21</sup> In this case, allocation bias could exist as the treatment and control groups are already observed in the data. This can lead to potentially biased estimations of treatment effects with machine learning models.

**Attrition bias.** Attrition bias also applies to clinical trials and can occur if there are systematic differences in the way different groups of participants are recruited or are dropped from a study. When exploring an intervention, different rates of losses to follow-up in the exposure groups may alter the demographic composition of these groups.<sup>20</sup> Attrition bias will be more severe in observational studies, as patients may move to another place or be transferred to another hospital, which will impact the machine learning model looking to predict clinical events.

**Publication bias.** Publication bias occurs when the decision to publish a study depends on its own results.<sup>22</sup> Empirical studies consistently show that studies with positive or statistically significant results are easier and take less time to be published than studies without significant results.<sup>23,24</sup> This can make it difficult for decision makers to distinguish between sound evidence and overestimate the effectiveness of specific treatments or models.<sup>24</sup> For example, since the start of the COVID-19 pandemic, studies on COVID-19 is being published at a rapid rate. However, many peer-reviewed publications included only a limited number of patients included and showed a high risk of bias.<sup>25</sup>

### Measurement bias

Measurement bias is a systematic error that occurs when the data are labeled inconsistently, or study variables (e.g., disease, exposure) are collected or measured inaccurately.<sup>26</sup> A recent example is the large disparity in the quality of COVID-19 data reported across India.<sup>27</sup>

One of the common causes of measurement bias is response bias. In the clinical context, response bias usually occurs in studies involving surveys or self-reported data. When respondents tend to give inaccurate or even wrong answers on self-reported questions, the survey results will be affected.<sup>28</sup> An example of response bias is that people might tend to always rate themselves favorably or feel pressured to provide socially acceptable

answers.<sup>29</sup> In addition, misleading questions can lead to biased answers. In addition, demographic groups who are willing to answer survey questions are sometimes different from those who are not.<sup>30</sup> Consequently, this will impact the machine learning algorithms trained on surveys or patient reported outcomes.

### Algorithm bias

Another source of bias is from the algorithms themselves,<sup>31</sup> which can be algorithm specific or agnostic. Algorithm specific bias is linked to their intrinsic hypotheses.<sup>32</sup> For example, logistic regression models assume the relationships between input and target variables are linear, but this may not be true. Such a bias presents a challenge in capturing the actual input-output relationships in the data. The loss function measures the difference between the algorithm, output and the ground truth outcome. It is used to evaluate how well the machine learning algorithm fits the data. Typical machine learning algorithms attempt to minimize such prediction loss on the training data, which is typically measured by adding up all prediction losses on individual samples. However, if the loss function is biased towards a specific demographic group (e.g., white patients in a population),<sup>33</sup> the corresponding model will be better trained for this group.

### Fairness metrics

The previous section has summarized the various potential sources of computational bias. Another important question is how we can quantify such bias given a specific healthcare context or data set. In this section, we will review different ways that bias could be evaluated, which are referred to as fairness metrics. Mathematical notations that are used in this section are summarized in [Table 1](#).

To illustrate the use of fairness metrics, we make use of a case study to build an alerting algorithm in ICU setting (e.g., for developing sepsis<sup>34</sup>) with the machine learning algorithm based on the patient's EHR, and the patient's race. For the purpose of illustration, we consider only two demographic groups (e.g., Black or white)

Symbol	Description
$A \in \{0, 1\}$	Binary protected attribute
$X \in \mathbb{R}^d$	Other observable attributes
$U$	Relevant latent attributes not observed
$Y \in \{0, 1\}$	The outcome to be predicted
$\hat{Y} := f(X, A) \in \{0, 1\}$	The prediction of $Y$
$\hat{Y}_{A \leftarrow a}$	Counterfactual value, i.e., what would $\hat{Y}$ have been if $A$ had been equal to $a$

**Table 1: Notations and symbols.**

and we examine how the alerting algorithm can behave differently for patients from two demographic groups using various fairness metrics.<sup>35,36</sup>

## Fairness through unawareness

The simplest approach to achieve some degree of fairness is to remove the protected attribute (e.g., race in our case study) as an independent variable in the model.<sup>40–42</sup> This method has been shown to be ineffective because these protected attributes are often highly correlated to other parameters in the data set. For example, race may be related to zip code, socioeconomic status, or disease predisposition. Therefore, simply removing protected attribute is not enough to eliminate disparate results between the two demographic groups.

## Demographic parity

Another definition of Fairness is demographic parity, also known as statistical parity or independence, which requires that the overall proportion of individuals in a protected group predicted as positive (or negative) to be the same as that of the overall population.<sup>38</sup> Although it is intuitive to understand, prior studies<sup>43</sup> found that optimizing demographic parity may prevent the model from taking into account relevant clinical characteristics related to protected variables and outcomes, thereby reducing the performance of the model for all groups.

## Equalized odds

Unlike demographic parity, equalized odds is a definition of Fairness that<sup>39</sup> allows the prediction  $\hat{Y}$  to depend on protected attribute  $A$ , but only through the target variable  $Y$ . This encourages the use of features that are directly related to  $Y$ , rather than through  $A$ .<sup>39</sup> To achieve equalized odds, both true positive rates (TPR) and true negative rates (TNR) of all groups defined by  $A$  are equal up to a fixed tolerance  $T$ . Compared to demographic parity, equalized odds is more flexible as it does not prevent learning a predictor where there is a real association between the protected attribute and the outcome.<sup>43</sup>

## Equal opportunity

The “equal opportunity” definition of Fairness checks whether the positive label is equally and accurately predicted by classifier for all values of the protected attribute.<sup>39</sup> In contrast to equalized odds, it is stronger because it means that all possible thresholds are equally likely to be met and therefore requires that all groups produce the same ROC curve, but the decision threshold can be adjusted to satisfy equalized odds.<sup>43</sup>

## Individual fairness

The notion of individual fairness is based on the principle that any two individuals who are similar in the

context of a given task should be treated similarly.<sup>40,44</sup> Clearly, individual fairness is more restrictive than group fairness defined by the protected attribute. The practical use of this concept is often limited due to the challenges of defining an appropriate similarity metric to encode the desired concept of fairness.<sup>40,43</sup> In addition, there were also arguments that individual fairness is inadequate, as similar treatment are not enough to achieve fairness; thus, it should not be used alone to detect bias or evaluate whether algorithms are fair.<sup>45</sup> The formulation of individual fairness remains an active area of research.

## Counterfactual measures

Counterfactual fairness is a potential way to explain why bias occurs. It states that a model is fair if its predictions about a particular individual in the real world is the same as it would be in a counterfactual world (i.e., in this case, if the patient’s ethnic group was changed from Black to white).<sup>37</sup> We list the mathematical definition of counterfactual fairness in the last row of Table 2, where  $\hat{Y}_{A \leftarrow a}$  represents the prediction  $\hat{Y}$  if  $A$  had taken value  $a$ . This metric considers the predictor to be fair if its prediction remains unchanged when the protected attribute of each sample is flipped to its counterfactual value. A close concept of counterfactual fairness is counterfactual reasoning.<sup>46</sup> Some studies have shown that counterfactual reasoning is susceptible to similar biases as outcome bias (evaluating the quality of decisions when the outcome is known).<sup>47</sup> In addition, it has been suggested that counterfactual reasoning may negatively affect the process of causality identification.<sup>48</sup> These concerns raise questions about the practical applicability of counterfactual measures.

## Choice of Fairness Metric

As described above, different metrics have different characteristics. According to Kleinberg et al.,<sup>49</sup> these aforementioned fairness metrics cannot be achieved at the same time, except in highly restricted special cases. Specifically, both equalized odds and demographic parity focus on group fairness. Although their calculations and reasoning are simple and intuitive, the derived models may be discriminatory to structured subgroups with protected attributes, leading to fairness gerrymandering.<sup>43,50</sup> The concept of individual fairness potentially alleviate the issues of group fairness metrics by forcing any two individuals who are similar at a given task should be similarly classified. However, it is challenging to a domain-specific similarity measure, thus the practical use of individual fairness is often limited. Clinical prediction models may produce unfair results based on particular metrics. The choice of the specific fairness metric used by researchers and machine learning developers thus depends on the specific context. In

Type	Definition	In our case study
Fairness Through Unawareness <sup>37</sup>	No protected attribute $A$ is explicitly used in the decision-making process: $\hat{Y} = f(X, A) = f(X)$	Train the model without using race variable
Demographic Parity <sup>38</sup> / Statistical Parity / Independence	The outcomes must be equal: $P(\hat{Y} A = 0) = P(\hat{Y} A = 1)$	Both demographic groups developed sepsis at equal rates
Equalized Odds <sup>39</sup> / Separation	Different groups deal with similar odds, if $\hat{Y}$ and $A$ are independent conditional on $Y$ : $P(\hat{Y} = 1 A = 0, Y = y) = P(\hat{Y} = 1 A = 1, Y = y), y \in \{0, 1\}$	The true positive rates (of those who actually developed sepsis, how many were correctly predicted to be positive) and false positive rates in both demographic groups are equal
Equal Opportunity <sup>39</sup>	The true positive rates in the unprivileged group and privileged group are equal. $P(\hat{Y} = 1 A = 0, Y = 1) = P(\hat{Y} = 1 A = 1, Y = 1)$ ,	The true positive rates in both demographic groups are equal
Individual Fairness <sup>40</sup>	Similar individuals have similar predictions. Formally, given a metric $d(\cdot, \cdot)$ , if individuals $i$ and $j$ are similar under this metric (i.e., $d(i, j)$ is small), then their predictions should be similar: $\hat{Y}(X^{(i)}, A^{(i)}) \approx \hat{Y}(X^{(j)}, A^{(j)})$	Similar patients have a similar chance of developing sepsis
Counterfactual Fairness <sup>37</sup>	Predictor $\hat{Y}$ is counterfactually fair if under any context $X = x$ and $A = a$ , $Pr(\hat{Y}_{A \leftarrow a'}(U) = y X = x, A = a) = Pr(\hat{Y}_{A \leftarrow a}(U) = y X = x, A = a)$ , for all $y$ and for any value $a'$ attainable by $A$	The predicted outcome does not change if a patient from one demographic group is assigned to the other demographic group

**Table 2: Summary of fairness metrics.**

addition to these computational aspects, a more fundamental consideration is whether the bias should be attributed to machine learning algorithms at all. Biological and socioeconomic factors can contribute to inherent bias as well. Therefore, it is important to work with domain and legal experts to first understand the problem context and decide whether a machine learning algorithm should be used at all (e.g., for ethical concerns) and whether it can induce potential bias, and then choose an appropriate fairness metric.

## Bias mitigation

With the various sources of bias and different fairness metrics, in this section we will summarize different bias mitigation approaches for achieving algorithmic fairness. These methods can be categorized as pre-processing,<sup>51</sup> in-processing,<sup>52–55</sup> and post-processing methods,<sup>56</sup> which are detailed below.

### Pre-processing

Data pre-processing refers to the procedures of cleaning and preparing raw data for building machine learning models.<sup>57</sup> Pre-processing methods can potentially remove the bias from the data.

**Choice of sampling.** Resampling is a popular pre-processing method to ensure the datasets are balanced across different groups.<sup>58</sup> In the context of algorithmic

fairness, the use of resampling is not to address class imbalance, but rather to ensure that all demographic groups are properly and proportionately represented in the training dataset. If the data set is large, the majority group can be randomly undersampled so that it is approximately the same size as the minority group without much information loss. However, since the data is often limited, it is more common to oversample the minority groups in the training data. Popular algorithms, like synthetic minority oversampling technique (SMOTE)<sup>59</sup> or its variations, such as SMOTE-ENC,<sup>60</sup> Borderline-SMOTE,<sup>61</sup> can be used to oversample or synthetically expand the size of the data from an under-represented demographic group. However, healthcare data (such as EHRs or questionnaires) are typically complicated, and it is thus challenging to generate synthetic data without producing overfitting.<sup>12</sup> In addition to resampling, collecting more data with good planning is always the best solution.<sup>33</sup>

**Reweighting.** Another method to train an algorithm to place a greater emphasis on an under-represented group is to use reweighting. This approach places different weights on each group-class combination based on the conditional probability of class by protected attribute, so that the protected attribute is independent of the outcome.<sup>51</sup> As a representative method, inverse propensity score weighting (IPW)<sup>62</sup> is often adopted to adjust poorly sampled data. It involves estimating the probability of individual participants in particular



groups and then analyzing the re-weighted samples of these participants.<sup>63</sup> However, IPW adjusts the distributions of all variables simultaneously, which may potentially increase imbalances and bias.<sup>64</sup> Borland et al.<sup>65</sup> presented dynamic reweighting (DR) to correct selection bias with interactive visual analysis.

## In-processing

In-processing methods aims at developing unbiased models directly from the data. A straightforward approach to achieve this goal is to remove the protected attribute from the model as we introduced in Section 3.1. However, if there are strong correlations between the protected attribute and other covariates, the information of the sensitive attributes will naturally introduce bias into the decision.

**Prejudice remover.** Prejudice refers to the fact that there is statistical dependence between the protected attribute and the predicted outcome or other independent variables.<sup>66</sup> Prejudice remover is a method that attempts to train a predictor whose predictions are independent of the protected attribute. For example, Kamiran and Calders et al.<sup>67</sup> proposed the concept of discrimination-aware classification and developed an algorithm to “clear away” such dependencies by “massaging the data” before applying traditional classification algorithms. Calders and Verwer<sup>52</sup> proposed a discrimination-free naive-Bayes through post-hoc processing, independent model training and balancing across different protected groups, or latent variable modeling. Kamishima et al.<sup>54</sup> proposed a prejudice remover regularization to enforce the prediction’s independence on the protected attribute. Zafar et al.<sup>53</sup> proposed the concept of “disparate mistreatment” as different misclassification rates across different protected groups, and introduced a measure for decision boundary based classifiers, which further can be incorporated into the classifier optimization objectives as constraints to remove prejudice. With increasing numbers of machine learning models being developed for clinical risk prediction, there have also been intense discussions on the corresponding ethical concerns.<sup>68,69</sup> These prejudice remover approaches can potentially make these algorithms fair.

**Adversarial learning.** Adversarial learning<sup>70</sup> is a learning paradigm that was originally designed for generating false samples to confuse the model. Typically, there is a generator guaranteeing the generated fake samples which are close to real samples, and a discriminator to discriminate the fake samples from the real ones. The goal of adversarial learning is to learn a generator to generate samples that the discriminator cannot really tell they faked or no. Pfohl et al.<sup>43</sup> applied adversarial learning for developing an “equitable” risk prediction model for atherosclerotic cardiovascular disease

(ASCVD) with EHR. They used the generator to build the risk predictor and discriminator to enforce equalized odds for the predicted risks across different protected groups.

**Other learning strategies.** Another closely related topic is interpretable learning,<sup>71</sup> as interpretable models can allow the decision makers to better understand why certain predictions are made and make necessary modifications. Recent work at the FICO Data Science Lab<sup>72</sup> has shown that interpretable neural networks can help uncover and eliminate data biases in models. Even in cases where the data is deliberately biased toward one subset of the population over another, the method minimizes the pickup of signals that are biased toward the core relationship.<sup>72</sup> Similar argument has also been made by Rudin<sup>73</sup> that interpretable models are more preferred in high stakes decision making scenarios such as healthcare than black-box models.

Independent learning is another bias mitigation strategy which trains a machine learning model for each protected group.<sup>74</sup> However, this approach may sacrifice the training data sample sizes and reduce the model performance.<sup>74</sup> Gao and Cui<sup>74</sup> introduced a transfer learning approach to align the sample distributions across different protected groups. They demonstrated that their method could achieve improved performance in underrepresented groups and effectively reduce disparity with cancer multiomics data.<sup>74</sup>

## Post-processing

The post-processing approach treats off-the-shelf predictors as black boxes and achieves fairness through adjustment of their predictions. For example, Hardt et al.<sup>39</sup> proposed equalized odds post-processing and calibrated equalization odds post-processing, which aims to solve for the probabilities of changing output labels to achieve the equalized odds objective. Kallus et al.<sup>75</sup> proposed to adjust the risk scores of the instances in the disadvantaged group with a parameterized monotonically increasing function to minimize the performance disparity. Cui et al.<sup>76</sup> proposed to adjust the ranking order of the samples across different protected groups according to their predicted scores with a dynamic programming procedure to achieve fairness without sacrificing prediction accuracy. One practical challenge for post-processing methods is that the involved adjustments are typically not explainable. Pan et al.<sup>77</sup> proposed a causal analysis approach that can quantitatively attribute algorithm performance disparity onto different causal decision paths, so that the paths with large contributions can be removed as post-processing.

In practice, these three types of methods work at different stages of a machine learning pipeline: pre-processing manipulates the data through sampling or

weighting before building the model, in-processing enforces fairness constraints during model building, and post-processing makes adjustments after the model was built. Different strategies have different assumptions; therefore, it is challenging to have a golden standard. Recent research from Park et al.<sup>78</sup> compared different risk mitigation methods in the context of post-partum risk prediction and found that these methods could indeed reduce bias, but different methods can lead to different results. Therefore, the practitioners should try to test different approaches and evaluate their impact in the particular context they were applied to.

### Popular software libraries

Over the past few years, a variety of software tools and libraries have emerged to help developers and users of machine learning algorithms to better explore the issue of fairness and bias. Some of these libraries include tools to visualize and measure the amount of bias in the training data. Other libraries provide tools that can evaluate the algorithm results based on various fairness metrics. We summarize existing popular algorithmic fairness research software libraries in [Table 3](#). Detailed comparison of some software libraries can be referred to recent articles.<sup>79,80</sup>

### Open questions

As data is the source for building machine learning models, it is critical to be aware of the potential bias and

improve the diversity and inclusiveness during the data collection process. In addition, we list some probably encountered directions or open questions in this section.

### Multiform fairness

Different types of fairness are sometimes incompatible. For example, a model could be fair for equal positive and negative predictive values, but unfair for equalized odds (and vice versa). It is important to understand which types of fairness are achievable under which scenarios. Therefore, fairness in computational medicine requires not only machine learning/computer scientists to understand, but also experts across disciplines to work together to come up with definitions that fit a particular model and apply them to a given context.

### Algorithm explainability

Explainable models can reveal how a machine learning algorithm works and thus potentially alleviate decision bias. However, on the other hand, interacting with incorrect recommendations paired with explanations that contain limited but easily interpretable information can adversely affect the clinician's treatment choices.<sup>90</sup> Understanding such interaction between algorithm explainability and bias is important for medical machine learning.

### Model generalization

Fairness in machine learning goes beyond preventing models from harming protected populations. It can also

Project Name	Developer	Description
FairMLHealth <sup>81</sup>	KenSci	Tools and tutorials for evaluating bias in healthcare machine learning.
AIF360 <sup>82</sup>	IBM	Fairness metrics for datasets and machine learning algorithms, interpretation of the metrics, and approaches for reducing bias in datasets and models. It is available in both Python and R.
Fairlearn <sup>83</sup>	Microsoft	A Python package to evaluate fairness and mitigate any observed inequities. Fairlearn includes mitigation algorithms and metrics for model evaluation. It also contains Jupyter notebooks with examples of Fairlearn usage.
Fairness-comparison <sup>84</sup>	Sorelle et al.	Compare fairness-aware machine learning techniques. It aims to facilitate benchmarking of fairness-aware machine learning algorithms.
MEASURES <sup>85</sup>	Cardoso et al.	A benchmark framework for assessing discrimination-aware models.
Fairness Indicators <sup>86</sup>	Google	A suite of tools built on top of TensorFlow Model Analysis that enable regular evaluation of fairness metrics in product pipelines.
ML-fairness-gym <sup>87</sup>	Google	A general framework for studying and exploring long-term equity effects in carefully constructed simulation scenarios where learning subjects interact with the environment over time.
themis-ml <sup>88</sup>	Niels Bantilan	A Python library built on top of pandas and sklearn that implements fairness-aware machine learning algorithms.
FairML <sup>89</sup>	Julius Adebayo	A Python toolkit for auditing machine learning model deviations.

**Table 3: Popular library for fairness research.**

help focus care where it is really needed. The data used to develop the model may not be generalized to the data used during the deployment of the model (training-serving skew).<sup>10</sup> Thus, besides model design and evaluation, fairness should also be incorporated into the scenario where the model is going to be deployed.<sup>91</sup>

## Conclusions

In this review, we summarized the current research on algorithmic fairness in computational medicine. We first described the three types of computational bias: data bias, measurement bias, and model bias. Then we presented the fairness quantification metrics that are used in various literature. Additionally, we introduced three types of bias mitigation methods, namely, pre-processing, in-processing and post-processing, and listed the popular software libraries and tools for bias evaluation and mitigation. Fairness is not just the result of rigorous and thoughtful research, but rather the social and political processes needed to advance health equity.<sup>92</sup> With machine learning and artificial intelligence models gaining more attention, we should be aware of these issues when designing the models and appropriately mitigate them.

## Search strategy and selection criteria

We searched PubMed and Google Scholar from inception of the database to Jul 30, 2021, for research articles using the search terms (“bias” OR “disparity” OR “fairness” OR “fair” OR “inequality” OR “equality”) AND (“machine learning” OR “artificial intelligence”) AND (“medical” OR “medicine” OR “healthcare”) in English. We independently reviewed the title and abstracts for inclusion. We also reviewed the reference lists of eligible texts.

## Contributors

All authors read and approved the final version of the manuscript. JX drafted the manuscript. FW made thorough revisions to the draft. YX performed an initial literature review on computational bias. WW, YN and ES performed the literature review and data abstraction on bias mitigation methods. All authors contributed to the writing and editing of the manuscript. JB and FW conceived the idea.

## Declaration of interests

None declared.

## Acknowledgements

FW is supported by NSF 1750326, NIH R01AG076234, R01MH124740 and R01AG072449. YX is supported by

CORONAVIRUSHUB-S-21-00188. YN is supported by NSF 1948432 and 2047843. WHW is supported by NSF 2029038 and 2135988. JB is supported by NIH R01AG076234, R01CA246418, R21CA253394, R21AG068717, and R21CA245858. The funders did not play any role in paper design, data collection, data analysis, interpretation, or writing of the paper.

## References

- 1 Esteva A, Kuprel B, Novoa RA, et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*. 2017;542:115–118.
- 2 Tomašev N, Glorot X, Rae JW, et al. A clinically applicable approach to continuous prediction of future acute kidney injury. *Nature*. 2019;572:116–119.
- 3 Komorowski M, Celi LA, Badawi O, Gordon AC, Faisal AA. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nat Med*. 2018;24:1716–1720.
- 4 Wang F, Preiner A. AI in health: state of the art, challenges, and future directions. *Yearb Med Inf*. 2019;28:016–026.
- 5 Gijssels CM, Groenewegen KA, Hoefer IE, et al. Race/ethnic differences in the associations of the framingham risk factors with carotid int and cardiovascular events. *PLoS One*. 2015;10:e0132321.
- 6 Zou J, Schiebinger L. AI can be sexist and racist—it's time to make it fair. *Nature*. 2018;324–326.
- 7 Kadambi A. Achieving fairness in medical devices. *Science*. 2021;372:30–31.
- 8 Sjoding MW, Dickson RP, Iwashyna TJ, Gay SE, Valley TS. Racial bias in pulse oximetry measurement. *New Engl J Med*. 2020;383:2477–2478.
- 9 Mehrabi N, Morstatter F, Saxena N, Lerman K, Galstyan A. A survey on bias and fairness in machine learning. *ACM Comput Surv (CSUR)*. 2021;54:1–35.
- 10 Rajkomar A, Hardt M, Howell MD, Corrado G, Chin MH. Ensuring fairness in machine learning to advance health equity. *Ann Internal Med*. 2018;169:866–872.
- 11 Gianfrancesco MA, Tamang S, Yazdany J, Schmajuk G. Potential biases in machine learning algorithms using electronic health record data. *JAMA Int Med*. 2018;178:1544–1547.
- 12 Fletcher RR, Nakeshimana A, Olubeko O. Addressing fairness, bias, and appropriate use of artificial intelligence and machine learning in global health. *Front Artif Intell*. 2021;3:116.
- 13 Mhasawade V, Zhao Y, Chunara R. Machine learning and algorithmic fairness in public and population health. *Nat Mach Intell*. 2021;1–8.
- 14 Jordan MI, Mitchell TM. Machine learning: Trends, perspectives, and prospects. *Science*. 2015;349:255–260.
- 15 Ng JH, Ye F, Ward LM, Haffer SC, Scholle SH. Data on race, ethnicity, and language largely incomplete for managed care plan members. *Heal Aff*. 2017;36:548–552.
- 16 Waite S, Scott J, Colombo D. Narrowing the gap: imaging disparities in radiology. *Radiology*. 2021;299:27–35.
- 17 Heckman JJ. Sample selection bias as a specification error. *Applied Econometrics*. 2013;31(3):129–137.
- 18 Adamson AS, Smith A. Machine learning and health care disparities in dermatology. *JAMA Dermatol*. 2018;154:1247–1248.
- 19 Lumaka A, Cosmans N, Lulebo Mampasi A, et al. Facial dysmorphism is influenced by ethnic background of the patient and of the evaluator. *Clin Genet*. 2017;92:166–171.
- 20 Nunan D, Aronson J, Bankhead C. Catalogue of bias: attrition bias. *BMJ Evid-Based Med*. 2018;23:21–22.
- 21 Hernán MA, Robins JM. Using big data to emulate a target trial when a randomized trial is not available. *Am J Epidemiol*. 2016;183:758–764.
- 22 Jennions MD, Lortie CJ, Rosenberg MS, Rothstein HR. Publication and related biases. *Handb Meta-Anal Ecol Evol*. 2013;207–236.
- 23 Dickersin K, Min Y-I. NIH clinical trials and publication bias. *Online J Curr Clin Trials*. 1993;31:4967.
- 24 Scherer RW, Meerpohl JJ, Pfeifer N, Schmucker C, Schwarzer G, von Elm E. Full publication of results initially presented in abstracts. *Cochrane Database of Systematic Reviews (Online)*. 2018.
- 25 Raynaud M, Zhang H, Louis K, et al. COVID-19-related medical research: a meta-research and critical appraisal. *BMC Medical Res Methodol*. 2021;21:1–11.



- 26 Coggon D, Barker D, Rose G. *Epidemiology for the Uninitiated*. John Wiley & Sons; 2009.
- 27 Vasudevan V, Gnanasekaran A, Sankar V, Vasudevan SA, Zou J. Disparity in the quality of covid-19 data reporting across india. *Bmc Public Health*. 2021;21:1–12.
- 28 Glen, S. Response bias: Definition and examples. *From StatisticsHowTo.com: elementary Statistics for the rest of us!* <https://www.statisticshowto.com/response-bias/>.
- 29 Paulhus DL. Measurement and control of response bias. *Meas Personal Soc Psychol Attitudes*. 1991.
- 30 van den Akker M, Buntinx F, Metsemakers J, Knottnerus J. Morbidity in responders and non-responders in a register-based population survey. *Fam practice*. 1998;15:261–263.
- 31 Hooker S. Moving beyond “algorithmic bias is a data problem”. *Patterns*. 2021;2:100241.
- 32 Carbonell JG, Michalski RS, Mitchell TM. An overview of machine learning. *Mach Learn*. 1983;1:3–23.
- 33 Chen IY, Johansson FD, Sontag D. Why is my classifier discriminatory? In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems*. 2018:3543–3554.
- 34 Wong A, et al. External validation of a widely implemented proprietary sepsis prediction model in hospitalized patients. *JAMA Intern Med*. 2021;181:1065–1070.
- 35 Ahmad MA, Patel A, Eckert C, Kumar V, Teredesai A. Fairness in machine learning for healthcare. In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2020:3529–3530.
- 36 Verma S, Rubin J. Fairness definitions explained. 2018 *IEEE/ACM International workshop on software fairness (fairware)*. IEEE; 2018:1–7.
- 37 Kusner M, Loftus J, Russell C, Silva R. Counterfactual fairness. *Adv Neural Inf Process Syst* 30 (NIPS 2017). 2017;30:4069–4079.
- 38 Calders T, Kamiran F, Pechenizkiy M. Building classifiers with independency constraints. 2009 *IEEE International Conference on Data Mining Workshops*. IEEE; 2009:13–18.
- 39 Hardt M, Price E, Srebro N. Equality of opportunity in supervised learning. *Adv Neural Inf Process Syst*. 2016;29:3323–3331.
- 40 Dwork C, Hardt M, Pitassi T, Reingold O, Zemel R. Fairness through awareness. In: *Proceedings of the 3rd innovations in theoretical computer science conference*. 2012:214–226.
- 41 Luong BT, Ruggieri S, Turini F. k-NN as an implementation of situation testing for discrimination discovery and prevention. In: *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2011:502–510.
- 42 Grgic-Hlaca N, Zafar MB, Gummadi KP, Weller A. The case for process fairness in learning: Feature selection for fair decision making. *NIPS Symposium on Machine Learning and the Law*. vol. 1. 2016:2.
- 43 Pfohl S, Marafino B, Coulet A, Rodriguez F, Palaniappan L, Shah NH. Creating fair models of atherosclerotic cardiovascular disease risk. In: *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. AAAI; 2019:271–278.
- 44 Zemel R, Wu Y, Swersky K, Pitassi T, Dwork C. Learning fair representations. *International Conference on Machine Learning*. PMLR; 2013:325–333.
- 45 Will Fleisher W. What’s fair about individual fairness? *AAAI/ACM Conference on AI, Ethics, and Society (AIES)*. ACM; 2021.
- 46 Lewis D. Causation. *J Philosophy*. 1974;70:556–567.
- 47 Baron J, Hershey JC. Outcome bias in decision evaluation. *J Personal Soc Psychol*. 1988;54:569.
- 48 Dawid AP. Causal inference without counterfactuals. *J Am Statistical Assoc*. 2000;95:407–424.
- 49 Kleinberg J, Mullainathan S, Raghavan M. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*. 2016.
- 50 Kearns M, Neel S, Roth A, Wu ZS. Preventing fairness gerrymandering: auditing and learning for subgroup fairness. *International Conference on Machine Learning*. PMLR; 2018:2564–2572.
- 51 Kamiran F, Calders T. Data preprocessing techniques for classification without discrimination. *Knowl Inf Syst*. 2012;33:1–33.
- 52 Calders T, Verwer S. Three naive bayes approaches for discrimination-free classification. *Data Mining Knowl Discovery*. 2010;21:277–292.
- 53 Zafar MB, Valera I, Gomez Rodriguez M, Gummadi KP. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In: *Proceedings of the 26th International Conference on World Wide Web*. 2017:1171–1180.
- 54 Kamishima T, Akaho S, Sakuma J. Fairness-aware learning through regularization approach. 2011 *IEEE 11th International Conference on Data Mining Workshops*. IEEE; 2011:643–650.
- 55 Goodfellow I, Pouget-Abadie J, Mirza M, et al. Generative adversarial nets. *Communications of the ACM*. 2020;63(11):139–144.
- 56 Tang Z, Zhang K. Attainability and optimality: the equalized-odds fairness revisited. *arXiv preprint arXiv:2202.11853*. 2020.
- 57 Zhang S, Zhang C, Yang Q. Data preparation for data mining. *Appl Artificial Intell*. 2003;17:375–381.
- 58 Kamiran F, Calders T. Classification with no discrimination by preferential sampling. In: *Proc. 19th Machine Learning Conf. Belgium and The Netherlands*. Citeseer; 2010:1–6.
- 59 Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. Smote: synthetic minority over-sampling technique. *J Artificial Intell Res*. 2002;16:321–357.
- 60 Mukherjee M, Khushi M. Smote-enc: A novel smote-based method to generate synthetic data for nominal and continuous features. *Appl Syst Innov*. 2021;4:18.
- 61 Han H, Wang W-Y, Mao B-H. Borderline-smote: a new over-sampling method in imbalanced data sets learning. *International conference on intelligent computing*. Springer; 2005:878–887.
- 62 Caliendo M, Kopeinig S. Some practical guidance for the implementation of propensity score matching. *J Econ Surveys*. 2008;22:317–32.
- 63 Nilsson A, Bonander C, Stromberg U, Canivet C, Ostergren P-O, Bjork J. Reweighting a swedish health questionnaire survey using extensive population register and self-reported data for assessing and improving the validity of longitudinal associations. *Plos One*. 2021;16:e0253969.
- 64 King G, Nielsen R. Why propensity scores should not be used for matching. *Polit Anal*. 2019;27:435–454.
- 65 Borland D, Zhang J, Kaul S, Gotz D. Selection-bias-corrected visualization via dynamic reweighting. *IEEE Trans Vis Comput Graph*. 2020;27:1481–1491.
- 66 Kamishima T, Akaho S, Asoh H, Sakuma J. Fairness-aware classifier with prejudice remover regularizer. *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer; 2012:35–50.
- 67 Kamiran F, Calders T. Classifying without discriminating. 2009 *2nd International Conference on Computer, Control and Communication*. IEEE; 2009:1–6.
- 68 Char DS, Shah NH, Magnus D. Implementing machine learning in health care—addressing ethical challenges. *New Engl J Med*. 2018;378:981.
- 69 Cohen IG, Amarasingham R, Shah A, Xie B, Lo B. The legal and ethical concerns that arise from using complex predictive analytics in health care. *Heal Affairs*. 2014;33:1139–1147.
- 70 Huang L, Joseph AD, Nelson B, Rubinstein BI, Tygar JD. Adversarial machine learning. In: *Proceedings of the 4th ACM workshop on Security and artificial intelligence*. 2011:43–58.
- 71 Wang F, Kaushal R, Khullar D. Should health care demand interpretable artificial intelligence or accept “black box” medicine? *Ann Internal Med*. 2020;172:59–60.
- 72 Zoldi, S. Fighting bias: How interpretable latent features remove bias in neural networks. 2001. <https://www.fico.com/blogs/fighting-bias-how-interpretable-latent-features-remove-bias-neural-net-works>.
- 73 Rudin C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell*. 2019;1:206–215.
- 74 Gao Y, Cui Y. Deep transfer learning for reducing health care disparities arising from biomedical data inequality. *Nat Commun*. 2020;11:1–8.
- 75 Kallus N, Zhou A. The fairness of risk scores beyond classification: bipartite ranking and the xauc metric. *Adv Neural Inf Process Syst*. 2019:32.
- 76 Cui S, Pan W, Zhang C, Wang F. Towards model-agnostic post-hoc adjustment for balancing ranking fairness and algorithm utility. In: *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2021:207–217.
- 77 Pan W, Cui S, Bian J, Zhang C, Wang F. Explaining algorithmic fairness through fairness-aware causal path decomposition. In: *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2021:1287–1297.
- 78 Park Y, Hu J, Singh M, et al. Comparison of methods to reduce bias from clinical prediction models of postpartum depression. *JAMA Netw Open*. 2021;4(4):e213909.

- 79 Swan Lin. Comparing VerifyML, AI Fairness 360 and Fairlearn. Medium. 2021. <https://medium.com/cylynx/verifyml-where-it-stands-among-other-ai-fairness-toolkits-8e6cad149b2>.
- 80 Pandey H. Comparison of the usage of Fairness Toolkits amongst practitioners: AIF360 and Fairlearn. TUDelft. 2022. <http://resolver.tudelft.nl/uuid:4ef11035-2f60-436f-85f9-7b9bed73b66d>.
- 81 Allen C, Ahmad MA, Eckert C, Hu J, Kumar V, Teredesai A. fairML-Health: Tools and tutorials for fairness evaluation in healthcare machine learning. 2020; <https://github.com/KenSciResearch/fairMLHealth>.
- 82 Bellamy RK, Dey K, Hin M, et al. AI Fairness 360: an extensible toolkit for detecting and mitigating algorithmic bias. *IBM J Res Dev*. 2019;63(4/5):1–4.
- 83 Bird S, Dudík M, Edgar R, et al. Fairlearn: a toolkit for assessing and improving fairness in AI. *Microsoft, Tech Rep*. 2020;MSR-TR-2020-32.
- 84 Friedler SA, Scheidegger C, Venkatasubramanian S, Choudhary S, Hamilton EP, Roth D. A comparative study of fairness-enhancing interventions in machine learning. In: *Proceedings of the conference on fairness, accountability, and transparency*. ACM; 2019:329–338.
- 85 Cardoso RL, Meira Jr. W, Almeida V, Zaki MJ. A framework for benchmarking discrimination-aware models in machine learning. In: *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. ACM; 2019:437–444.
- 86 Google. Tensorflow fairness indicators. [https://www.tensorflow.org/responsible\\_ai/fairness\\_indicators/tutorials/Fairness\\_Indicators\\_Example\\_Colab](https://www.tensorflow.org/responsible_ai/fairness_indicators/tutorials/Fairness_Indicators_Example_Colab).
- 87 Google. ML-fairness-gym: a tool for exploring long-term impacts of machine learning systems. <https://ai.googleblog.com/2020/02/ml-fairness-gym-tool-for-exploring-long.html> (2020).
- 88 Bantilan, N. A library that implements fairness-aware machine learning algorithms. <https://themis-ml.readthedocs.io/en/latest/>.
- 89 Adebayo, J. FairML - is a python toolbox auditing the machine learning models for bias. <https://github.com/adebayoj/fairml>.
- 90 Jacobs M, Pradier MF, McCoy TH, Perlis RH, Doshi-Velez F, Gajos KZ. How machine-learning recommendations influence clinician treatment selections: the example of antidepressant selection. *Transl Psychiatry*. 2021;11:1–9.
- 91 Cui S, Pan W, Liang J, Zhang C, Wang F. Addressing algorithmic disparity and performance inconsistency in federated learning. *Adv Neural Inf Process Syst*. 2021;34.
- 92 Sikstrom L, Maslej MM, Hui K, Findlay Z, Buchman DZ, Hill SL. Conceptualising fairness: three pillars for medical algorithms and health equity. *BMJ Heal Care Inf*. 2022;29.