

# Racial Bias in Clinical and Population Health Algorithms: A Critical Review of Current Debates

Madison Coots<sup>1</sup>

Kristin A. Linn<sup>2</sup>

Sharad Goel<sup>1</sup>

Amol S. Navathe<sup>2</sup>

Ravi B. Parikh<sup>3</sup>

<sup>1</sup> Harvard Kennedy School, Harvard University, Cambridge, Massachusetts, USA

mcoots@g.harvard.edu

sgoel@hks.harvard.edu

<sup>2</sup> Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

klinn@pennmedicine.upenn.edu

amol@pennmedicine.upenn.edu

<sup>3</sup> Emory University School of Medicine, Atlanta, Georgia, USA

ravi.parikh@pennmedicine.upenn.edu

## **Keywords**

Race, bias, algorithms, fairness, equity, healthcare

## **Abstract**

Among healthcare researchers, there is increasing debate over how best to assess and ensure the fairness of algorithms used for clinical decision support and population health—particularly concerning potential racial bias. Here we first distill concerns over the fairness of healthcare algorithms into four broad categories: (1) the explicit inclusion (or, conversely, the exclusion) of race and ethnicity in algorithms, (2) unequal algorithm decision rates across groups, (3) unequal error rates across groups, and (4) potential bias in the target variable used in prediction. With this taxonomy, we critically examine seven prominent and controversial healthcare algorithms. We show that popular approaches that aim to improve the fairness of healthcare algorithms can in fact worsen outcomes for individuals across all racial and ethnic groups. We conclude by offering an alternative, consequentialist framework for algorithm design that mitigates these harms by instead foregrounding outcomes and clarifying tradeoffs in the pursuit of equitable decision making.

## **I. INTRODUCTION**

Over the past two decades, organizations across sectors have developed and deployed algorithms to enhance decision-making. In healthcare, algorithms are increasingly used for guiding high-stakes decisions, including disease screening and treatment, as well as for allocating limited healthcare resources. Algorithmic decision-making promises to increase efficiency and reduce subjectivity, but researchers and clinicians have raised concerns that healthcare algorithms exacerbate inequities, particularly relating to race and ethnicity.

In this article, we survey a wide range of prominent healthcare algorithms used for population health and individual decision support, all of which have been subject of extensive debates over their fairness. While “fairness” itself is a contested term, we use it here to describe an algorithm’s tendency to enhance rather than diminish the equity of decisions. The clinical algorithms we consider span several medical fields, including oncology, obstetrics, cardiology, and nephrology, and concern both shared decision-making and the allocation of limited resources. Using these algorithms as case studies, we distill the myriad fairness concerns that have been raised for healthcare algorithms into four categories: (1) the inclusion (or, conversely, the exclusion) of race and ethnicity in algorithm inputs, (2) unequal algorithm decision rates

across groups, (3) unequal error rates across groups, and (4) potential bias in the target variable used in prediction. In Table 1, we situate each of the algorithms we consider into this structure.

We briefly summarize these four broad concerns before discussing them in depth in subsequent sections. The use of race and ethnicity in health algorithms is heavily debated (62-64, 18-20, 57, 69). Advocates for race-aware algorithms argue that explicitly considering race improves the accuracy of decisions for all groups (38, 39, 42, 1, 69). Conversely, proponents of race-unaware algorithms—commonly called “race-blind” algorithms—argue that using race perpetuates pernicious racial attitudes and exacerbates racial and ethnic inequities (23, 62-64, 33, 57). In addition to scrutinizing an algorithm’s inputs, researchers have sought to assess the fairness of algorithms by examining differences in decision and error rates across racial and ethnic groups. For example, in the context of lung cancer, researchers have noted that common algorithms recommend different screening rates across racial groups and also exhibit racial gaps in missed referrals for screening (i.e., unequal error rates) (2, 49). As a result, many advocate for designing algorithms to equalize decision and error rates across groups. That strategy seeks to ensure that the benefits and burdens of algorithmic decision-making are shared equally across groups—though also tacitly ignores differences in need and individual preferences. Finally, there is concern over the mismeasurement of target variables used for prediction in healthcare algorithms, a problem also known as “label bias” (68). For example, algorithms trained to predict healthcare expenditure as a proxy of healthcare need may under-allocate healthcare resources to disadvantaged groups due to racial disparities in healthcare access and expenditure (48).

Many of the fairness concerns raised for healthcare algorithms focus on *how* an algorithm makes decisions—for example, does it use race, or does it equalize decisions across groups? That focus stands in contrast to a perspective that foregrounds *outcomes*, an approach to ethical decision-making often called consequentialism. From a consequentialist perspective, what renders an algorithm fair is its impact on individuals and society, not its set of inputs or some particular statistical summary. Consequentialism suggests an algorithm be designed to maximize aggregate “utility,” the overall desirability or benefit that comes from a specific decision or policy. In the healthcare context, utility might encapsulate quantitative measures such as life-years gained, or improvement in quality of life, as well as the monetary and non-monetary costs associated with unnecessary testing.

**Table 1. Taxonomy of healthcare algorithms and their fairness concerns.** We situate each of the seven algorithms we consider in this review within our taxonomy of the four main fairness concerns for clinical and population health algorithms. LYFS refers to life-years from screening, VBAC refers to vaginal birth after cesarean, CVD refers to cardiovascular disease, and eGFR refers to estimated glomerular filtration rate.

		Use of race	Unequal decision rates	Unequal error rates	Label bias
<b>Non-resource constrained</b>	<b>Lung cancer incidence risk model</b>		✓	✓	
	<b>Lung cancer LYFS model</b>	✓	✓	✓	✓
	<b>VBAC success calculator</b>	✓	✓		
	<b>CVD incidence risk model</b>	✓	✓		✓
<b>Resource constrained</b>	<b>CVD hospital mortality risk model</b>		✓		
	<b>Kidney function (eGFR) equation</b>	✓	✓		
	<b>Healthcare need prediction models</b>				✓

In this article, we argue that the dominant approach to designing fair healthcare algorithms—one that, for example, seeks to equalize decision or error rates—can often harm the groups it seeks to protect. Instead, we advocate for adopting a consequentialist perspective to algorithm design. The remainder of our article is structured as follows. Beginning with an extended case study of risk algorithms for lung cancer, we unpack the four broad categories of fairness concerns for healthcare algorithms outlined above. We then expand our discussion by reviewing several more prominent healthcare algorithms that have been the subject of recent fairness debates. Finally, we conclude with an example that illustrates the value of a consequentialist framework, offering a path forward for designing equitable algorithms in healthcare and beyond.

## **II. FAIRNESS CONCERNS IN LUNG CANCER SCREENING**

In the United States, lung cancer is the third most commonly diagnosed cancer and is the leading cause of cancer-related death (8). Black men in the United States have higher rates of lung cancer incidence and mortality than men in any other racial or ethnic group (55, 56). Researchers have attributed these disparities to differences in social determinants of health—such as access to healthcare and exposure to carcinogens—that are correlated with race and ethnicity, with race and ethnicity acting as surrogates for these factors in clinical models for lung cancer (59).

Screening in the form of low-dose computed tomography (CT) scans remains the most effective method for diagnosing and informing treatment for lung cancer. However, screening comes with both monetary and non-monetary costs, such as the direct cost of the scan, taking time away from work, and psychological stress associated with screening and false positives. Given these tradeoffs, the United States Preventive Services Task Force (USPSTF) only recommends annual screening for high-risk individuals: adults aged 50 to 80 years who have a 20 pack-year smoking history and currently smoke or have quit within the past 15 years (35). Researchers have also developed risk models in an attempt to produce more personalized estimates of risk (10, 31, 61). To identify high-risk individuals, clinicians use thresholds on these risk scores to inform recommendations for follow-up screening.

The algorithms used to produce these risk scores are the subject of extensive debate, particularly with respect to their inclusion of race (37), and how the use of race interacts with the

choice of target variable, tying into discussions on label bias (59). Further, the fairness of these risk models has been characterized in terms of differences in decision and error rates across demographic groups (2, 49). Using lung cancer as an extended case study, we discuss in more depth the four categories of fairness concerns listed above, and offer empirical evidence on the consequences of employing different approaches to fairness popular in the literature.

## **A. THE USE OF RACE**

In an effort to develop the most statistically accurate models, some researchers have recommended the use of race-aware algorithms that account for racial and ethnic disparities in lung cancer incidence and mortality (37, 10, 31, 61). Medical societies, such as the American Thoracic Society, have also recommended augmenting the USPSTF guidelines with race-aware predictive models to guide decisions, in part to identify more high-risk racial minorities for screening (52). However, adjustments for race in clinical algorithms are widely debated across numerous disease contexts, as well as in medicine more generally (62-64, 30, 23, 57, 69). Many researchers have advocated for eliminating race and ethnicity in similar risk estimation tools due to concerns that including race and ethnicity reifies race as biologically meaningful and may in turn result in inappropriately racialized medical treatment (62-64, 57).

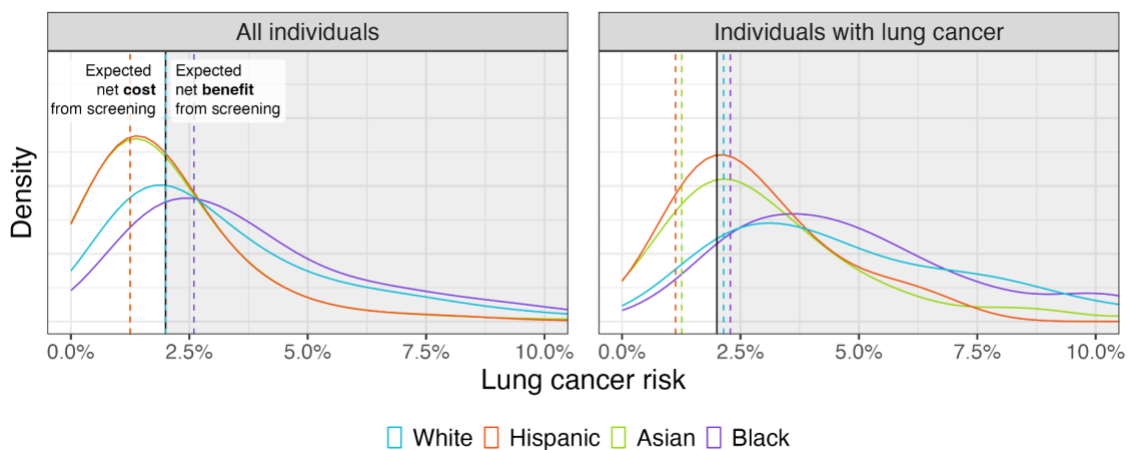
## **B. UNEQUAL DECISION RATES ACROSS GROUPS**

Researchers have sought to evaluate the fairness of lung cancer risk models by comparing screening recommendation rates across groups. For example, Landy et al. (37) propose a (necessarily) race-aware model that ensures screening rates for racial and ethnic minorities are equal to or greater than that of White individuals. While perhaps intuitively appealing, equalizing such decision rates can harm members of all groups. To see this, we note that in shared decision-making contexts such as this, the decision threshold—in this case, the screening threshold—is set at the point of indifference, where the costs of the decision are expected to equal the benefits (50). As a result, individuals with true risk above the decision threshold are expected to have positive utility from being screened. Conversely, individuals with true risk below the threshold are expected to have negative utility from being screened (e.g., from incurring the costs of screening while expecting little benefit). However, equalizing decision rates does not consider an

individual’s utility for being screened, and may thus impose utility losses from over- or under-screening individuals.

To see how equalizing decision rates can produce these undesired consequences, consider the left plot of Figure 1, which shows distributions of lung cancer risk as estimated by the Lung Cancer Risk Assessment Tool (LCRAT), disaggregated by race and ethnicity, using data from the National Lung Screening Trial (NLST) (46). (Figure 1 mirrors analyses performed by Chohlas-Wood et al. (11)). In this population, 62% of White individuals have a risk score above the recommended screening threshold of 2% (31, 32)—indicated by the vertical black line—35% of Asian individuals, 36% of Hispanic individuals, and 74% of Black patients are also above the threshold. As evident in the plot, these differences in decision rates arise primarily due to differences in the group-level distributions of risk.

**Figure 1.** *The distribution of lung cancer risk for all individuals and individuals with lung cancer. Estimates of risk were generated using the lung cancer risk assessment tool (LCRAT). The dashed vertical lines correspond to screening thresholds that equalize decision rates across groups (left) and false negative rates across groups (right).*



Now consider a policy that sets group-specific thresholds on risk so that Asian, Hispanic and Black individuals are recommended for screening at the same rate as White individuals (62%). Under this approach, we would screen Asian and Hispanic individuals with risk above 1.2%, which includes many relatively low-risk individuals for whom we expect screening to yield negative utility. Examples of sources of negative utility from screening include pain and

complications from invasive biopsies or surgeries. Yet we would only screen Black individuals who have relatively high risk of lung cancer, above approximately 2.6%, failing to screen many Black individuals for whom we expect screening to have net benefits, also resulting in lost utility. By equalizing screening rates across groups, we would in fact harm Asian, Hispanic, and Black individuals by failing to screen some individuals who are expected to benefit from screening and screening others for whom the costs are expected to outweigh the benefits.

### C. UNEQUAL ERROR RATES ACROSS GROUPS

Other work has sought to assess the fairness of the USPSTF screening criteria in terms of group-level sensitivity<sup>1</sup> (2, 49). This analysis amounts to a comparison of false negative rates (FNR), given that FNR is one minus sensitivity. Those studies have shown that the USPSTF criteria fail to recommend Black individuals with lung cancer for screening at higher rates than White individuals (2, 49). Consequently, some researchers recommend lowering the USPSTF criteria on smoking history to increase screening among Black individuals with lung cancer (2, 49). This argument, however, similarly fails to account for differences in risk distributions across groups.

To see how equalizing error rates—like equalizing decision rates—can produce undesired consequences, consider the right plot of Figure 1. We depict distributions of estimated lung cancer risk (as estimated by LCRAT) among people who have lung cancer. Under a policy of screening patients above the recommended 2% threshold (indicated by the vertical black line), Hispanic individuals have the highest false negative rate at 44%, Asian individuals have an FNR of 28%, White individuals have an FNR of 13%, and Black individuals have the lowest FNR at 10%. To equalize false negative rates across groups (while leaving risk estimates unchanged), we must set group-specific screening thresholds. For example, if we screen Hispanic individuals above a 1.1% threshold, screen Asian individuals above a 1.3% threshold, screen White individuals above a 2.1% threshold, and screen Black individuals above a 2.3% threshold, then all groups would have a false negative rate of approximately 15%. However, such a policy would mean that we recommend screening for some relatively low-risk Hispanic and Asian individuals,

---

<sup>1</sup> The sensitivity equals  $TP/(TP + FN)$ , where TP is the number of true positives, and FN is the number of false negatives.



resulting in expected utility losses for some members of these groups. In general, equalizing error rates across groups has the potential to harm members of all racial and ethnic groups.

#### **D. THE RISK OF LABEL BIAS**

Not all lung cancer risk models are designed to predict the same outcome. The LCRAT and PLCOm2012 models predict lung cancer incidence, whereas the life-years gained from screening-computed tomography (LYFS-CT) model predicts expected benefit from screening (31, 10, 61). However, researchers have expressed concern that models—particularly race-aware models—that predict life-years gained from screening are susceptible to label bias due to differences in life expectancy across racial and ethnic groups (59). Black individuals have lower average life expectancy than White individuals, with researchers attributing these gaps to differences in social determinants of health, such as income and geography (54). A model that predicts life-years gained from screening may therefore risk predicting patterns in social determinants of health as opposed to absolute benefits of screening.

As we have seen in the case of lung cancer, fairness concerns in algorithmic decision-making are nuanced, and divorcing measures of fairness from outcomes may produce unintended and undesired consequences. These issues are not isolated to lung cancer screening but are pervasive across algorithms used in healthcare, which we explore further in the following sections.

### **III. CASE STUDIES ON THE FAIRNESS OF HEALTHCARE ALGORITHMS**

We next examine several prominent—and, in some cases, controversial—healthcare algorithms used for risk estimation and resource allocation. For each algorithm, we briefly review its genesis and intended applications, and discuss the relevant fairness concerns they invoke. We first consider three risk scores used in obstetrics and cardiology. As shown in Table 1, two of these risk scores are used in a shared decision-making context, and one is used in a resource-constrained setting. We then discuss two algorithms used to allocate limited resources in nephrology and care management.

## **A. RISK ESTIMATION**

### *A.1 VAGINAL BIRTH AFTER CESAREAN (VBAC)*

Among pregnant women who previously have had cesarean sections (C-sections), there are trade-offs to a trial of labor (or TOLAC, “trial of labor after cesarean”). Vaginal births have well-established advantages over repeat C-sections, including shorter recovery times, lower risks of infection and hemorrhage, and better outcomes in future pregnancies (15). However, for some women, these benefits may not outweigh the risks. Vaginal birth after cesarean (VBAC) success calculators help healthcare providers and pregnant women weigh this decision. The first VBAC calculator—called the Grobman calculator—was developed to predict the likelihood of a successful VBAC by considering factors such as age, BMI, any prior vaginal delivery, previous VBAC, reason for prior cesarean, and race and ethnicity (36, 26, 27).

VBAC calculators have been criticized for their use of race, and the resulting disparities in TOLAC recommendation rates across racial and ethnic groups. (63, 64). The Grobman calculator produces systematically higher VBAC success probabilities for White women than for Black or Hispanic women, consistent with studies showing that Black and Hispanic women are less likely to achieve successful VBAC than White women (43). These racial disparities in estimated VBAC success rates likely exist because, conditional on other covariates, race and ethnicity capture a mix of unobserved clinical and non-clinical factors, such as income, education or healthcare access.

Researchers have hypothesized that using race-aware VBAC calculators, like the Grobman calculator, led to differences in how women of different racial and ethnic groups were counseled by their doctors to attempt TOLAC (63, 64). Fairness critiques of VBAC calculators have predominantly centered on unequal recommendations rates to attempt TOLAC across groups. Akin to our lung cancer example, equalizing VBAC recommendation rates may similarly require using race-specific thresholds on VBAC success probability, due to possible differences in risk distributions across groups. This approach, however, would likely result in not recommending TOLAC for some women above the recommended threshold and recommending TOLAC for some women below that threshold—a policy that may impose utility losses for women across all groups. These utility losses would stem from both women who would have benefited from attempting TOLAC but were not counseled to do so, as well as from women who

were inappropriately recommended TOLAC and suffered avoidable complications during delivery.

Grobman et al. (28, 29) have since developed a race-unaware VBAC calculator, which exhibits comparable overall accuracy to the previous race-aware version, though the authors did not report accuracy or recommendation rates across racial and ethnic groups (6, 27, 29). As researchers continue to evaluate the fairness of VBAC calculators and other decision aids, we believe they should move away from scrutinizing differential recommendation rates across groups, an approach that fundamentally ignores potential differences in underlying risk distributions and which may result in utility losses for members of all groups.

## *A.2 CARDIOVASCULAR DISEASE*

Clinical algorithms are commonly used to estimate cardiovascular risks, but they have been criticized for explicitly considering race, exhibiting unequal decision rates across groups, and being impacted by label bias. Cardiovascular disease encompasses a broad range of conditions and events such as coronary heart disease (CHD), coronary heart failure (CHF), heart attack, and stroke. To reduce CVD-related morbidity and mortality, clinicians have long prescribed statins, a class of cholesterol-lowering drugs, to prevent CVD events. Risk calculators for CVD were developed to help physicians determine when to recommend statins as a prophylactic.

The Framingham Risk Score (FRS), a race-unaware model trained to predict incidence of CHD, was among the first such risk estimation algorithms (66). Studies found that the FRS performed reasonably well at predicting CHD events in Black individuals, but concerns emerged over the limited scope of the FRS in only predicting CHD and not other CVD events, like CHF (17), especially because Black individuals exhibited higher rates of CHF than other racial and ethnic groups (5, 3). Racial and ethnic differences in CHF incidence rates may lead to label bias when CHD events are used as a proxy for general cardiovascular risks. To address this concern, researchers subsequently developed an expanded (but still race-unaware) CVD risk model that additionally predicted risk of stroke, peripheral artery disease, and heart failure (16).

Persistent racial and ethnic disparities in incidence rates across these events motivated the development of a race-aware CVD model: the pooled cohort equations (PCE), also known as the American College of Cardiology/American Heart Association (ACC/AHA) risk calculator (24).

The PCE outperformed other risk scores at predicting initial ASCVD events among Black individuals, as well as others (24, 25). However, like earlier risk models, the PCE failed to include CHF in its set of considered events, again raising concerns of label bias (7). In response, the AHA released the PREVENT equations for predicting risk of CVD and CVD subtypes, including CHF (34). Notably, the PREVENT equations are also race-unaware, due to a previous decision to exclude race as a predictor (33). The model instead includes a zip-code level social deprivation index, which helps account for CVD risk factors that are likely correlated with race and ethnicity (33).

The authors of PREVENT report the equations are suitably calibrated for Black individuals, even without including race or ethnicity—though some have questioned the decision to exclude race and ethnicity from the model. Diao et al. (20) characterized the expected changes in statin and antihypertensive therapy eligibility from the switch to the race-unaware PREVENT equations, finding that their use would decrease eligibility for statin and antihypertensive therapy for approximately 17 million U.S. adults and that these changes would affect a greater proportion of Black adults than White adults. It is, however, possible that the PREVENT equations yield more clinically appropriate decisions, as the PCE have been found to overestimate CVD risks for members of all racial and ethnic groups—a pattern researchers have attributed to changes in the prevalence of risk factors over time (such as smoking), and advances in care and prevention (44, 29). Further, those whose treatment recommendations change would likely have cardiovascular risks near the decision threshold, meaning they might not have benefitted substantially from treatment (14). Without explicitly considering utility (e.g., in terms of quality-adjusted life-years), it is hard to fully assess the impact of excluding race and ethnicity from the model .

While PREVENT excludes race and ethnicity, other cardiovascular risk models are still race-aware, such as the AHA Get with the Guidelines–Heart Failure Risk Score, drawing debate. That score informs triage decisions on whether to admit patients to intensive and specialty care, which are often limited resources in many hospitals (58). The model was designed to predict risk of in-hospital mortality using data from a cohort of patients hospitalized with heart failure (51). In the training data, in-hospital mortality was *lower* for Black patients, counter to expectations (51). One possible explanation for this pattern relates to racial disparities in access to cardiology care. Past work has found that one of the strongest predictors of admission to the cardiology service is whether a patient was previously seen by an outpatient cardiologist at the hospital—

and there were significant racial and ethnic disparities in the proportion of patients who had seen a cardiologist within the past year (67, 22). Because of those disparities, Black patients admitted to intensive care might have had better access to healthcare and lower risk of mortality as a result.

Due to the lower estimated risks for Black patients, researchers have raised concerns that the AHA Get with the Guidelines–Heart Failure Risk Score may misdirect intensive care away from Black patients (63). But, as discussed previously, decision rate-based criticisms of model fairness ignore potential differences in risk across groups. It may in reality be the case that non-Black patients face higher risks in this setting, and would be better served by the additional care. We would, however, caution against using this risk score in settings where the patient population differs substantially from the population used to train the score—especially if those differences are correlated with race or ethnicity. For example, in safety net hospitals, the score might under-predict mortality risk for Black individuals, since those patients might have higher risks than the relatively healthy Black individuals in the training data.

## **B. RESOURCE ALLOCATION**

### *B.1 KIDNEY TRANSPLANTS*

Like algorithms for cardiovascular disease, those for kidney disease are both widespread and controversial, having been criticized for their use of race and for exhibiting unequal decision rates across groups (63). In the United States, resources for kidney transplantation are highly constrained. In 2021, the average wait time for a deceased donor kidney was five years, and more than half of listed transplant candidates were expected to die or be removed from the list before receiving a transplant (65). In an attempt to make the best use of the constrained supply of donor kidneys, researchers and clinicians have turned to algorithms to estimate a patient’s need for transplant. Patients are typically recommended for transplants based on estimates of their kidney function, as measured by their glomerular filtration rate (GFR). Given challenges with measuring GFR directly, GFR has traditionally been estimated (eGFR) using an algorithm based on factors like age, sex, race, body size (usually weight or surface area), and serum creatinine (13, 60). In addition to resource constrained kidney allocation, eGFR equations have also been used to make non-resource constrained chronic kidney disease diagnoses and recommendations for drug treatments or other therapies (18).

Early eGFR equations—such as the Modification of Diet in Renal Disease (MDRD) Study equation and the Chronic Kidney Disease Epidemiology Collaboration (CKD-EPI) equation—used race to estimate GFR (39-41). All else being equal, these equations estimated *higher* GFR for Black individuals, meaning better kidney function. Race-aware eGFR equations drew concerns that they made Black patients appear healthier than they were in reality (63). As a result, many hospitals began using ad hoc race-unaware estimates of kidney function by reporting the “White/other” eGFR prediction for all patients (18). This strategy was found to slightly increase chronic kidney disease diagnoses among Black adults (18).

This critique of race-aware eGFR equations and the corresponding policy response mirror calls to equalize decision rates across groups. However, as discussed above, equalizing decision rates fails to allow for group differences and it risks harming members of all groups. Indeed, race-unaware eGFR values typically under-estimate GFR for Black individuals (30)—the race adjustment in the original equations was included precisely to ensure estimates were calibrated across groups. Consequently, the race-unaware estimates make Black patients appear *less* healthy than they likely are in reality. That pattern may have led to inappropriate diagnoses, potentially resulting in net harm to Black patients. Researchers have since revised the CKD-EPI equation to replace race with cystatin C as a predictor of eGFR, leading to race-unaware estimates that are approximately calibrated (30).

Race-aware eGFR equations have been similarly criticized for deprioritizing Black patients for kidney transplantation and specialist care (63)—since they estimate higher kidney function for Black individuals compared to otherwise similar White individuals. However, as discussed above, the race-aware equations produce largely accurate estimates of GFR for both Black and non-Black patients. Policymakers may well prefer to enforce some degree of parity in kidney allocation rates, even if that means prioritizing a healthier Black patient over a less healthy White patient. But we believe that these difficult tradeoffs should be confronted directly (12). Seeking to increase eligibility for kidney transplants for Black individuals by using an eGFR equation that underestimates their kidney function risks other adverse consequences in the form of over-treatment.

## *B.2 HEALTHCARE COSTS AS A PROXY OF NEED*

Healthcare providers in the United States offer specialty care management programs to improve the care of high-risk patients with complex health needs. These programs aim to help individuals better manage their health by offering additional support from teams of dedicated nurses, social workers, and community health workers. However, these programs are expensive and healthcare systems consequently use algorithms to identify patients for whom the benefits justify the additional costs (4). A common strategy is to predict patients' future medical expenses, and then direct specialty care management to those expected to incur the largest costs (48).

However, past work has demonstrated that algorithms trained to predict healthcare costs can fail to allocate resources to high-need racial minorities. Obermeyer et al. (48) evaluated the fairness of a commercial algorithm widely used by healthcare systems to guide patient referrals to specialty care programs. The algorithm was trained to predict future costs—as a proxy for complex healthcare needs—based on insurance claims (e.g., diagnoses, procedures, medications) made by an individual in the prior year. The researchers found that the algorithm's generated risk scores were well-calibrated across race groups for predicting healthcare costs. Conditional on risk score, both Black and White individuals had approximately the same costs in the following year. However, the researchers found that the algorithm was poorly calibrated for predicting realized health. Conditional on risk score, Black patients had significantly more illness burden than White patients. For example, at the 97th percentile of risk—the threshold for allocating resources—Black patients had 26% more chronic illnesses than White patients. Due to this miscalibration, resources were diverted away from high-need Black patients to healthier White patients.

This example illustrates the problem of label bias. Healthcare costs are a poor proxy of healthcare needs given disparities in healthcare access and Medicaid enrollment that are correlated with race. Past work has shown that, conditional on need, healthcare spending is lower for Black individuals than for White individuals (47). Consequently, accurate prediction of healthcare costs necessarily leads to racially biased allocation of healthcare resources. Obermeyer et al. (48) estimated that changing the target of prediction to an index variable that incorporates health alongside cost prediction would lead to more resources being allocated to

Black patients. This result suggests that algorithmic label bias, at least in some circumstances, is both fixable and preventable by thoughtfully selecting prediction targets.

#### **IV. TOWARDS CONSEQUENTIALIST ALGORITHM DESIGN**

The algorithm case studies in the previous sections reveal problems with popular approaches to fairness, which often fail to consider the impact of decisions. These issues highlight the need for a design approach that foregrounds the consequences of an algorithm’s use—a challenge we take on here. To guide our discussion, we consider risk models for type 2 diabetes. Researchers have proposed using race-based models for estimating diabetes risk in an effort to address known racial disparities in diabetes diagnoses (1). Current guidelines recommend using a 1.5% threshold on estimated risk for recommending follow-up screening in the form of a blood test (1). Diabetes screening is not a resource constrained practice, and so, in line with our discussion above, we do not evaluate potential differences in decision rates (or error rates) across groups, as they tell us little about the fairness of an algorithm. Rather, we address the other two fairness concerns considered in this article: label bias and the inclusion (or exclusion) of race and ethnicity. Using a consequentialist framework, we show how to arbitrate between race-aware and race-unaware risk models, following Coats et al. (14).

For our analysis, we use data from the National Health and Nutrition Examination Survey (NHANES) (four cycles from 2011 to 2018) (9). We restricted our sample to approximately 18,000 patients who were not pregnant, were 18–70 years old, and had a body mass index (BMI) between 18.5 kg/m<sup>2</sup> and 50.0 kg/m<sup>2</sup>. Using this data, we trained two linear diabetes risk models based on age and BMI, and which differed only in their inclusion of race and ethnicity as a predictive variable. This example is for illustrative purposes only and we caution against using these models to guide clinical decisions.

##### **A. SELECTING APPROPRIATE PREDICTION TARGETS**

When designing or evaluating a model, it is imperative to scrutinize the prediction target for label bias. Researchers should be careful to consider the ways in which the proposed label may mismeasure the true outcome through systemic mechanisms of inequality, such as inequitable healthcare access. In our diabetes example, the label used to train the risk models was constructed by combining the results of a blood test administered by NHANES with the response



to a question on whether the respondent had ever been diagnosed with diabetes by a doctor—and so the label likely accurately captures disease status, the true outcome of interest. If, however, our label were only constructed on the individual’s response to the diagnosis question, our diabetes risk models would more closely predict diabetes diagnoses, as opposed to diabetes incidence. That misalignment between label and outcome could result in underestimating diabetes risk in groups that have less access to healthcare.

## **B. FOREGROUNDING UTILITY**

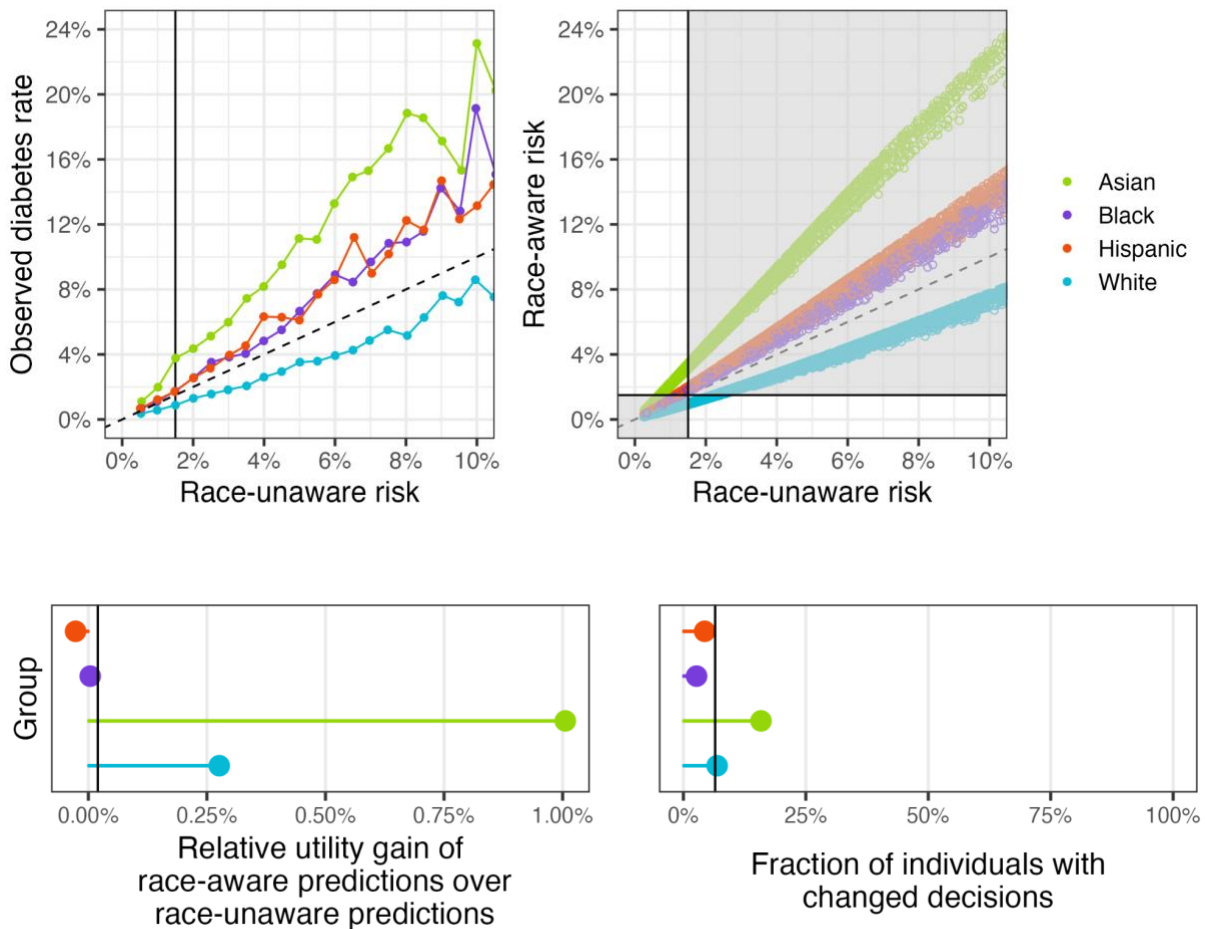
One of the major shortcomings of common algorithmic fairness arguments is that they do not consider individual utilities. In particular, many have advocated for using specific risk algorithms on the grounds that they increase decision rates (often for screening or treatment) for racial and ethnic minority groups (37, 18). Yet, these arguments often do not consider expected changes in utility from these changed decisions, and generally overlook the possibility that increasing decision rates can lead to utility *losses* for members of all groups—as we demonstrated with our lung cancer example. A consequentialist approach to algorithm design instead foregrounds individual utilities in fairness evaluations.

We first examine the accuracy of race-unaware risk models. In the top left plot of Figure 2, we visualize the estimated risk from a race-unaware model for diabetes against the observed rates of diabetes across racial and ethnic groups. This plot reveals discrepancies between the model predictions and true diabetes incidence for all groups. This miscalibration would consequently lead the model to fail to recommend screening for some high-risk Asian individuals—whose risk is underestimated by the model—and inappropriately recommend screening for some low-risk White individuals—whose risk is overestimated by the model. In both cases, the miscalibration would impose net costs on members of both groups from suboptimal screening recommendations. Past work has suggested using race-aware risk models to correct miscalibration across groups and obtain more accurate predictions (1). However, that work has stopped short of considering not just differences in accuracy, but differences in utility between race-aware and race-unaware models.

To compare the utility of race-aware versus race-unaware models, we follow Coots et al. (14) and first construct a simplified utility function to aid our comparison of risk models. Our utility function assumes a constant cost of screening and a constant benefit of correctly detecting

**Figure 2. A consequentialist approach to the design of a diabetes risk prediction model**

In line with analyses in Coats et al. (14), (top left) estimated risk from a race-unaware model for diabetes against the observed rates of diabetes across racial and ethnic groups. The vertical line corresponds to the recommended diabetes risk threshold of 1.5%, above which the typical patient can expect to benefit from a screening. The diagonal dashed lines represent hypothetical risk scores that are perfectly calibrated to empirical diabetes rates; (top right); scatter plots showing race-unaware risk plotted against race-aware risk for each individual in the data. Individuals in the shaded regions receive the same recommendation under both models; (bottom left) the relative gain in utility from the use of race-aware predictions to make screening recommendations across racial and ethnic groups; (bottom right) the fraction of individuals with different recommendations under race-unaware and race-aware models.



diabetes. (For further detail on the utility function, see Coots et al. (14).) Applying this utility function to the decisions produced under a race-aware and race-unaware model, we find that the relative gain in utility from using a race-aware model is smaller than expected in light of the substantial improvements in accuracy offered by the race-aware model over the race-unaware model. Relative to a baseline policy of no screening, we estimate that the race-aware model would improve overall utility by 0.2% over the race-unaware model. The lower left panel of Figure 2 shows that gains are similarly small across race groups. Our simplified utility function is for illustrative purposes only, to estimate the magnitude of the expected benefits, and is not intended to capture all the complex clinical considerations.

The modest utility gains stem from two factors. First, the vast majority of individuals (94%) would receive the same screening recommendation under both the race-aware and race-unaware models. In the top right plot of Figure 2, we plot the race-aware risk estimate for an individual against their race-unaware risk estimate. The dots in the shaded regions of the plot correspond to individuals for whom both models produce the same recommendation. In the lower right plot of Figure 2, we show the fraction of individuals with different decisions under each model by race and ethnicity. The second factor driving this result is that those individuals whose decisions do change are typically close to the decision threshold, and therefore accrue relatively small utility gains from the use of a race-aware model. In other words, the small subset of individuals with changed decisions should be largely ambivalent to being screened.

The race-unaware model is starkly miscalibrated but results in smaller than expected utility losses relative to the race-aware model. By foregrounding utility, our analysis helps clarify the expected benefits from using a race-aware model. However, the *costs* of using race—from risk of stigmatization or reinforcing pernicious attitudes on biological determinism, for example—remain an open question. Ultimately, these costs must be weighed against the estimated benefits in selecting the most appropriate model for decision-making.

#### IV. CONCLUSION

As algorithms are increasingly used to guide healthcare decisions, discussions around algorithmic fairness have come to the forefront, with racial equity being a particular focus of attention. By critically examining contemporary debates, we have argued for reframing what it means for an algorithm to be fair. Past fairness approaches—grounded in narrow summary

statistics, like decision rates and error rates—fail to anticipate the outcomes produced by algorithms, thereby risking unintended harm, including to those in racial and ethnic minority groups. With a consequentialist approach to algorithm design, we advocate for explicitly considering the utility of decisions produced by candidate algorithms to better understand the impact of design choices. This is no easy task. A consequentialist approach requires defining an appropriate utility function—a complex assignment that may also require aligning differing values across stakeholders. This challenge has led some scholars to critique consequentialist approaches to policy (53). But we believe that in many cases of practical importance, it is both feasible and useful to articulate one’s values, estimate the impacts of different algorithms, and confront the resulting trade-offs. We hope that our discussion helps researchers, clinicians, and policymakers better understand the common threads underlying ongoing debates, and illuminates a path forward for designing more equitable healthcare algorithms.

## **DATA AND CODE AVAILABILITY**

Data and code to reproduce our analysis are available at: <https://github.com/madisoncoots/racial-bias-in-healthcare-algs>.

## **LITERATURE CITED**

1. Aggarwal R, Bibbins-Domingo K, Yeh RW, et al. Diabetes screening by race and ethnicity in the United States: equivalent body mass index and age thresholds. *Annals of Internal Medicine*. 2022 Jun;175(6):765-73.
2. Aldrich, Melinda C, Sarah F Mercaldo, Kim L Sandler, William J Blot, Eric L Grogan, and Jeffrey D Blume. 2019. “Evaluation of USPSTF Lung Cancer Screening Guidelines Among African American Adult Smokers.” *JAMA Oncology*5 (9): 1318–24. <https://doi.org/10.1001/jamaoncol.2019.1402>.
3. Bahrami, Hossein, Richard Kronmal, David A Bluemke, Jean Olson, Steven Shea, Kiang Liu, Gregory L Burke, and João A. C Lima. 2008. “Differences in the Incidence of Congestive Heart Failure by Ethnicity: The Multi-Ethnic Study of Atherosclerosis.” *Archives of Internal Medicine* (1960) 168 (19): 2138–45. <https://doi.org/10.1001/archinte.168.19.2138>.

4. Bates DW, Saria S, Ohno-Machado L, Shah A, Escobar G. Big data in health care: using analytics to identify and manage high-risk and high-cost patients. *Health Aff (Millwood)*. 2014 Jul;33(7):1123-31. doi: 10.1377/hlthaff.2014.0041. PMID: 25006137.
5. Bibbins-Domingo, Kirsten, Mark J Pletcher, Feng Lin, Eric Vittinghoff, Julius M Gardin, Alexander Arynchyn, Cora E Lewis, O. Dale Williams, and Stephen B Hulley. 2009. “Racial Differences in Incident Heart Failure Among Young Adults.” *The New England Journal of Medicine* 360 (12): 1179–90. <https://doi.org/10.1056/NEJMoa0807265>.
6. Buckley, Ayisha, Stephanie Sestito, Tonia Ogundipe, Jacqueline Roig, Henri Mitchell Rosenberg, Natalie Cohen, Kelly Wang, et al. 2022. “Racial and Ethnic Disparities Among Women Undergoing a Trial of Labor After Cesarean Delivery: Performance of the VBAC Calculator with and Without Patients’ Race/Ethnicity.” *Reproductive Sciences (Thousand Oaks, Calif.)* 29 (7): 2030–38. <https://doi.org/10.1007/s43032-022-00959-2>.
7. Carnethon, Mercedes R, Jia Pu, George Howard, Michelle A Albert, Cheryl A.M Anderson, Alain G Bertoni, Mahasin S Mujahid, et al. 2017. “Cardiovascular Health in African Americans: A Scientific Statement From the American Heart Association.” *Circulation* 136 (21): e393–e423. <https://doi.org/10.1161/cir.0000000000000534>.
8. Centers for Disease Control and Prevention. Lung Cancer Statistics. US Department of Health and Human Services; 2023. <https://www.cdc.gov/cancer/lung/statistics/index.htm>
9. Centers for Disease Control and Prevention (CDC); National Center for Health Statistics (NCHS). National Health and Nutrition Examination Survey Data, 2011-2018. <https://www.cdc.gov/nchs/nhanes/index.htm>
10. Cheung LC, Berg CD, Castle PE, Katki HA, Chaturvedi AK. Life-gained-based versus risk-based selection of smokers for lung cancer screening. *Ann Intern Med*. 2019; 171(9):623-632. doi:10.7326/M19-1263
11. Chohlas-Wood A, Coots M, Goel S, et al. Designing equitable algorithms. *Nat Comput Sci*. 2023;3:601-610. doi: 10.1038/s43588-023-00485-4.
12. Chohlas-Wood A, Coots M, Zhu H, Brunskill E, Goel S. 2023. Learning to be fair: a consequentialist approach to equitable decision-making. arXiv:2109.08792v4 [cs.LG]
13. Cockcroft, Donald W., and Henry Gault. 1976. “Prediction of Creatinine Clearance from Serum Creatinine.” *Nephron* 16 (1): 31–41. <https://doi.org/10.1159/000180580>.

14. Coots, M., Saghafian, S., Kent, D. & Goel, S. A framework for considering the role of race and ethnicity in estimating disease risk. Under review (2024)
15. Curtin SC, Gregory KD, Korst LM, Uddin SF. Maternal Morbidity for Vaginal and Cesarean Deliveries, According to Previous Cesarean History: New Data From the Birth Certificate, 2013. *Natl Vital Stat Rep.* 2015 May 20;64(4):1-13, back cover. PMID: 26046963.
16. D'AGOSTINO, Ralph B, Ramachandran S VASAN, Michael J PENCINA, Philip A WOLF, Mark COBAIN, Joseph M MASSARO, and William B KANNEL. 2008. "General Cardiovascular Risk Profile for Use in Primary Care The Framingham Heart Study." *Circulation (New York, N.Y.)* 117 (6): 743–53. <https://doi.org/10.1161/CIRCULATIONAHA.107.699579>.
17. D'Agostino, Ralph B, Scott Grundy, Lisa M Sullivan, and Peter Wilson. 2001. "Validation of the Framingham Coronary Heart Disease Prediction Scores: Results of a Multiple Ethnic Groups Investigation." *JAMA : the Journal of the American Medical Association* 286 (2): 180–87. <https://doi.org/10.1001/jama.286.2.180>.
18. Diao, James A, Gloria J Wu, Herman A Taylor, John K Tucker, Neil R Powe, Isaac S Kohane, and Arjun K Manrai. 2021. "Clinical Implications of Removing Race From Estimates of Kidney Function." *JAMA : the Journal of the American Medical Association* 325 (2): 184–86. <https://doi.org/10.1001/jama.2020.22124>.
19. Diao JA, He Y, Khazanchi R, Nguemeni Tiako MJ, Witonsky JI, Pierson E, Rajpurkar P, Elhawary JR, Melas-Kyriazi L, Yen A, Martin AR, Levy S, Patel CJ, Farhat M, Borrell LN, Cho MH, Silverman EK, Burchard EG, Manrai AK. Implications of Race Adjustment in Lung-Function Equations. *N Engl J Med.* 2024 Jun 13;390(22):2083-2097. doi: 10.1056/NEJMsa2311809. Epub 2024 May 19. PMID: 38767252; PMCID: PMC11305821.
20. Diao JA, Shi I, Murthy VL, et al. Projected Changes in Statin and Antihypertensive Therapy Eligibility With the AHA PREVENT Cardiovascular Risk Equations. *JAMA.* Published online July 29, 2024. doi:10.1001/jama.2024.12537
21. Doubeni CA, Simon M, Krist AH. Addressing systemic racism through clinical preventive service recommendations from the US Preventive Services Task Force. *JAMA.* 2021;325(7):627-628. doi:10.1001/jama. 2020.26188

22. Eberly, Lauren A., Bram Wispelwey, Aaron Richterman, Anne G. Beckett, Emily C. Cleveland Manchanda, Regan H. Marsh, Cindy Y. Chang, et al. 2020. "Response by Eberly et Al to Letter Regarding Article, "Identification of Racial Inequities in Access to Specialized Inpatient Heart Failure Care at an Academic Medical Center"." *Circulation. Heart Failure* 13 (6): e007193–e007193.  
<https://doi.org/10.1161/CIRCHEARTFAILURE.120.007193>.
23. Eneanya, Nwamaka Denise, Wei Yang, and Peter Philip Reese. 2019. "Reconsidering the Consequences of Using Race to Estimate Kidney Function." *JAMA : the Journal of the American Medical Association* 322 (2): 113–14. <https://doi.org/10.1001/jama.2019.5774>.
24. Goff DC Jr, Lloyd-Jones DM, Bennett G, et al. 2013 ACC/AHA guideline on the assessment of cardiovascular risk: a report of the American College of Cardiology/American Heart Association Task Force on Practice Guidelines. *J Am Coll Cardiol.* 2014;63(25 Pt B):2935-2959. doi:10.1016/j.jacc.2013.11.005
25. Goff, David C, and Donald M Lloyd-Jones. 2016. "The Pooled Cohort Risk Equations—Black Risk Matters." *JAMA Cardiology* 1 (1): 12–14.  
<https://doi.org/10.1001/jamacardio.2015.0323>.
26. Grobman WA, Lai Y, Landon MB, et al. Does information available at admission for delivery improve prediction of vaginal birth after cesarean? *Am J Perinatol* 2009;26:693–701
27. GROBMAN, William A, Yinglei Lai, Ronald J WAPNER, Yoram SOROKIN, Menachem MIODOVNIK, Marshall CARPENTER, Mary J O'SULLIVAN, et al. 2007. "Development of a Nomogram for Prediction of Vaginal Birth after Cesarean Delivery." *Obstetrics and Gynecology (New York. 1953)* 109 (4): 806–12.  
<https://doi.org/10.1097/01.AOG.0000259312.36053.02>.
28. Grobman, William A., Grecio J. Sandoval, Madeline Murguia Rice, Suneet P. Chauhan, Rebecca G. Clifton, Maged M. Costantine, Kelly S. Gibson, et al. 2024. "Prediction of Vaginal Birth after Cesarean Using Information at Admission for Delivery: a Calculator Without Race or Ethnicity." *American Journal of Obstetrics and Gynecology* 230 (3): S804–S806. <https://doi.org/10.1016/j.ajog.2023.02.008>.
29. Grobman, William A., Grecio Sandoval, Madeline Murguia Rice, Jennifer L. Bailit, Suneet P. Chauhan, Maged M. Costantine, Cynthia Gyamfi-Bannerman, et al. 2021.

- “Prediction of Vaginal Birth after Cesarean Delivery in Term Gestations: a Calculator Without Race and Ethnicity.” *American Journal of Obstetrics and Gynecology* 225 (6): 664.e1–664.e7. <https://doi.org/10.1016/j.ajog.2021.05.021>.
30. Inker, Lesley A, Nwamaka D Eneanya, Josef Coresh, Hocine Tighiouart, Dan Wang, Yingying Sang, Deidra C Crews, et al. 2021. “New Creatinine- and Cystatin C–Based Equations to Estimate GFR without Race.” *New England Journal of Medicine* 385 (19): 1737–49. <https://doi.org/10.1056/nejmoa2102953>.
31. Katki HA, Kovalchik SA, Berg CD, et al. Development and Validation of Risk Models to Select Ever-Smokers for CT Lung Cancer Screening. *JAMA*. 2016;315(21):2300–2311. doi:10.1001/jama.2016.6255
32. Katki HA, Kovalchik SA, Petito LC, et al. Implications of Nine Risk Prediction Models for Selecting Ever-Smokers for Computed Tomography Lung Cancer Screening. *Ann Intern Med*. 2018 Jul 3;169(1):10-19. doi: 10.7326/M17-2701. Epub 2018 May 15. PMID: 29800127; PMCID: PMC6557386.
33. Khan, Sadiya S, Josef Coresh, Michael J Pencina, Chiadi E Ndumele, Janani Rangaswami, Sheryl L Chow, Latha P Palaniappan, et al. 2023. “Novel Prediction Equations for Absolute Risk Assessment of Total Cardiovascular Disease Incorporating Cardiovascular-Kidney-Metabolic Health: A Scientific Statement From the American Heart Association.” *Circulation (New York, N.Y.)* 148 (24): 1982–2004. <https://doi.org/10.1161/CIR.0000000000001191>.
34. Khan, Sadiya S, Kunihiro Matsushita, Yingying Sang, Shoshana H Ballew, Morgan E Grams, Aditya Surapaneni, Michael J Blaha, et al. 2024. “Development and Validation of the American Heart Association's PREVENT Equations.” *Circulation (New York, N.Y.)* 149 (6): 430–49. <https://doi.org/10.1161/CIRCULATIONAHA.123.067626>.
35. Krist AH, Davidson KW, Mangione CM, et al; US Preventive Services Task Force. Screening for lung cancer: US Preventive Services Task Force recommendation statement. *JAMA*. 2021;325(10):962-970. doi:10.1001/jama. 2021.1117
36. Landon, Mark B., Sharon Leindecker, Catherine Y. Spong, John C. Hauth, Steven Bloom, Michael W. Varner, Atef H. Moawad, et al. 2005. “The MFMU Cesarean Registry: Factors Affecting the Success of Trial of Labor after Previous Cesarean



- Delivery.” *American Journal of Obstetrics and Gynecology* 193 (3): 1016–23.  
<https://doi.org/10.1016/j.ajog.2005.05.066>.
37. Landy R, Gomez I, Caverly TJ, et al. Methods for Using Race and Ethnicity in Prediction Models for Lung Cancer Screening Eligibility. *JAMA Netw Open*. 2023 Sep;6(9):e2331155. doi: 10.1001/jamanetworkopen.2023.31155.
38. Lett E, Asabor E, Beltrán S, et al. Conceptualizing, Contextualizing, and Operationalizing Race in Quantitative Health Sciences Research. *Ann Fam Med*. 2022 Mar-Apr;20(2):157-163. doi: 10.1370/afm.2792. Epub 2022 Jan 19. PMID: 35045967; PMCID: PMC8959750.
39. Levey, A. S, J. P Bosch, J. B Lewis, T Greene, N Rogers, and D Roth. 1999. “A More Accurate Method To Estimate Glomerular Filtration Rate from Serum Creatinine: A New Prediction Equation.” *Annals of Internal Medicine* 130 (6): 461–70.  
<https://doi.org/10.7326/0003-4819-130-6-199903160-00002>.
40. LEVEY, Andrew S, Josef CORESH, Tom GREENE, Lesley A STEVENS, Yaping Zhang, Stephen HENDRIKSEN, John W KUSEK, and Frederick VAN LENTE. 2006. “Using Standardized Serum Creatinine Values in the Modification of Diet in Renal Disease Study Equation for Estimating Glomerular Filtration Rate.” *Annals of Internal Medicine* 145 (4): 247–54. <https://doi.org/10.7326/0003-4819-145-4-200608150-00004>.
41. LEVEY, Andrew S, Lesley A STEVENS, Josef CORESH, Christopher H SCHMID, Yaping Zhang, Alejandro F CASTRO, Harold I FELDMAN, et al. 2009. “A New Equation to Estimate Glomerular Filtration Rate.” *Annals of Internal Medicine* 150 (9): 604–12. <https://doi.org/10.7326/0003-4819-150-9-200905050-00006>.
42. Manski CF, Mullahy J, Venkataramani AS. Using measures of race to make clinical predictions: decision making, patient health, and fairness. *Proc Natl Acad Sci U S A*. 2023;120(35):1-e2303370120. doi: 10.1073/pnas.2303370120.
43. Martin, J., Hamilton, B., Osterman, M., Driscoll, A., & Drake, P. (2018). Births: Final data for 2017. *National Vital Statistics Reports*, 67(8), 34.
44. Mullainathan S, Obermeyer Z. Does Machine Learning Automate Moral Hazard and Error? *Am Econ Rev*. 2017 May;107(5):476-480. doi: 10.1257/aer.p20171084. PMID: 28781376; PMCID: PMC5540263.

45. Muntner, Paul, Lisandro D Colantonio, Mary Cushman, David C Goff, George Howard, Virginia J Howard, Brett Kissela, Emily B Levitan, Donald M Lloyd-Jones, and Monika M Safford. 2014. "Validation of the Atherosclerotic Cardiovascular Disease Pooled Cohort Risk Equations." *JAMA* 311 (14): 1406–15.  
<https://doi.org/10.1001/jama.2014.2630>.
46. National Lung Screening Trial Research Team; Aberle DR, Berg CD, Black WC, et al. The National Lung Screening Trial: overview and study design. *Radiology*. 2011 Jan;258(1):243-53. doi: 10.1148/radiol.10091808. Epub 2010 Nov 2. PMID: 21045183; PMCID: PMC3009383.
47. National Research Council (US) Committee on Population. Racial and Ethnic Differences in the Health of Older Americans. Martin LG, Soldo BJ, editors. Washington (DC): National Academies Press (US); 1997. PMID: 23115810.
48. Obermeyer et al. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366,447-453(2019).DOI:10.1126/science.aax2342
49. Pasquinelli, Mary M., Martin C. Tammemägi, Kevin L. Kovitz, Marianne L. Durham, Zanë Deliu, Kayleigh Rygalski, Li Liu, Matthew Koshy, Patricia Finn, and Lawrence E. Feldman. 2021. "Brief Report: Risk Prediction Model Versus United States Preventive Services Task Force 2020 Draft Lung Cancer Screening Eligibility Criteria—Reducing Race Disparities." *JTO Clinical and Research Reports* 2 (3): 100137–100137.  
<https://doi.org/10.1016/j.jtocrr.2020.100137>.
50. Pauker SG, Kassirer JP. The threshold approach to clinical decision making. *N Engl J Med*. 1980 May 15;302(20):1109-17. doi: 10.1056/NEJM198005153022003. PMID: 7366635.
51. Peterson PN, Rumsfeld JS, Liang L, et al. A validated risk score for in-hospital mortality in patients with heart failure from the American Heart Association Get with the Guidelines program. *Circ Cardiovasc Qual Outcomes* 2010;3:25-32.
52. Rivera, M. Patricia, Hormuzd A. Katki, Nichole T. Tanner, Matthew Triplette, Lori C. Sakoda, Renda Soylemez Wiener, Roberto Cardarelli, et al. 2020. "Addressing Disparities in Lung Cancer Screening Eligibility and Healthcare Access. An Official American Thoracic Society Statement." *American Journal of Respiratory and Critical Care Medicine* 202 (7): e95–112. <https://doi.org/10.1164/rccm.202008-3053st>.

53. Sandel, Michael J. (ed.) (2009). *Justice: What's the Right Thing to Do?*. New York: Farrar, Straus and Giroux.
54. Schwandt, Hannes, Janet Currie, James Banks, Marlies Bär, Paola Bertoli, Aline Bütikofer, Sarah Cattan, et al. 2021. "Inequality in Mortality Between Black and White Americans by Age, Place, and Cause and in Comparison to Europe, 1990 to 2018." *IDEAS Working Paper Series from RePEc* 118 (40): 1. <https://doi.org/10.1073/pnas.2104684118>.
55. SEER\*Explorer: An interactive website for SEER cancer statistics [Internet]. Surveillance Research Program, National Cancer Institute; 2024 Apr 17. [cited 2024 May 8]. Available from: <https://seer.cancer.gov/statistics-network/explorer/>. Data source(s): SEER Incidence Data, November 2023 Submission (1975-2021), SEER 22 registries.
56. SEER\*Explorer: An interactive website for SEER cancer statistics [Internet]. Surveillance Research Program, National Cancer Institute; 2024 Apr 17. [cited 2024 May 8]. Available from: <https://seer.cancer.gov/statistics-network/explorer/>. Data source(s): U.S. Mortality Data (1969-2022), National Center for Health Statistics, CDC.
57. Shaikh N, Lee MC, Stokes LR, et al. Reassessment of the Role of Race in Calculating the Risk for Urinary Tract Infection: A Systematic Review and Meta-analysis. *JAMA Pediatr.* 2022;176(6):569–575. doi:10.1001/jamapediatrics.2022.0700
58. Smith, Wally R., Roy M. Poses, Donna K. McClish, Elizabeth C. Huber, F. Lynne W. Clemo, Donna Alexander, and Brian P. Schmitt. 2002. "Prognostic Judgments and Triage Decisions for Patients With Acute Congestive Heart Failure." *Chest* 121 (5): 1610–17. <https://doi.org/10.1378/chest.121.5.1610>.
59. Stevens, Elizabeth R., Tanner Caverly, Jorie M. Butler, Polina Kukhareva, Safiya Richardson, Devin M. Mann, and Kensaku Kawamoto. 2023. "Considerations for Using Predictive Models That Include Race as an Input Variable: The Case Study of Lung Cancer Screening." *Journal of Biomedical Informatics* 147: 104525–104525. <https://doi.org/10.1016/j.jbi.2023.104525>.
60. Stevens, Lesley A, Josef Coresh, Tom Greene, and Andrew S Levey. 2006. "Assessing Kidney Function — Measured and Estimated Glomerular Filtration Rate." *New England Journal of Medicine* 354 (23): 2473–83. <https://doi.org/10.1056/nejmra054415>.

61. Tammemägi MC, Katki HA, Hocking WG, Church TR, Caporaso N, Kvale PA, Chaturvedi AK, Silvestri GA, Riley TL, Commins J, Berg CD. Selection criteria for lung-cancer screening. *N Engl J Med*. 2013 Feb 21;368(8):728-36. doi: 10.1056/NEJMoa1211776. Erratum in: *N Engl J Med*. 2013 Jul 25;369(4):394. PMID: 23425165; PMCID: PMC3929969.
62. Vyas, Darshali A, Aisha James, William Kormos, and Utibe R Essien. 2022. "Revising the Atherosclerotic Cardiovascular Disease Calculator Without Race." *The Lancet Digital Health* 4 (1): e4–e5. [https://doi.org/10.1016/S2589-7500\(21\)00258-2](https://doi.org/10.1016/S2589-7500(21)00258-2).
63. Vyas, Darshali A, Leo G Eisenstein, and David S Jones. 2020. "Hidden in Plain Sight — Reconsidering the Use of Race Correction in Clinical Algorithms." *The New England Journal of Medicine* 383 (9): 874–82. <https://doi.org/10.1056/NEJMms2004740>.
64. Vyas, Darshali A., David S. Jones, Audra R. Meadows, Khady Diouf, Nawal M. Nour, and Julianna Schantz-Dunn. 2019. "Challenging the Use of Race in the Vaginal Birth after Cesarean Section Calculator." *Women's Health Issues* 29 (3): 201–4. <https://doi.org/10.1016/j.whi.2019.04.007>.
65. Wang JH, Hart A. Global Perspective on Kidney Transplantation: United States. *Kidney360*. 2021 Aug 19;2(11):1836-1839. doi: 10.34067/KID.0002472021. PMID: 35373000; PMCID: PMC8785833.
66. Wilson, P. W. F, R. B D'Agostino, D Levy, A. M Belanger, H Silbershatz, and W. B Kannel. 1998. "Prediction of Coronary Heart Disease Using Risk Factor Categories." *Circulation* 97 (18): 1837–47. <https://doi.org/10.1161/01.cir.97.18.1837>.
67. Wispelwey, Bram, Cindy Y Chang, Katherine C Brooks, Rose Kakoza, Joseph Loscalzo, Eldrin F Lewis, Viswatej Avutu, et al. 2019. "Identification of Racial Inequities in Access to Specialized Inpatient Heart Failure Care at an Academic Medical Center." *Circulation. Heart Failure* 12 (11): e006214–e006214. <https://doi.org/10.1161/CIRCHEARTFAILURE.119.006214>.
68. Zanger-Tishler et al. Risk scores, label bias, and everything but the kitchen sink. *Sci Adv*. 10,eadi8411(2024).DOI:10.1126/sciadv.adi8411
69. Zink, A., Obermeyer, Z., & Pierson, E. (2024). Race adjustments in clinical algorithms can help correct for racial disparities in data quality. *Proceedings of the National Academy of Sciences*, 121(34), e2402267121.