

# Net benefit, calibration, threshold selection, and training objectives for algorithmic fairness in healthcare

Stephen R. Pfohl<sup>\*,1</sup>, Yizhe Xu<sup>1</sup>, Agata Foryciarz<sup>2</sup>,  
Nikolaos Ignatiadis<sup>3</sup>, Julian Genkins<sup>1</sup>, Nigam H. Shah<sup>1</sup>

<sup>1</sup>Stanford Center for Biomedical Informatics Research, Stanford University, Stanford, California 94305, USA

<sup>2</sup>Department of Computer Science, Stanford University, Stanford, California 94305, USA

<sup>3</sup>Department of Statistics, Stanford University, Stanford, California 94305, USA

\*Correspondence to: [spfohl@stanford.edu](mailto:spfohl@stanford.edu)

## Abstract

A growing body of work uses the paradigm of algorithmic fairness to frame the development of techniques to anticipate and proactively mitigate the introduction or exacerbation of health inequities that may follow from the use of model-guided decision-making. We evaluate the interplay between measures of model performance, fairness, and the expected utility of decision-making to offer practical recommendations for the operationalization of algorithmic fairness principles for the development and evaluation of predictive models in healthcare. We conduct an empirical case-study via development of models to estimate the ten-year risk of atherosclerotic cardiovascular disease to inform statin initiation in accordance with clinical practice guidelines. We demonstrate that approaches that incorporate fairness considerations into the model training objective typically do not improve model performance or confer greater net benefit for any of the studied patient populations compared to the use of standard learning paradigms followed by threshold selection concordant with patient preferences, evidence of intervention effectiveness, and model calibration. These results hold when the measured outcomes are not subject to differential measurement error across patient populations and threshold selection is unconstrained, regardless of whether differences in model performance metrics, such as in true and false positive error rates, are present. In closing, we argue for focusing model development efforts on developing calibrated models that predict outcomes well for all patient populations while emphasizing that such efforts are complementary to transparent reporting, participatory design, and reasoning about the impact of model-informed interventions in context.

## 1 Introduction

The use of machine learning to guide clinical decision-making and resource allocation can introduce or perpetuate inequities in care access and quality, ultimately contributing to health disparities [1, 2]. Aiming to detect and mitigate such harms, recent works leverage the *algorithmic fairness* paradigm [3] to define evaluation criteria and model development procedures that quantify and constrain the magnitude of statistical differences in model behavior or performance across patient subgroups [4–13]. Within this paradigm, numerous criteria, metrics, and algorithms have been proposed, and both major and minor incompatibilities and trade-offs among them have been identified [14–19].

The purpose of this work is to synthesize, contextualize, and validate underappreciated limitations of the algorithmic fairness paradigm to contribute to the development of best practices for appropriately operationalizing algorithmic fairness principles in healthcare [20]. We do so in a setting where observational data stored in an electronic health records or claims database is used to fit a patient-level predictive model for a clinical outcome where the score output by the model informs the allocation of a clinical intervention, typically through comparison of the score to a decision threshold. For our analysis, we assume that the observed outcomes are not subject to unobserved differential measurement error across patient subgroups [1, 21], that the choice of decision threshold used to allocate a clinical intervention on the basis of the output

of a predictive model is not constrained by resource or operational constraints [22], and that the values embedded in the data collection and problem formulation processes are transparently reported and reflect those of the patient populations affected by the use of the model [23–28].

For the development of predictive models to inform clinical decision-making, we argue for aiming to maximize the expected utility that the model-informed intervention confers to each patient subgroup of interest. The notion of expected utility that we consider depends on the values and preferences of affected stakeholders and can be quantified in terms of the expected costs or utilities associated with false positive and false negative errors in binary classification settings or in terms of the expected benefits and harms of the intervention conditioned on risk in more general settings [29, 30]. We hypothesize that, in practice, model development strategies that nominally promote fairness, by constraining for parity in model performance metrics across subgroups or by maximizing worst-case model performance over subgroups, do not confer greater expected utility for any patient subgroup than the approach of identifying a set of calibrated models that predict the outcome well for each subgroup, followed by threshold selection reflecting the contextual assessment of the benefits and harms of the intervention. The key observations motivating this hypothesis are detailed in section 2 and largely follow directly from related work [14–16, 18, 19, 29–36].

We evaluate our hypothesis through a case study of estimators of the risk of atherosclerotic cardiovascular disease (ASCVD) within ten years to inform the initiation of cholesterol-lowering statin therapy [37–41]. We conduct experiments to assess which model development strategies confer maximal expected utility for subgroups defined in terms of race, ethnicity, sex, or co-morbidities (type 1 and type 2 diabetes, chronic kidney disease (CKD), or rheumatoid arthritis (RA)). We compare pooled and stratified unconstrained empirical risk minimization (ERM) to regularized fairness objectives and distributionally robust optimization (DRO) objectives that aim to minimize differences in or improve the worst-case area under the receiver operating characteristic curve (AUC) or log-loss across subgroups. We further conduct an analysis to investigate the impact of constraints on differences in true and false positive rates. To evaluate the utility that the model confers, we use the notion of *net benefit* [30, 34] to define normalized expected utility measures that parameterize the relative value of the harms and benefits of statin initiation on the basis of decision thresholds recommended by clinical practice guidelines. To evaluate net benefit in this setting, we adopt the assumption that the intervention induces constant relative risk reduction (section 3.2.2). Furthermore, we use an inverse probability of censoring weighting (IPCW) approach to extend each of the training objectives and evaluation metrics used to account for censoring in ten-year ASCVD outcomes.

## 2 Background and problem formulation

### 2.1 Supervised learning for binary outcomes

Here, we introduce the formal notation and key assumptions used throughout the work. Let  $X \in \mathcal{X} = \mathbb{R}^m$  be a variable designating a vector of covariates and  $Y \in \mathcal{Y} = \{0, 1\}$  be a binary indicator of an outcome. We consider data that may be partitioned on the basis of a discrete indicator of a categorical attribute  $A \in \mathcal{A} = \{A_k\}_{k=1}^K$  with  $K$  categories. In some cases,  $A$  may correspond to an attribute that describes partitions of the population, where the value of  $A = A_k$  refers to a specific partition defined by the attribute. Examples of attributes used to partition the population include demographic attributes (*e.g.* race, ethnicity, gender, sex, age subgroup) or strata defined by complex clinical phenotypes or comorbidity profiles. We use the shorthand  $\mathcal{D}_{A_k}$ , when referring to the subset of the data  $\mathcal{D}$  corresponding to the subgroup  $A_k$ .

The objective of supervised learning with binary outcomes is to use data  $\mathcal{D} = \{(x_i, y_i, a_i)\}_{i=1}^N \sim P(X, Y, A)$  to learn a function  $f_\theta \in \mathcal{F} : \mathbb{R}^m \rightarrow [0, 1]$  parameterized by  $\theta$ . The function  $f_\theta$  can be considered to be a risk estimator that, when optimal, estimates  $\mathbb{E}[Y | X] = P(Y = 1 | X)$ . We designate the random variable resulting from the application of the model  $f_\theta$  to  $X$  to be given by  $S$ , such that  $S = f_\theta(X)$ . Given  $S$ , a predictor  $\hat{Y}$  may be derived by comparing  $S$  to a threshold  $\tau_y \in [0, 1]$  to produce binary predictions  $\hat{Y}(X) = \mathbb{1}[f_\theta(X) \geq \tau_y] \in \{0, 1\}$ .

The *calibration curve*  $c : [0, 1] \rightarrow [0, 1]$  is defined as a function that describes the expected value of  $Y$  given  $S$ , such that  $c(s) = E[Y | S = s] = P(Y = 1 | S = s)$ . A model is said to be calibrated if  $c(s) = s$  for all  $s$ . The calibration curve can be used to assess the extent to which a model over or underestimates the risk of the outcome  $Y$ . For instance, if  $c(s') > s'$  then the observed event rate for the set of patients with scores of

$s'$  is greater than  $s$ , implying that the model underestimates risk for patients with scores of  $s'$ . Analogously,  $c(s') < s'$  implies overestimation of risk for patients with scores of  $s'$ .

## 2.2 Algorithmic fairness criteria

Assessments of algorithmic fairness rely on *fairness criteria*, *i.e.* statistical properties reflecting moral or normative judgements as to the principles that constitute fairness. A broad class of fairness criteria can be described in terms of *metric parity* ( $g_j(\cdot) \perp A$ ), which requires that one or more metrics  $g_j : \mathcal{F} \times (\mathcal{X}, \mathcal{Y}) \rightarrow \mathbb{R}^+$  be equal across the subgroups defined by  $A$ . Common instantiations of metric parity include *equalized odds* ( $\hat{Y} \perp A \mid Y$  or  $S \perp A \mid Y$ ) [42], which requires both the true positive rates and the false positive rates to be equal across subgroups, *demographic parity* ( $\hat{Y} \perp A$  or  $S \perp A$ ) [43], which requires the rate at which patients are classified as belonging to the positive class is equal across subgroups, *predictive parity* ( $Y \perp A \mid \hat{Y} = 1$ ) [17], which requires parity in the positive predictive values, as well as criteria defined over other performance metrics [44, 45], including the AUC [46, 47] or the average log-loss or empirical risk [48]. Another important class of fairness criteria is defined over the calibration curve. Within that class, we focus on the *sufficiency* condition ( $Y \perp A \mid S$ ) [3, 19], which requires the calibration curves for each subgroup be equal, and the *group calibration* condition ( $\mathbb{E}[Y \mid S = s, A] = s$ ) [16, 49], which requires the model to be calibrated for each subgroup.

## 2.3 Assessing the utility and net benefit of decision-making at a threshold

To contextualize the presentation of algorithmic fairness, we present a utility-theoretic perspective on clinical decision-making. For this framing, we consider a decision rule that implies intervention allocation on the basis of a binary predictor  $\hat{Y}(X) = \mathbb{1}[f_\theta(X) \geq \tau_y]$ .

We define

$$U_{\text{cond}}(s) = U_{\text{cond}}^1(s) - U_{\text{cond}}^0(s) \quad (1)$$

as the *conditional* expected utility of the decision rule, where  $U_{\text{cond}}^1(s)$  designates the expected utility associated with treating patients whose predicted scores  $S = f_\theta(X)$  are  $s$ , and  $U_{\text{cond}}^0(s)$  is the expected utility of *not* treating patients whose scores are  $s$ . We define the *aggregate* expected utility  $U_{\text{agg}}(\tau_y)$  of the decision to be the average utility over the population, given that the intervention is allocated for all patients with scores at or above the threshold  $\tau_y$ :

$$U_{\text{agg}}(\tau_y) = \mathbb{E}[U_{\text{cond}}^1 \mid S \geq \tau_y]P(S \geq \tau_y) + \mathbb{E}[U_{\text{cond}}^0 \mid S < \tau_y]P(S < \tau_y). \quad (2)$$

The optimal decision rule for a fixed predictive model is one where the intervention is allocated to patients with scores for which  $U_{\text{cond}}(s) > 0$  and not allocated to those for which  $U_{\text{cond}}(s) < 0$ . If  $U_{\text{cond}}(s)$  is strictly monotonically increasing in  $s$  and has a root in  $[0, 1]$  then the optimal threshold  $\tau_y^*$  is given by the point at which  $U_{\text{cond}}(s = \tau_y^*) = 0$ . When  $U_{\text{cond}}(s)$  is strictly monotonic but has no root in  $[0, 1]$ , then either the treat-all ( $\tau_y = 0$ ) or treat-none ( $\tau_y = 1$ ) strategies is optimal.

In some cases,  $U_{\text{cond}}$  can be written as a simple function of the calibration curve. For example, if the costs and benefits of decision-making can be written as fixed expected costs or utilities of true positive ( $u_{\text{TP}}$ ), false positive ( $u_{\text{FP}}$ ), true negative ( $u_{\text{TN}}$ ), and false negative ( $u_{\text{FN}}$ ) classification, then

$$U_{\text{cond}}(s) = (u_{\text{TP}} - u_{\text{FN}})c(s) + (u_{\text{FP}} - u_{\text{TN}})(1 - c(s)) \quad (3)$$

and the optimal threshold is given by [18, 29]

$$\tau_y^* = c^{-1}\left(\frac{u_{\text{TN}} - u_{\text{FP}}}{u_{\text{TN}} - u_{\text{FP}} + u_{\text{TP}} - u_{\text{FN}}}\right). \quad (4)$$

It follows that when a model is calibrated, the optimal threshold is given by  $\tau_y' = \frac{u_{\text{TN}} - u_{\text{FP}}}{u_{\text{TN}} - u_{\text{FP}} + u_{\text{TP}} - u_{\text{FN}}}$ . When the model is miscalibrated, but the calibration curve is strictly monotonic, the optimal threshold is given the point at which the calibration curve intersects  $\tau_y'$ . Furthermore, given the relationship between the  $c(s)$  and  $U_{\text{cond}}$ , monotonicity in the calibration curve implies monotonicity in the conditional utility, and setting a threshold on the basis of the calibration curve can be interpreted as setting a threshold on  $U_{\text{cond}}$ .

To assess the expected utility of the decision rule over a population, it is typically not necessary to evaluate  $U_{\text{agg}}(\tau_y)$  with equation (2). Instead, a chosen decision threshold can be used to parameterize the *net benefit* [30, 34] of the decision rule under the assumption that the chosen threshold is optimal, for a calibrated model, based on the values of the decision maker and the effectiveness of the intervention. The net benefit under the assumption of fixed costs or utilities of classification errors is given by [30, 34]

$$\text{NB}(\tau_y; \tau_y^*) = P(S \geq \tau_y | Y = 1)P(Y = 1) - P(S \geq \tau_y | Y = 0)P(Y = 0) \frac{\tau_y^*}{1 - \tau_y^*}, \quad (5)$$

where  $\tau_y$  is the evaluated decision threshold and  $\tau_y^*$  parameterizes the net benefit. This metric is fundamental to *decision curve analysis* [30, 34], as a decision curve is the curve that results from evaluating net benefit for a range of thresholds for which  $\tau_y = \tau_y^*$ . Both the net benefit and  $U_{\text{agg}}$  are maximized at the threshold that results from the application of equation (4) when the assumptions outlined above are met.

We introduce the notion of the *calibrated net benefit* (cNB) to assess the net benefit under the assumption that the decision threshold used is adjusted on the basis of observed miscalibration. If  $c(s)$  is the calibration curve, then the calibrated net benefit evaluated at a threshold  $\tau_y$  is given by the net benefit evaluated at a threshold  $\tau_c = c^{-1}(\tau_y)$  on the score  $S$ . The calibrated net benefit under the assumption of fixed classification costs is given by

$$\text{cNB}(\tau_y; \tau_y^*) = P(S \geq c^{-1}(\tau_y) | Y = 1)P(Y = 1) - P(S \geq c^{-1}(\tau_y) | Y = 0)P(Y = 0) \frac{\tau_y^*}{1 - \tau_y^*}. \quad (6)$$

## 2.4 Implications for algorithmic fairness

A key consequence of the analysis presented thus far is that, subject to the assumptions detailed in section 2.3, the optimal threshold rule applied to a predictive model that outputs a continuous-valued risk score is based directly on the calibration characteristics of the model and the assumed expected costs or utilities of classification errors that encapsulate the effectiveness of the intervention and the preferences for downstream benefits and harms. As has been argued in related work [15, 18, 31, 50], it follows that if the model is calibrated for each subgroup, the decision threshold that maximizes expected utility and net benefit for each subgroup is the same when the expected utilities associated with each classification error do not change across subgroups. We verify this claim in simulation in supplementary section A1 (Supplementary Figure A1). Furthermore, in this case, sufficiency implies that the use of a consistent threshold on the risk score for all subgroups corresponds to the use of a consistent threshold on the conditional utility  $U_{\text{cond}}$  across subgroups, corresponding to an intuitive notion of *fairness* even in the case that the chosen decision threshold is not necessarily optimal [18, 31]. However, we note that this can still be a misleading notion of fairness given that it does not account for heterogeneity in the outcome not accounted for by the model under consideration [18].

As is described in prior work [16–19, 32], one should expect models that minimize the empirical risk for the population overall, with respect to a data distribution containing features  $X$  that encode  $A$ , to satisfy fairness criteria defined in terms of the calibration curve, including sufficiency and group calibration, and to violate equalized odds, demographic parity, and predictive parity when such models exhibit differences in the distribution of the risk score  $S$  or in the prevalence or incidence of the outcome  $Y$ . As discussed in Liu et al. [19], the use of ERM with a sufficiently large training dataset drawn from the target population is consistent with learning a model satisfying that criteria, implying that such models are expected to be calibrated overall and for each patient subgroup, but are expected to have non-trivial differences in true and false positive error rates, as well as in other performance metrics, when the incidence of the outcome or distribution of risk varies over those subgroups [16, 19, 32].

Consequently, approaches undertaken to constrain the model training objective [14, 44, 45, 51–53] to minimize violation of fairness criteria such as equalized odds or demographic parity typically reduce utility through some combination of explicit threshold adjustment [42] towards a threshold unrelated to the one selected on the basis of preference solicitation in the context of the intervention, induced miscalibration that analogously implies decision-making at a threshold unrelated to the utility-maximizing one [15], or reduction in model fit [14]. Given the relationship between the calibration curve and the conditional utility described in section 2.3, induced miscalibration that results in sufficiency violation implies that the use of a consistent threshold on the score across subgroups results in the use of different thresholds on  $U_{\text{cond}}$  across subgroups.

## 2.5 Algorithmic fairness training objectives

We evaluate training objectives that incorporate algorithmic fairness goals and constraints into their specification. We do so not to advocate for the use of their use, but rather to develop evidence as to the extent to which theoretical properties and trade-offs manifest empirically. We focus our efforts on “in-processing” approaches [44, 45, 51–53] rather than on pre- [54–58] or post-processing [3, 42] (*e.g.* threshold-adjustment) approaches because in-processing approaches are well-suited to learning models that achieve the minimum achievable trade-off between measures of model performance and fairness in practical finite-sample settings [59] and further allow for exploration of smooth trade-offs induced by relaxation of the constraint [45, 51]. We specifically focus on scalable gradient-based learning procedures that use regularized objectives to penalize violation of fairness criteria in a minibatch setting, to enable the use of these procedures for deep neural network models learned with large-scale datasets. We investigate approaches that, rather than constraining for parity in a metric across subgroups, attempts to improve the worst-case value of the metric over subgroups using distributionally robust optimization (DRO) [36, 60–63].

Following Pfohl et al. [14], the regularized training objective is ERM that incorporates a non-negative penalty term  $R$  that assesses the extent to which a fairness criterion of interest is violated and a non-negative parameter  $\lambda$  that may be tuned to control the extent to which violation of the criteria is penalized:

$$\min_{\theta \in \Theta} \sum_{i=1}^N w_i \ell(y_i, f_{\theta}(x_i)) + \lambda R, \quad (7)$$

where  $w_i$  are sample weights. In our experiments, we use this formulation to penalize violation of equalized odds and differences in AUC and log-loss across subgroups. To penalize violation of equalized odds, we primarily use a term that penalizes the Maximum Mean Discrepancy (MMD) [64] between the distribution of scores between each patient subgroup and the overall population conditioned on the observed values of the outcome  $Y$ , as in Pfohl et al. [14]:

$$\min_{\theta \in \Theta} \sum_{i=1}^N w_i \ell(y, f_{\theta}(x)) + \lambda \frac{1}{K} \sum_{Y_j \in \mathcal{Y}} \sum_{A_k \in \mathcal{A}} \hat{D}_{\text{MMD}}(P(f(X) | A = A_k, Y = Y_j) || P(f(X) | Y = Y_j)). \quad (8)$$

A full specification of the MMD-based training objective is included in supplementary section A.3.

We further use a regularized objective defined on the basis of a penalty that assesses violation of metric parity to penalize differences in the AUC or log-loss between each subgroup with the overall population:

$$\min_{\theta \in \Theta} \sum_{i=1}^N w_i \ell(y, f_{\theta}(x)) + \lambda \sum_{j=1}^J \sum_{A_k \in \mathcal{A}} (g_j(f_{\theta}, \mathcal{D}_{A_k}) - g_j(f_{\theta}, \mathcal{D}))^2. \quad (9)$$

We also evaluate the use of this objective to penalize violation of equalized odds at relevant thresholds by plugging surrogates of the true and false positive rates into equation (9). A full specification of the relevant objectives is provided in supplementary section A.3.

Beyond regularized objectives for algorithmic fairness, we evaluate distributionally robust optimization [60, 65, 66] procedures that encode the goal of maximizing worst-case performance over subgroups as one of learning to be robust over marginal shifts in the proportion of data available from each subgroup. The use of these objectives reflects a shift in perspective from the goal of requiring that some statistic be equal across subgroups towards one of aiming to identify models that perform well for each subgroup [36, 60–62, 66]. In this work, we leverage the *GroupDRO* framework (hereafter referred to as DRO) developed in Sagawa et al. [60] and extended in Pfohl et al. [36]. The algorithm is implemented as the following alternating updates conducted over minibatches:

$$\lambda_k \leftarrow \lambda_k \exp(\eta g(f_{\theta}, \mathcal{D}_{A_k})) / \sum_{k=1}^K \exp(\eta g(f_{\theta}, \mathcal{D}_{A_k})) \quad (10)$$

and

$$\min_{\theta \in \Theta} \sum_{k=1}^K \lambda_k \sum_{i=1}^{n_k} w_i \ell(y_i, f_{\theta}(x_i)), \quad (11)$$

where  $\eta$  is a non-negative scalar hyperparameter,  $\{\lambda_k\}_{k=1}^K$  are non-negative scalars that sum to 1, and  $g$  is a performance metric where lower values of the metric indicate better performance. In our experiments, we evaluate the use of the log-loss and  $1 - \text{AUC}$  as the choice of metric  $g$ , as in Pfohl et al. [36].

### 3 Case study in atherosclerotic disease risk estimation

#### 3.1 Background on ASCVD risk estimation for statin initiation

Clinical practice guidelines for the primary prevention of cardiovascular disease recommend the use of estimates of ten-year atherosclerotic cardiovascular disease (ASCVD) risk to inform the initiation of cholesterol-lowering statin therapy [37–41]. These guidelines primarily recommend the use of risk estimates provided by the Pooled Cohort Equations [37] and its extensions [67]. However, these estimates have been reported to systematically over-estimate or under-estimate risk in ways that are consequential for the appropriateness of downstream treatment decisions. This misestimation has been reported to occur both overall [68–71] and for subgroups defined on the basis of race/ethnicity [72–74], sex [68, 69, 75], socioeconomic status [41], or for patients with comorbidities which influence ASCVD risk or the expected benefit and harms of statin therapy, including diabetes [71, 74], chronic kidney disease (CKD) [74, 76, 77], and rheumatoid arthritis (RA) [78, 79]. Approaches undertaken to address these issues include the development of new risk estimators from large, diverse observational cohorts using modern machine learning methods [10, 67, 80–82], revisions to guidelines to encourage follow-up testing when the benefits of statin therapy are unclear and shared patient-clinician decision-making to incorporate patient preferences and other context [41], and the incorporation of fairness constraints into the model development process [9, 10, 15].

#### 3.2 Extending the approach

##### 3.2.1 Supervised learning with censored binary outcomes

When the binary outcome  $Y = \mathbb{1}[T \leq \tau_t]$  is defined as the occurrence of the outcome event at a time  $T$  at or before a time horizon  $\tau_t \in \mathbb{R}^+$ , it is important to account for censoring. The presence of censoring implies that either the outcome event time  $T \in \mathbb{R}^+$  or the censoring time  $C \in \mathbb{R}^+$  will be observed, but not both. The outcome data in an observed dataset  $\mathcal{D} = \{(x_i, u_i^t, \delta_i^t, a_i)\}_{i=1}^N$  is represented by an observed follow-up time  $U^t = \min(T, C)$  and an indicator  $\Delta^t = \mathbb{1}[T \leq C]$  that reflects whether the observed follow-up time corresponds to an outcome or a censoring event. The binary outcome  $Y$  is said to be censored if the censoring time  $C$  occurs prior to both the observed follow-up time and the time horizon, *i.e.*  $C < T$  and  $C < \tau_t$ . We define a composite observed follow-up time  $U^y = \min(T, C, \tau_t)$  for the binary outcome and an indicator  $\Delta^y = 1 - \mathbb{1}[C < T] * \mathbb{1}[C < \tau_t]$  that reflects whether a patient’s binary outcome is uncensored. A visual depiction of the relationship between the outcome and censoring event times and the value and censoring status of the binary outcome is shown in supplementary figure C1.

The use of inverse probability of censoring weighting (IPCW) allows for the derivation and evaluation of predictive models for censored binary outcomes [83–87], analogous to propensity score weighting procedures used for causal effect estimation [88]. The appropriate weights are those that are proportional to the inverse probability of remaining uncensored at the time of the composite observed follow-up time. Specifically, for an estimate of the censoring survival function  $G(s, x) = P(C > s \mid X = x)$  we define normalized weights

$$w_i = \frac{\delta_i^y}{G(u_i^y, x_i)} \left( \sum_{i=1}^N \frac{\delta_i^y}{G(u_i^y, x_i)} \right)^{-1} \quad (12)$$

that for a patient  $i$  reflects the reciprocal of the conditional probability of remaining uncensored at the time  $u_i^y$  given features  $x_i$ . To enable this approach, we make the following assumptions: (1) *coarsening at random* [83, 89] where the outcome event time is independent of the censoring time conditioned on the features, *i.e.*  $T \perp C \mid X$ , (2)  $G(U, X) > 0$  for all data with uncensored binary outcomes (for which  $\Delta^y = 1$ ), and (3) that  $f_\theta$  is a deterministic transformation.

The IPCW weights may be derived with any procedure that allows for learning a conditional model for the censoring survival function. In our experiments, we use flexible neural network models in discrete time,

such as those described in Kvamme and Borgan [90]. Given these weights, the unconstrained model fitting procedure is weighted ERM. We extend each of the metrics used for evaluation and or as components of the training objectives presented in section 2.5 to account for censoring by incorporating IPCW weights. A full specification of the relevant metrics and training objectives is provided in supplementary section A.3.

### 3.2.2 Assessing net benefit in terms of risk reduction

For the evaluation of models that predict the risk of ASCVD to inform statin initiation, we introduce an alternative formulation of the net benefit that is defined in terms of the population absolute risk reduction after subtracting out harms represented on the same scale. We use the guideline-concordant thresholds of 7.5% and 20%, which correspond to the bounds of the intermediate and high-risk categories, respectively in clinical practice guidelines [37, 40, 41]. We do so to parameterize the net benefit in terms of clinically-plausible benefit-harm trade-offs. Here, we summarize the key aspects of the formulation, but include a full derivation in supplementary section A.2.

For this case, the relevant utilities are defined by the absence ( $u_0^y$ ) and presence ( $u_1^y$ ) of an ASCVD event within ten years. The expected event rates conditioned on the score  $s$  are given by  $p_y^0(s)$  and  $p_y^1(s)$  in the absence and presence of treatment, respectively. The conditional absolute risk reduction is given by  $\text{ARR}(s) = p_y^0(s) - p_y^1(s)$ . We assume that the expected harm of the intervention can be represented as a constant  $k_{\text{harm}}$  that is independent of the risk estimate. With these assumptions,  $U_{\text{cond}}(s) = (u_0^y - u_1^y)\text{ARR}(s) - k_{\text{harm}}$  and the optimal threshold is given by  $\tau_y^* = \text{ARR}^{-1}\left(\frac{k_{\text{harm}}}{u_0^y - u_1^y}\right)$ .

We further assume that the intervention induces constant *relative* risk reduction, such that  $\text{ARR}(s) = rc(s)$  for a constant  $r \in (0, 1)$  and the conditional expected utility and optimal threshold are simple transformations of the calibration curve, as was the case for the fixed-cost setting. In this case,  $U_{\text{cond}}(s) = (u_0^y - u_1^y)rc(s) - k_{\text{harm}}$  and  $\tau_y^* = c^{-1}\left(\frac{k_{\text{harm}}}{r(u_0^y - u_1^y)}\right)$ . We derive a formulation of the net benefit in this setting as

$$\text{NB}(\tau_y; \tau_y^*) = -(1 - \text{NPV}(\tau_y))P(S < \tau_y) - P(S \geq \tau_y)\left((1 - r)\text{PPV}(\tau_y) + r\tau_y^*\right) + P(Y = 1), \quad (13)$$

where  $\text{NPV}(\tau_y)$  and  $\text{PPV}(\tau_y)$  are the negative and positive predictive values evaluated at a threshold  $\tau_y$ . The calibrated net benefit is defined analogously:

$$\begin{aligned} \text{cNB}(\tau_y; \tau_y^*) &= -(1 - \text{NPV}(c^{-1}(\tau_y)))P(S < c^{-1}(\tau_y)) \dots \\ &\quad - P(S \geq c^{-1}(\tau_y))\left((1 - r)\text{PPV}(c^{-1}(\tau_y)) + r\tau_y^*\right) + P(Y = 1). \end{aligned} \quad (14)$$

To operationalize this notion of net benefit, we use a simple model for the treatment effect of statin initiation presented in Soran et al. [91]. This model relates the expected reduction in ASCVD risk to the reduction in low-density lipoprotein cholesterol (LDL-C) that results from statin initiation. It assumes that each 1 mmol/L reduction in LDL-C results in an expected 22% proportional reduction in the ten-year risk of ASCVD, based on evidence from a meta-analysis of randomized control trials [92]. This implies that if the absolute reduction in LDL-C in mmol/L is given by  $\kappa$ , the relative reduction in ten-year ASCVD risk is given by  $r = 1 - (1 - 0.22)^\kappa$  [91]. Therefore, the task of describing the expected reduction in risk as a function of the risk estimate can be reduced to the task of describing the expected reduction in LDL-C as a function of the risk estimate.

To describe the expected reduction in LDL-C from statin therapy for the cohort, we separately consider the evidence for the extent to which statins reduce LDL-C as a function of LDL-C alongside the relationship between observed LDL-C values and the risk estimates for the cohort. As in Soran et al. [91], we assume the use of moderate intensity statin therapy that reduces LDL-C by 43% on average, independent of the pre-treatment level of LDL-C, consistent with the usage of 20 mg of atorvastatin [91, 93, 94]. We extract the most recent historical LDL-C result, if present, for each patient in the test set whose binary outcome was uncensored, filtering out extreme results of  $< 10$  or  $> 500$  mg/dL LDL-C, resulting in 32,366 results. We note that the risk estimates produced by the selected model learned with ERM appear to be uncorrelated with observed untreated LDL-C levels in the cohort ( $R^2 = 0.004$ ; Supplementary Figure C2), suggesting that both the expected absolute reduction in LDL-C and the relative risk reduction  $r$  may be modeled as constants that are independent of the risk estimates. We extract a risk-score-independent estimate of

the mean LDL-C in the cohort as 3.01 mmol/L, using an IPCW-weighted mean over the extracted LDL-C values. The value for the expected relative risk reduction that follows from statin initiation is given by  $r = 1 - (1 - 0.22)^{(3.01*0.43)} = 0.275 = 27.5\%$ .

**Table 1:** Characteristics of the cohort drawn from the Optum CDM database. Data are grouped based on sex, racial and ethnic categories, and the presence of type 2 and type 1 diabetes, rheumatoid arthritis (RA), and chronic kidney disease (CKD). Shown, for each subgroup, is the number of patients extracted, the rate at which the ten-year ASCVD outcome is censored, and an inverse probability of censoring weighted estimate of the incidence of the ten-year ASCVD outcome.

Group	Count	Censoring rate	Incidence
Female	3,253,609	0.816	0.105
Male	2,549,256	0.821	0.120
Asian	165,198	0.814	0.0829
Black	438,144	0.786	0.136
Hispanic	433,238	0.800	0.104
Other	880,116	0.936	0.115
White	3,886,169	0.797	0.110
Asian, female	88,100	0.806	0.0793
Asian, male	77,098	0.823	0.0874
Black, female	262,559	0.784	0.128
Black, male	175,585	0.788	0.150
Hispanic, female	235,736	0.792	0.102
Hispanic, male	197,502	0.810	0.107
Other, female	522,369	0.938	0.108
Other, male	357,747	0.932	0.125
White, female	2,144,845	0.794	0.102
White, male	1,741,324	0.802	0.119
Type 2 diabetes absent	5,388,193	0.817	0.104
Type 2 diabetes present	414,672	0.835	0.20
Type 1 diabetes absent	5,741,282	0.818	0.110
Type 1 diabetes present	61,583	0.825	0.240
RA absent	5,733,505	0.819	0.110
RA present	69,360	0.782	0.185
CKD absent	5,758,773	0.819	0.110
CKD present	44,092	0.767	0.253

### 3.3 Cohort definition

All data are derived from Optum’s de-identified Clinformatics® Data Mart Database (Optum CDM), a statistically de-identified large commercial and medicare advantage claims database containing records from 2007 to 2019. We utilize version 8.1 of the database mapped to the Observational Medical Outcomes Partnership Common Data Model (OMOP CDM) version 5.3.1 [95–97]. Approval for the use of this data for this study was granted by the Stanford Institutional Review Board protocol #46829.

We apply criteria to extract a cohort for learning estimators of ten-year ASCVD risk that mirrors the population eligible for risk-based allocation of statins based on clinical practice guidelines [40]. The characteristics of the extracted cohort are provided in Table 1. We consider as candidate index events all office visits and outpatient encounters for patients between 40 and 75 years of age at the time of the visit for patients without a prior statin prescription or history of cardiovascular disease (Supplementary Table B1). We restrict the set of candidate index events to those recorded as occurring at or before December 31, 2008

for which least one year of historical data is available, and randomly sample one of the resulting candidate index events per patient for inclusion in the final cohort.

The times of ASCVD and censoring events are identified relative to the index event dates. ASCVD events are defined as the occurrence of a diagnosis code for myocardial infarction, stroke, or fatal coronary heart disease (Supplementary Table B1). We consider coronary heart disease to be fatal if death occurs within a year of the recording of the diagnosis code. Censoring events are identified as the earliest date of statin prescription (Supplementary Table B1), death, or the end of the latest enrollment period. From the extracted ASCVD and censoring times, we construct composite binary outcomes and censoring indicators at ten years, following the logic of section 3.2.1.

### 3.3.1 Subgroup definitions

We define discrete subgroups on the basis of (1) a combined race and ethnicity variable based on reported racial and ethnic categories, (2) patient sex, (3) intersectional categories describing intersections of racial and ethnic categories with sex, (4) history of either type 2 diabetes, type 1 diabetes, RA, or CKD at the index date. To construct the race and ethnicity attribute, we assign “Hispanic” if the recorded OMOP CDM concept for ethnicity is recorded as “Hispanic or Latino”, and the value of the recorded OMOP CDM racial category otherwise. This resulted in a final categorization of “Asian”, “Black or African American”, “Hispanic or Latino”, “Other”, and “White”, which we shortened to “Asian”, “Black”, “Hispanic”, “Other”, and “White” for succinctness in the presentation of the results. We identify patients with a history of type 2 diabetes, type 1 diabetes, rheumatoid arthritis, or chronic kidney disease using the presence of a concept identifier indicative of the condition recorded prior to the index date (Supplementary Table B1). The selected concept identifiers used for identifying type 2 and type 1 diabetes are adapted from Reys and Rijnbeek [98]; those used to identify chronic kidney disease are adapted from Suchard et al. [99].

## 3.4 Feature extraction

We apply a procedure similar to the one described in Pfohl et al. [36] to extract a set of clinical features to use as input to fully-connected feedforward neural networks and logistic regression models. This procedure concatenates features representing unique OMOP CDM concepts recorded prior to each patient’s selected index date. We use OMOP CDM concepts corresponding to time-agnostic demographic features (race, ethnicity, sex, and age discretized in five-year intervals) as well as longitudinal recorded diagnoses, medication orders, medical device usage, encounter types, laboratory test orders, flags indicating whether the test results were normal or abnormal based on reference ranges, and other coded clinical observations binned in three time intervals corresponding to 29 to 1 days prior to the index date, 365 days to 30 days prior to the index, and any time prior to the index date.

## 3.5 Data partitioning

We partition the cohort such that 62.5% is used as a training set, 12.5% is used as a validation set, and 25% of the data is used as a test set. We subsequently partition the training data into five equally-sized partitions. We train five models for each hyperparameter configuration, holding out one of the partitions of the training set for use as a development set to assess early stopping criteria, and perform model selection based on algorithm-specific model selection criteria defined over the average performance of the five models on the validation set.

## 3.6 Derivation of inverse probability of censoring weights

We consider the estimation of the risk of ASCVD at a fixed time horizon as an example of a supervised learning problem with a censored binary outcome, using the procedures described in section 3.2.1. To derive IPCW weights, we utilize neural networks trained with the discrete-time likelihood [90, 100, 101] to estimate the censoring survival function conditioned on the full set of features used to fit the model for ten-year ASCVD. For each cohort, we derive five such models using the training set partitioning strategy described in section 3.5. We use a fixed model architecture with one hidden layer of 128 hidden units that predicts the discrete-time hazard in twenty intervals whose boundaries are determined by the quantiles of the observed

censoring times in the union of the four training set partitions that are not held-out. We train these models in a minibatch setting and perform early stopping if the discrete-time likelihood does not improve for twenty-five epochs of 100 minibatches. Subsequently, we define IPCW weights for each patient in the training set by taking the inverse of the predicted censoring survival function at the minimum of the time of censoring, the ASCVD outcome event, or ten years, for each patient, using the model trained on the set of training set partitions that exclude the patient. The weights for patients in the validation and test sets are derived as the reciprocal of the average estimate of the censoring survival function derived from the five models.

### 3.7 Experiments

Here, we outline the structure of the experiments. To serve as baseline comparators for all experiments, we train models using unconstrained IPCW-weighted ERM without stratification. We refer to this setting as *pooled ERM*. The first experiment aims to evaluate strategies to learn models that predict the outcome well for subgroups defined following stratification by race, ethnicity, and sex, including intersectional categories, and for patients with ASCVD-promoting comorbidities. The second experiment aims to assess the implications of penalizing violation of the equalized odds criterion across subgroups defined on the basis of race, ethnicity, and sex. In each case, we evaluate the net benefit of statin initiation on the basis of the risk estimates under the assumption that the observed relationship with the benefits of using the ASCVD risk estimator to initiate moderate-intensity statin therapy can be modeled as inducing constant relative risk reduction (section 3.2.2), the expected harm of treatment is assumed not to vary on the basis of the risk estimate, and that the trade-off between benefits and harms reflects the choice of a decision threshold of either 7.5% or 20%.

#### 3.7.1 Unconstrained empirical risk minimization without stratification

We evaluate feedforward neural networks and logistic regression models trained with pooled ERM in a minibatch setting using stochastic gradient descent. We conduct a grid search over model-specific and algorithm-specific hyperparameters. For feedforward neural networks, we evaluate a grid of hyperparameters that include learning rates of  $1 \times 10^{-4}$  and  $1 \times 10^{-5}$ , one and three hidden layers of size 128 or 256 hidden units, and a dropout probability of 0.25 or 0.75. For logistic regression models, we use weight decay regularization [102] drawn from a grid of values containing 0, 0.01, and 0.001. The training procedure is conducted in a minibatch setting of up to 150 iterations of 100 minibatches of size 512 using the Adam [103] optimizer in the Pytorch framework [104]. We use an early-stopping rule that returns the model with the lowest log-loss evaluated on the development set when the log-loss has not improved for twenty-five epochs of 100 minibatches. The procedure is repeated separately for each of the five training/development set partitions. Following training, we apply each model derived from the training procedure to the validation set and select hyperparameters on the basis of the best average log-loss evaluated in the validation set across all training partitions.

#### 3.7.2 Approaches to improve model performance over patient subgroups

To compare with pooled ERM, we evaluate models trained with ERM separately on each subgroup (*stratified ERM*), models trained with IPCW-weighted regularized training objectives that penalize differences in the log-loss or AUC between each subgroup and the overall population, and IPCW-weighted DRO objectives that target the worst-case log-loss or AUC across subgroups. The hyperparameter grid, early stopping, and model selection procedures conducted for the stratified ERM experiments exactly match those used for the pooled ERM experiments. For models trained with regularized objectives or DRO, we use a feedforward neural network with hyperparameters fixed to three hidden layers with 256 hidden units, a dropout probability of 0.25, and a learning rate of  $1 \times 10^{-4}$ . For the regularized models, we evaluate a grid of five  $\lambda$  values distributed log-uniformly from  $1 \times 10^{-2}$  to 10 and conduct early-stopping on the basis of the value of the penalized loss. For the DRO experiments, we evaluate unmodified and balanced sampling by subgroup, as well as a grid of value of  $\eta$  given by 0.01, 0.1, and 1.

As in the case of the unpenalized ERM experiments, we fix the batch size to be 512 and evaluate early-stopping criteria in intervals of 100 minibatches and terminate when the criterion has not improved for 25 iterations. For the fairness-regularized models, we perform early stopping on the basis of the penalized loss that incorporates the regularization term. To conduct early-stopping for DRO experiments, we use

worst-case early-stopping criteria [36] that returns the model with the best worst-case subgroup AUC or log-loss observed over the training procedure when the worst-case value has not improved for twenty-five epochs of 100 minibatches. We use the worst-case subgroup AUC for early-stopping when the AUC-based training objective is used and the worst-case subgroup log-loss when the DRO objective defined over the log-loss is used.

For model selection on the validation set, we use criteria defined in terms of the worst-case performance (either AUC or log-loss) for both regularized and DRO experiments, over the full set of hyperparameter configurations. We use the worst average performance produced by averaging validation set performance over the training replicates. As was the case for early-stopping, we use the worst-case AUC for model selection for the regularized and DRO experiments that incorporate the AUC into their objective, and use the worst-case log-loss for model selection for objectives that incorporate the log-loss into their objective.

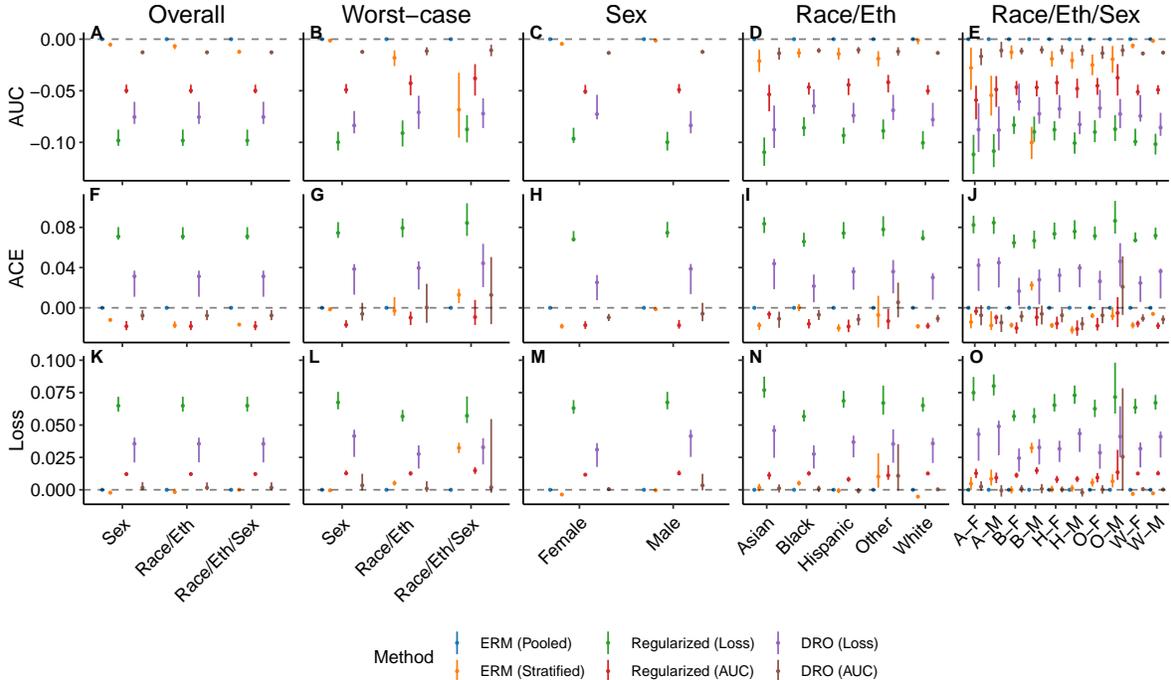
### 3.7.3 Regularized fairness objectives for equalized odds

To evaluate the effect of penalizing violation of equalized odds, we use regularized training objectives that incorporate an IPCW-weighted MMD penalty to penalize differences in the outcome-conditioned distribution of the risk score between each subgroup and the marginal population (equations (8) and (31)), as well as a penalty that penalizes differences in the true positive and false negative rates between each subgroup and the marginal population at the guideline-relevant thresholds of 7.5% and 20% [40] using an IPCW-weighted objective that uses a softplus relaxation to the indicator function (equation (32)). To simplify the experiment, we conduct this analysis only with the intersectional categories defined by race, ethnicity, and sex, but separately evaluate the models on the intersectional categories and for race/ethnicity and sex separately. Furthermore, we fix hyperparameters to those used for the regularized and DRO models in section 3.7.2 and evaluate five values of the regularization penalty  $\lambda$  distributed log-uniformly from  $1 \times 10^{-2}$  to 10. As before, we fix the batch size to be 512 and evaluate early-stopping criteria in intervals of 100 minibatches and terminate when the value of the penalized loss has not improved for 25 iterations. For these models, we do not conduct explicit model selection over the regularization path on the basis of validation set performance given that it was of interest to evaluate each value of  $\lambda$  separately.

### 3.7.4 Evaluation of model performance

To compute 95% confidence intervals for model performance metrics, we draw 1,000 bootstrap samples from the test set, stratified by levels of the outcome and subgroup attribute relevant to the evaluation, compute the IPCW-weighted performance metrics for the set of five derived models on each bootstrap sample, and take the 2.5% and 97.5% empirical quantiles of the resulting distribution that results from pooling over both the models and bootstrap replicates. We construct analogous confidence intervals for the difference in the model performance relative to pooled ERM by computing the difference in the performance computed on the same bootstrap sample and taking the 2.5% and 97.5% empirical quantiles of the distribution of the differences. To construct confidence intervals for the worst-case performance over subgroups, we extract the worst-case performance for each bootstrap sample.

We assess model performance in the test set in terms of IPCW-weighted variants of the AUC, the average log-loss, the absolute calibration error (ACE) [14, 105, 106], true positive rate, false positive rate, calibration curve, and the net benefit. Estimates of the calibration curve used to compute the ACE and calibrated net benefit rely on an estimate of the calibration curve learned via a logistic regression estimator trained on the test data to predict the outcome from a logit-transformed outputs of the predictive model as inputs. The inverse of the calibration curve used to compute the calibrated net benefit is derived analytically based on the coefficients of the learned logistic regression model. To compute the ACE, we use an IPCW-weighted average of the absolute value of the differences between the model output and the calibration curve.



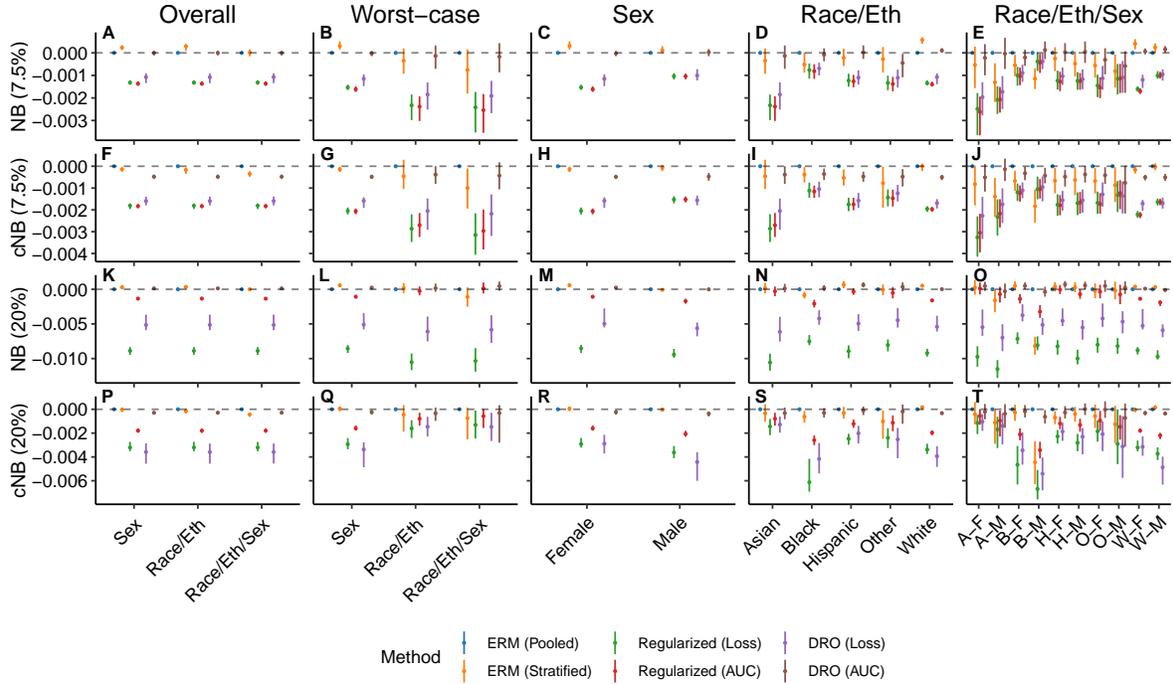
**Figure 1:** The performance of models that estimate ten-year ASCVD risk, stratified by race, ethnicity, and sex, relative to the results attained by the application of unpenalized ERM to the overall population. Results shown are the relative AUC, absolute calibration error (ACE), and log-loss assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unconstrained pooled or stratified ERM, regularized objectives that penalize differences in the log-loss or AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss or AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.

## 4 Results

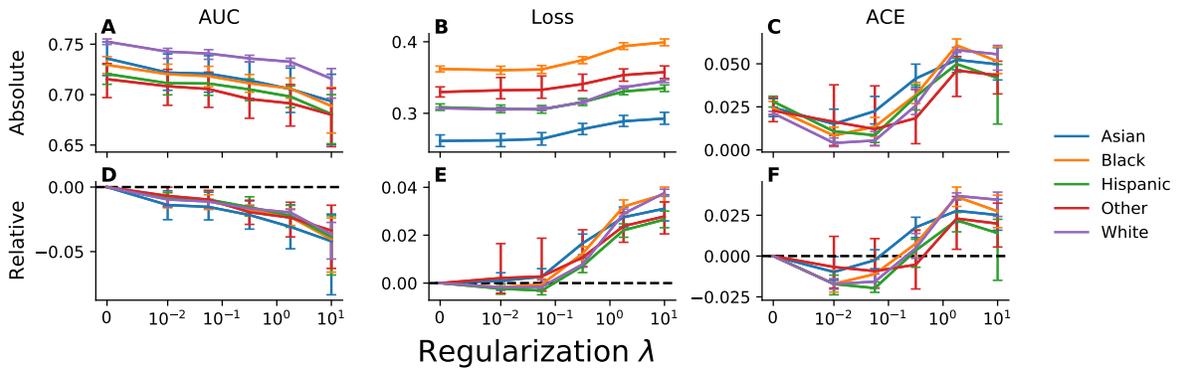
### 4.1 Approaches to improve model performance over subgroups

We conducted an experiment to assess whether approaches that penalize differences in AUC or log-loss across subgroups or optimize for the worst-case value of these metrics improve upon empirical risk minimization approaches in terms of the model performance and net benefit measures. In the main text, we report the results assessed relative to those derived from unpenalized ERM applied to the entire population for subgroups defined in terms of race, ethnicity, and sex (Figure 1), as well as for subgroups with ASCVD-promoting comorbidities (Supplementary Figure C3). Absolute performance estimates are reported in the supplementary material (Supplementary Figure C4 and Supplementary Figure C5).

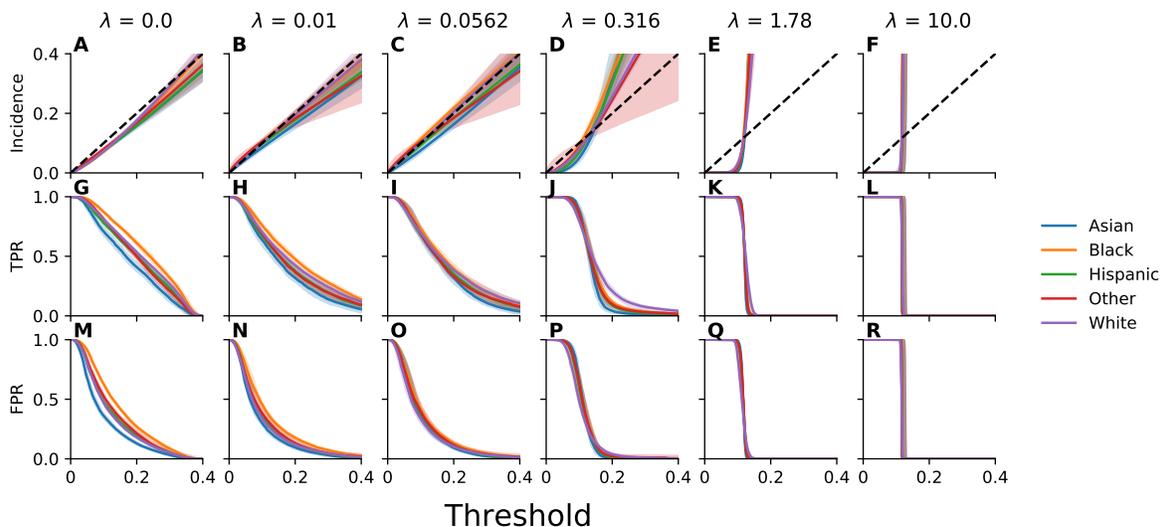
We find that the use of unconstrained empirical risk minimization using data from the entire population typically results in models with the greatest AUC for each subgroup, but stratified ERM procedures that train a separate model for each subgroup achieve an AUC that does not differ substantially in some cases, particularly for majority subgroups (Figure 1D,E and Supplementary Figure C3C,D,E,F). The models trained with regularized fairness objectives or DRO and selected on the basis of the worst-case AUC or log-loss do not improve on the AUC assessed for each subgroup, and typically perform substantially worse, with the least extreme degradation observed for those models trained with the AUC-based DRO training objective (Figure 1C,D,E and Supplementary Figure C3C,D,E,F). Despite the lack of improvement in AUC, we observe that subgroup-specific ERM and both regularized and DRO-based objectives that incorporate the AUC into their training objective often result in improved model calibration for some subgroups (1F,G,H,I,J and Supplementary Figure C3G,I,J,K,L). Similarly, subgroup-specific training does result in minor improvements in the log-loss for some subgroups relative to ERM applied to the entire population, but these results are



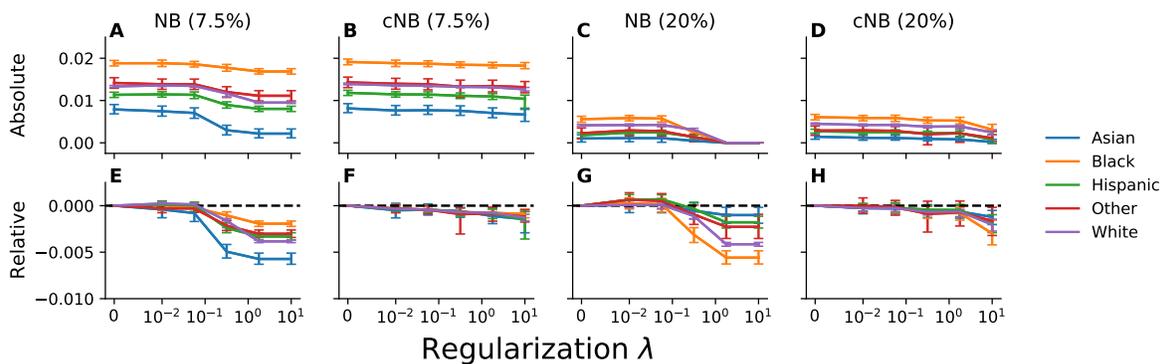
**Figure 2:** The net benefit of models that estimate ten-year ASCVD risk, stratified by race, ethnicity, and sex, relative to the results attained by the application of unpenalized ERM to the overall population. Results shown are the net benefit (NB) and calibrated net benefit (cNB), parameterized by the choice of a decision threshold of 7.5% or 20%, assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unconstrained pooled or stratified ERM, regularized objectives that penalize differences in the log-loss or AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss or AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



**Figure 3:** Model performance evaluated across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the area under the receiver operating characteristic curve (AUC), log-loss, and absolute calibration error (ACE). Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



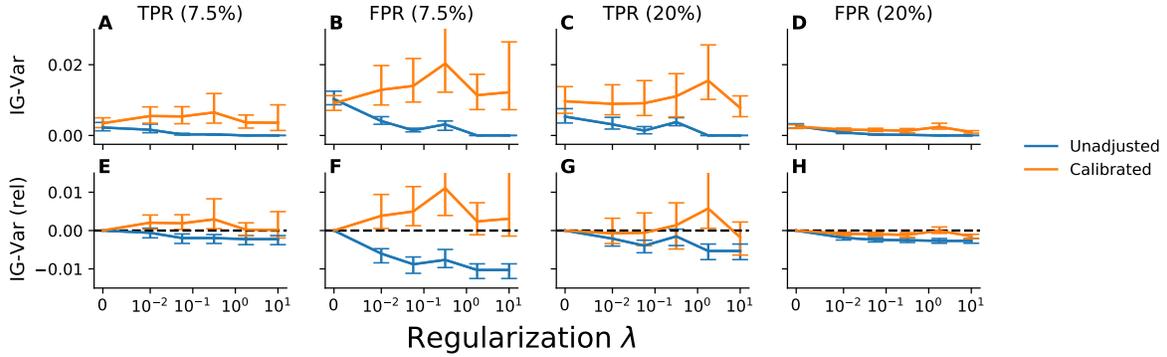
**Figure 4:** Calibration curves, true positive rates, and false positive rates evaluated for a range of thresholds across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , are the calibration curve (incidence), true positive rate (TPR), and false positive rate (FPR) as a function of the decision threshold. Error bands indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



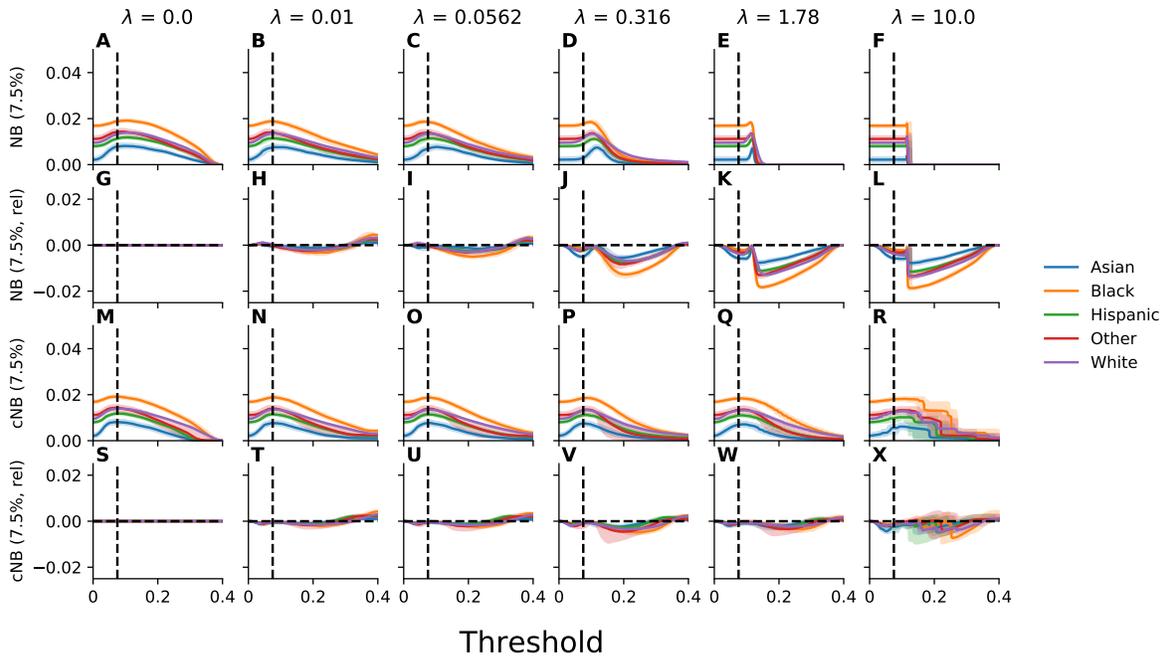
**Figure 5:** The net benefit evaluated across racial and ethnic subgroups, parameterized by the choice of a decision threshold of 7.5% or 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup, is the net benefit (NB) and calibrated net benefit (cNB) as a function of the value of the regularization parameter  $\lambda$ . Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.

typically observed only for larger subgroups when they are present (Figure 1M,D,O and Supplementary Figure C3O,R).

The implication of these effects can be understood holistically through an assessment of the net benefit of statin therapy initiated on the basis of the risk estimates. Overall, no approach consistently confers more net benefit than unpenalized ERM applied to the entire population for each subgroup, when the net benefit is assessed for the benefit-harm tradeoffs corresponding to either of the thresholds of 7.5% or 20%, but subgroup-specific training and AUC-based DRO approaches do lead to minor improvements in some cases (Figure 2C,D,E,M,N,O and Supplementary Figure C7C,F,O,R). However, we note that, for each subgroup, no



**Figure 6:** Satisfaction of equalized odds evaluated across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted is the intergroup variance (IG-Var) in the true positive and false positive rates at decision thresholds of 7.5% and 20%. Recalibrated results correspond to those attained for models for which the threshold has been adjusted to account for the observed miscalibration. Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



**Figure 7:** The net benefit evaluated for a range of thresholds across racial and ethnic subgroups, parameterized by the choice of a decision threshold of 7.5%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.

approach improves on the calibrated net benefit, *i.e.* the net benefit achieved following adjustment of the decision threshold to account for the observed miscalibration, relative to unpenalized ERM applied to the entire population (Figure 2H,I,J,R,S,T and Supplementary Figure C7I,J,K,L,U,V,W,X). This indicates that

for those cases where an alternative strategy results in an increase in the net benefit conferred relative to that which is achieved for the pooled ERM strategy, it is a consequence of the improvement in calibration at the threshold of interest.

## 4.2 Regularized fairness objectives for equalized odds

We further conducted an experiment to assess the implications of the use of a training objective that penalizes violation of equalized odds across intersectional subgroups defined by race, ethnicity, and sex. In the main text, we present the results corresponding to an MMD-based penalty evaluated over subgroups defined by race and ethnicity, but include in the supplementary material analogous results corresponding to evaluation over intersectional categories and for sex (Supplementary Figures C11 to C24). Furthermore, the supplementary material includes analogous results for experiments that penalize equalized odds at both of the thresholds of 7.5% and 20% using softplus relaxations of the true positive and false positive rates (Supplementary Figures C25 to C45).

We observe that as the strength of the penalty  $\lambda$  increases, the AUC assessed for each subgroup monotonically decreases (Figure 3A,D). With a minor degree of equalized-odds promoting regularization (*i.e.*  $\lambda = 0.01, 0.0562$ ), calibration actually improves relative to the result for unpenalized ERM (Figure 3C,F) and there is little to no change in the log-loss for each subgroup despite the reduction in AUC (Figure 3B,E). This is reflected in the calibration curves presented in Figure 4, where we observe modest miscalibration consistent with overestimation of risk for each subgroup for the unconstrained model (Figure 4A) with improvements in the calibration of the model with a minor degree of regularization (Figure 4B,C). However, for large degrees of regularization (*i.e.*  $\lambda = 1.78$  and  $\lambda = 10$ ), both the calibration and log-loss assessed for each subgroup deteriorates, although the reduction in AUC remains modest (Figure 3). In this case, the variability in the risk estimates sharply decreases to concentrate around the incidence of the outcome for larger degrees of regularization, which is reflected in the shape of the calibration curve and error rates as a function of the threshold (Figure 4F,L,R), consistent with overestimation for patients with risk lower than the incidence and underestimation for patients with risk greater than the incidence.

For the unconstrained model, the true positive rates and false positive rates at each threshold are ranked across subgroups in accordance with the observed incidence for each subgroup, such that the Black population has the largest true positive rate and false positive rate while the Asian population has the lowest true positive rate and false positive rate (Figure 4G,H). The regularized training objective is successful at enforcing the equalized odds constraint, in that the variability in false positive and true positives rates trends towards zero as the strength of the penalty increases (Figures 4 and Figure 6).

For the benefit-harm tradeoff implied by the use of either a threshold of 7.5% or 20%, we observe clear reductions in net benefit for each subgroup for large values of  $\lambda$  (Figure 5A,C,E,G). With minor amounts of regularization, we observe little to no reduction in net benefit parameterized by either a threshold of 7.5% or 20%, and the point estimates for 20% even suggest a relative increase in net benefit compared to unpenalized ERM (Figure 5E,G). However, for large degrees of regularization, we observe large reductions in net benefit relative to that which is attained from unpenalized ERM, but the magnitude of these differences are attenuated when the thresholds applied for each subgroup are adjusted to account for miscalibration (Figure 5B,D,F,H). We further observe that the calibrated net benefit for equalized odds penalized models does not improve on unpenalized ERM at any value of  $\lambda$  (Figure 5C,F,D,H). Overall, the reduction in net benefit observed directly due to operating at a suboptimal decision threshold, as a result of miscalibration, is generally larger than the reduction in net benefit that results due to the reduction in the AUC of the model at larger values of  $\lambda$ . Furthermore, we note that threshold adjustment to recover net benefit lost due to the miscalibration resulting from the use of the training objective that penalizes equalized odds violation does not preserve the satisfaction of the equalized odds fairness constraint, as the variability in error rates at the adjusted thresholds is observed to be similar to or more variable than that which results from unpenalized ERM (Figure 6).

To gain further insight into these phenomena, we plot the net benefit for a range of decision thresholds, assuming that the benefit-harm tradeoff is fixed to one implied by the use of a threshold of 7.5% (Figure 7). In the supplementary material, we include analogous results for the threshold of 20% (Supplementary Figure C9)), as well as standard decisions curves defined such that the net benefit plotted for each point on the curve corresponds to the benefit-harm tradeoff implied by the corresponding threshold on the x-axis (Supplementary

Figure C10)). As expected for the analysis corresponding to a threshold of 7.5%, the calibrated net benefit is maximized for each subgroup at a threshold on the risk estimates corresponding to the point where the observed incidence of the outcome conditioned on the risk estimate is 7.5% (Figure 7M,N,O,P,Q,R). Furthermore, when the model overestimates risk at a threshold of 7.5% due to miscalibration, such as was the case for the unpenalized ERM model and for the models trained with a large penalty on equalized odds violation, the threshold that maximizes the net benefit is one greater than 7.5% (Figure 7A,D,E,F). In these cases, adjusting the threshold on the penalized models to compensate for miscalibration recovers the majority of difference in net benefit relative to the model derived with unpenalized ERM.

## 5 Discussion

The results suggest that in settings where the observed model miscalibration may be adjusted for with subgroup-specific recalibration or via threshold-adjustment, no approach to learning an ASCVD risk estimator confers more net benefit for each subgroup than unpenalized ERM applied to the entire population. This claim follows from the observation that no alternative approach resulted in greater *calibrated* net benefit for any subgroup. We find that the net benefit for each population is maximized for each subgroup at a threshold on the risk score that is consistent with the analysis presented in section 2.3.

In cases where we observe improvements in the unadjusted net benefit over ERM, or little to no change despite a reduction in AUC, the differences directly follow from improvements in the calibration of the model derived from the alternative approach. We observe such effects for models trained with objectives that penalize equalized odds to a minor degree, those trained with stratified ERM procedures that train a separate model for each subgroup, as well as for regularized fairness objectives and DRO procedures that operate over the AUC assessed for each subgroup. Taken together, these results indicate that models derived from unpenalized ERM should not necessarily be assumed to be well-calibrated in practice, further highlighting the importance of model development, selection, and post-processing strategies that aims to identify the best-fitting, well-calibrated model for each subgroup.

Algorithmic fairness assessments in healthcare based on the equalized odds criterion are likely to be misleading. If the model is calibrated and fits well for each subgroup, differences in those error rates are expected when the observed outcome incidence differs [19]. Similarly, effort undertaken to minimize equalized odds violation is likely to introduce harm when it results in unrecognized miscalibration or reduction of model fit. If the differences in outcome incidence reflect measurement error that differs systematically across subgroups [1, 28], then violation of equalized odds may be present. However, for such a case, we argue for conducting a calibration-based fairness assessment with respect to a proxy of the targeted unobserved outcome that is not subject to differential measurement error across subgroups, as in Obermeyer et al. [1]. When those differences in outcome incidence, which may or may not be mismeasured or observable, are a result of population-level differences in disease burden across patient subgroups as a result of structural disparities [107, 108], we argue that the appropriate response is to endeavor to understand both the cause of those disparities and the impact of potential interventions on the structural factors that perpetuate health disparities [109–111].

While this work motivates the use of approaches that reason about algorithmic fairness in terms of calibration characteristics [9, 112], such assessments are not comprehensive. For instance, calibration-based assessments do not account for differences in benefit that arise due to differences in the discrimination performance nor in differences in unmodeled heterogeneity in the outcome across subgroups [18, 50]. The presence of measurement error can also mask consequential violation of sufficiency with respect to the targeted unobserved outcome that is not subject to measurement error [1, 28]. Furthermore, when a predictive model is used for referral to a clinical service that cannot process more than a fixed number of cases due to resource constraints, *e.g.* as in Jung et al. [22], then it may not be practical to operate at the utility-maximizing threshold. In that case, differences in the magnitude of the unrealized utility across subgroups are likely if the distribution of risk differs across subgroups, even if sufficiency holds and a consistent global threshold is applied across subgroups. Such a capacity constrained situation poses a set of ethical conflicts and trade-offs that should be navigated with participatory processes incorporating the preferences and attitudes of a diverse set of stakeholders [113].

## 6 Acknowledgements

We thank the Stanford Center for Population Health Sciences Data Core and the Stanford Research Computing Center for supporting the data and computing infrastructure used in this work. This work is supported by the National Heart, Lung, and Blood Institute R01 HL144555 and the Stanford Medicine Program for AI in Healthcare. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding bodies.

## 7 Data availability

Individuals wishing to access the data used in this work may sign a data use agreement with Stanford and Optum to access the data for replication or confirmatory studies on the Stanford Secure Data Ecosystem.

## 8 Code availability

We make all code available at [https://github.com/som-shahlab/net\\_benefit\\_ascvd](https://github.com/som-shahlab/net_benefit_ascvd).

## 9 Competing interests statement

The authors declare no competing interests.

## References

- [1] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.
- [2] Darshali A. Vyas, Leo G. Eisenstein, and David S. Jones. Hidden in Plain Sight — Reconsidering the Use of Race Correction in Clinical Algorithms. *New England Journal of Medicine*, page NEJMms2004740, jun 2020. ISSN 0028-4793. doi: 10.1056/NEJMms2004740.
- [3] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. fairmlbook.org, 2019. URL [fairmlbook.org](http://fairmlbook.org).
- [4] Alvin Rajkomar, Michaela Hardt, Michael D. Howell, Greg Corrado, and Marshall H. Chin. Ensuring Fairness in Machine Learning to Advance Health Equity. *Annals of Internal Medicine*, dec 2018. ISSN 0003-4819. doi: 10.7326/M18-1990.
- [5] Irene Y Chen, Peter Szolovits, and Marzyeh Ghassemi. Can ai help reduce disparities in general medical and mental health care? *AMA journal of ethics*, 21(2):167–179, 2019.
- [6] R Yates Coley, Eric Johnson, Gregory E Simon, Maricela Cruz, and Susan M Shortreed. Racial/ethnic disparities in the performance of prediction models for death by suicide after mental health visits. *JAMA psychiatry*, 2021.
- [7] Laleh Seyyed-Kalantari, Guanxiong Liu, Matthew McDermott, Irene Y Chen, and Marzyeh Ghassemi. Chexclusion: Fairness gaps in deep chest x-ray classifiers. In *BIOCOMPUTING 2021: Proceedings of the Pacific Symposium*, pages 232–243. World Scientific, 2020.
- [8] Yoonyoung Park, Jianying Hu, Moninder Singh, Issa Sylla, Irene Dankwa-Mullan, Eileen Koski, and Amar K Das. Comparison of methods to reduce bias from clinical prediction models of postpartum depression. *JAMA network open*, 4(4):e213909–e213909, 2021.
- [9] Noam Barda, Gal Yona, Guy N Rothblum, Philip Greenland, Morton Leibowitz, Ran Balicer, Eitan Bachmat, and Noa Dagan. Addressing bias in prediction models by improving subpopulation calibration. *Journal of the American Medical Informatics Association*, 28(3):549–558, 2021.

- [10] Stephen Pfohl, Ben Marafino, Adrien Coulet, Fatima Rodriguez, Latha Palaniappan, and Nigam H. Shah. Creating Fair Models of Atherosclerotic Cardiovascular Disease Risk. In *AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society*, 2019.
- [11] Anna Zink and Sherri Rose. Fair regression for health care spending. *Biometrics*, 76(3):973–982, 2020.
- [12] Imon Banerjee, Ananth Reddy Bhimireddy, John L Burns, Leo Anthony Celi, Li-Ching Chen, Ramon Correa, Natalie Dullerud, Marzyeh Ghassemi, Shih-Cheng Huang, Po-Chih Kuo, et al. Reading race: Ai recognises patient’s racial identity in medical images. *arXiv preprint arXiv:2107.10356*, 2021.
- [13] Laleh Seyyed-Kalantari, Haoran Zhang, Matthew McDermott, Irene Y Chen, and Marzyeh Ghassemi. Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. *Nature medicine*, pages 1–7, 2021.
- [14] Stephen R. Pfohl, Agata Foryciarz, and Nigam H. Shah. An empirical characterization of fair machine learning for clinical risk prediction. *Journal of Biomedical Informatics*, 113:103621, 2021. ISSN 1532-0464. doi: <https://doi.org/10.1016/j.jbi.2020.103621>.
- [15] Agata Foryciarz, Stephen R. Pfohl, Birju Patel, and Nigam H. Shah. Evaluating algorithmic fairness in the presence of clinical guidelines: the case of atherosclerotic cardiovascular disease risk estimation. *medRxiv*, 2021. doi: 10.1101/2021.11.08.21266076. URL <https://www.medrxiv.org/content/early/2021/11/10/2021.11.08.21266076>.
- [16] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent Trade-Offs in the Fair Determination of Risk Scores. *arXiv preprint arXiv:1609.05807*, sep 2016. ISSN 17409713. doi: 10.1111/j.1740-9713.2017.01012.x.
- [17] Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *ArXiv e-prints*, feb 2017. ISSN 2167-6461. doi: 10.1089/big.2016.0047.
- [18] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic Decision Making and the Cost of Fairness. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’17*, pages 797–806, New York, NY, USA, jan 2017. ACM. ISBN 978-1-4503-4887-4. doi: 10.1145/3097983.3098095.
- [19] Lydia T Liu, Max Simchowitz, and Moritz Hardt. The Implicit Fairness Criterion of Unconstrained Learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4051–4060, Long Beach, California, USA, 2019. PMLR. URL <http://proceedings.mlr.press/v97/liu19f.html>.
- [20] Judy Wawira Gichoya, Liam G McCoy, Leo Anthony Celi, and Marzyeh Ghassemi. Equity in essence: a call for operationalising fairness in machine learning for healthcare. *BMJ Health & Care Informatics*, 28(1), 2021. doi: 10.1136/bmjhci-2020-100289. URL <https://informatics.bmj.com/content/28/1/e100289>.
- [21] Abigail Z Jacobs and Hanna Wallach. Measurement and fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 375–385, 2021.
- [22] Kenneth Jung, Sehj Kashyap, Anand Avati, Stephanie Harman, Heather Shaw, Ron Li, Margaret Smith, Kenny Shum, Jacob Javitz, Yohan Vetteth, et al. A framework for making predictive models useful in practice. *Journal of the American Medical Informatics Association*, 28(6):1149–1158, 2021.
- [23] Irene Y Chen, Emma Pierson, Sherri Rose, Shalmali Joshi, Kadija Ferryman, and Marzyeh Ghassemi. Ethical machine learning in healthcare. *Annual Review of Biomedical Data Science*, 4, 2020.
- [24] Ruha Benjamin. Assessing risk, automating racism. *Science*, 366(6464):421–422, 2019.

- [25] Samir Passi and Solon Barocas. Problem formulation and fairness. In *FAT\* 2019 - Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency*, pages 39–48. Association for Computing Machinery, Inc, jan 2019. ISBN 9781450361255. doi: 10.1145/3287560.3287567.
- [26] Mark P Sendak, Michael Gao, Nathan Brajer, and Suresh Balu. Presenting machine learning model information to clinical end users with model facts labels. *NPJ digital medicine*, 3(1):1–4, 2020.
- [27] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*, 2018.
- [28] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229, 2019.
- [29] Harold C. Sox, Michael C. Higgins, and Douglas K. Owens. *Medical Decision Making*. John Wiley & Sons, Ltd, Chichester, UK, jun 2013. ISBN 9781118341544. doi: 10.1002/9781118341544. URL <http://doi.wiley.com/10.1002/9781118341544>.
- [30] Andrew J. Vickers, Ben Van Calster, and Ewout W. Steyerberg. Net benefit approaches to the evaluation of prediction models, molecular markers, and diagnostic tests. *BMJ (Online)*, 352, jan 2016. ISSN 17561833. doi: 10.1136/bmj.i6.
- [31] Chloé Bakalar, Renata Barreto, Stevie Bergman, Miranda Bogen, Bobbie Chern, Sam Corbett-Davies, Melissa Hall, Isabel Kloumann, Michelle Lam, Joaquin Quiñonero Candela, et al. Fairness on the ground: Applying algorithmic fairness approaches to production systems. *arXiv preprint arXiv:2103.06172*, 2021.
- [32] Camelia Simoiu, Sam Corbett-Davies, and Sharad Goel. The problem of infra-marginality in outcome tests for discrimination. *The Annals of Applied Statistics*, 11(3):1193–1216, 2017.
- [33] Laure Wynants, Maarten Van Smeden, David J. McLernon, Dirk Timmerman, Ewout W. Steyerberg, and Ben Van Calster. Three myths about risk thresholds for prediction models. *BMC Medicine*, 17(1):192, oct 2019. ISSN 17417015. doi: 10.1186/s12916-019-1425-3. URL <https://bmcmmedicine.biomedcentral.com/articles/10.1186/s12916-019-1425-3>.
- [34] Andrew J Vickers and Elena B Elkin. Decision curve analysis: a novel method for evaluating prediction models. *Medical Decision Making*, 26(6):565–574, 2006.
- [35] Andrew J Vickers, Michael W Kattan, and Daniel J Sargent. Method for evaluating prediction models that apply the results of randomized trials to individual patients. *Trials*, 8(1):1–11, 2007.
- [36] Stephen R Pfohl, Haoran Zhang, Yizhe Xu, Agata Foryciarz, Marzyeh Ghassemi, and Nigam H Shah. A comparison of approaches to improve worst-case predictive model performance over patient subpopulations. *arXiv preprint arXiv:2108.12250*, 2021.
- [37] David C Goff, Donald M Lloyd-Jones, Glen Bennett, Sean Coady, Ralph B D’agostino, Raymond Gibbons, Philip Greenland, Daniel T Lackland, Daniel Levy, Christopher J O’donnell, 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. *Journal of the American College of Cardiology*, 63 (25 Part B):2935–2959, 2014.
- [38] Neil J. Stone, Jennifer G. Robinson, Alice H. Lichtenstein, C. Noel Bairey Merz, Conrad B. Blum, Robert H. Eckel, Anne C. Goldberg, David Gordon, Daniel Levy, Donald M. Lloyd-Jones, Patrick McBride, J. Sanford Schwartz, Susan T. Shero, Sidney C. Smith, Karol Watson, and Peter W.F. F. Wilson. 2013 ACC/AHA guideline on the treatment of blood cholesterol to reduce atherosclerotic cardiovascular risk in adults: A report of the american college of cardiology/american heart association task force on practice guidelines. *Circulation*, 129(25 SUPPL. 1):S1–S45, jun 2014. ISSN 15244539. doi: 10.1161/01.cir.0000437738.63853.7a. URL <http://circ.ahajournals.org/lookup/doi/10.1161/01.cir.0000437738.63853.7a>.

- [39] Scott M. Grundy, Neil J. Stone, Alison L. Bailey, Craig Beam, Kim K. Birtcher, Roger S. Blumenthal, Lynne T. Braun, Sarah de Ferranti, Joseph Faiella-Tommasino, Daniel E. Forman, Ronald Goldberg, Paul A. Heidenreich, Mark A. Hlatky, Daniel W. Jones, Donald Lloyd-Jones, Nuria Lopez-Pajares, Chiadi E. Ndumele, Carl E. Orringer, Carmen A. Peralta, Joseph J. Saseen, Sidney C. Smith, Laurence Sperling, Salim S. Virani, and Joseph Yeboah. 2018 AHA/ACC/AACVPR/AAPA/ABC/ACPM/ADA/AGS/APhA/ASPC/NLA/PCNA Guideline on the Management of Blood Cholesterol: Executive Summary: A Report of the American College of Cardiology/American Heart Association Task Force on Clinical Practice Guidelines. *Journal of the American College of Cardiology*, 73(24):3168–3209, jun 2019. ISSN 15583597. doi: 10.1016/j.jacc.2018.11.002.
- [40] Donna K Arnett, Roger S Blumenthal, Michelle A Albert, Andrew B Buroker, Zachary D Goldberger, Ellen J Hahn, Cheryl Dennison Himmelfarb, Amit Khera, Donald Lloyd-Jones, J William McEvoy, et al. 2019 acc/aha guideline on the primary prevention of cardiovascular disease: a report of the american college of cardiology/american heart association task force on clinical practice guidelines. *Journal of the American College of Cardiology*, 74(10):e177–e232, 2019.
- [41] Donald M. Lloyd-Jones, Lynne T. Braun, Chiadi E. Ndumele, Sidney C. Smith Jr, Laurence S. Sperling, Salim S. Virani, and Roger S. Blumenthal. Use of Risk Assessment Tools to Guide Decision-Making in the Primary Prevention of Atherosclerotic Cardiovascular Disease: A Special Report From the American Heart Association and American College of Cardiology. *Circulation*, 139(25):E1162–E1177, jun 2019. doi: 10.1161/CIR.0000000000000638. URL <https://www.ahajournals.org/doi/abs/10.1161/CIR.0000000000000638>.
- [42] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of Opportunity in Supervised Learning. *Advances in Neural Information Processing Systems*, pages 3315–3323, 2016. ISSN 10495258. doi: 10.1109/ICCV.2015.169.
- [43] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. Building classifiers with independency constraints. In *2009 IEEE International Conference on Data Mining Workshops*, pages 13–18. IEEE, 2009. ISBN 9780769539027. doi: 10.1109/ICDMW.2009.83. URL [https://www.win.tue.nl/~sim\\$mp\\$pechen/publications/pubs/CaldersICDM09.pdf](https://www.win.tue.nl/~sim$mp$pechen/publications/pubs/CaldersICDM09.pdf).
- [44] L. Elisa Celis, Lingxiao Huang, Vijay Keswani, and Nisheeth K. Vishnoi. Classification with Fairness Constraints: A Meta-Algorithm with Provable Guarantees. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 319–328, jun 2018.
- [45] Andrew Cotter, Heinrich Jiang, Maya R Gupta, Serena Wang, Taman Narayan, Seungil You, Karthik Sridharan, Maya R Gupta, Seungil You, and Karthik Sridharan. Optimization with Non-Differentiable Constraints with Applications to Fairness, Recall, Churn, and Other Goals. *Journal of Machine Learning Research*, 20(172):1–59, sep 2019.
- [46] Alex Beutel, Jilin Chen, Tulsee Doshi, Hai Qian, Li Wei, Yi Wu, Lukasz Heldt, Zhe Zhao, Lichan Hong, Ed H. Chi, and Cristos Goodrow. Fairness in recommendation ranking through pairwise comparisons. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2212–2220, mar 2019. doi: 10.1145/3292500.3330745. URL <http://arxiv.org/abs/1903.00780>.
- [47] Harikrishna Narasimhan, Andrew Cotter, Maya Gupta, and Serena Wang. Pairwise fairness for ranking and regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5248–5255, 2020.
- [48] Robert C. Williamson and Aditya Krishna Menon. Fairness risk measures. In *36th International Conference on Machine Learning, ICML 2019*, volume 2019-June, pages 11763–11774. International Machine Learning Society (IMLS), jan 2019. ISBN 9781510886988.
- [49] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q. Weinberger. On fairness and calibration. In *Advances in Neural Information Processing Systems*, pages 5680–5689, sep 2017.

- [50] Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*, 2018.
- [51] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 60–69, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [52] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, Krishna P Gummadi, Manuel Gomez Roriguez, and Krishna P Gummadi. Fairness Constraints: Mechanisms for Fair Classification. In Aarti Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 962–970, Fort Lauderdale, FL, USA, 2017. PMLR.
- [53] Andrew Cotter, Maya Gupta, Heinrich Jiang, Nathan Srebro, Karthik Sridharan, Serena Wang, Blake Woodworth, and Seungil You. Training Well-Generalizing Classifiers for Fairness Metrics and Other Data-Dependent Constraints. In *International Conference on Machine Learning*, pages 1397–1405, jun 2019. URL <http://arxiv.org/abs/1807.00028>.
- [54] Richard S Zemel, Yu Wu, Kevin Swersky, Toniann Pitassi, and Cynthia Dwork. Learning Fair Representations. *Proceedings of the 30th International Conference on Machine Learning*, 28:325–333, 2013. ISSN 1938-7228.
- [55] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The Variational Fair Autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.
- [56] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning Adversarially Fair and Transferable Representations. *Proceedings of the 35th International Conference on Machine Learning*, 80:3384–3393, feb 2018. ISSN 1938-7228. URL <http://proceedings.mlr.press/v80/madras18a.html><http://arxiv.org/abs/1802.06309>.
- [57] Jiaming Song, Pratyusha Kalluri, Aditya Grover, Shengjia Zhao, and Stefano Ermon. Learning controllable fair representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2164–2173. PMLR, 2019.
- [58] Christina Ilvento. Metric learning for individual fairness. In *1st Symposium on Foundations of Responsible Computing (FORC 2020)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2020.
- [59] Blake Woodworth, Suriya Gunasekar, Mesrob I Ohannessian, and Nathan Srebro. Learning Non-Discriminatory Predictors. In Satyen Kale and Ohad Shamir, editors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pages 1920–1953, Amsterdam, Netherlands, 2017. PMLR.
- [60] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally Robust Neural Networks. In *International Conference on Learning Representations*, 2020.
- [61] Emily Diana, Wesley Gill, Michael Kearns, Krishnamurthy Suresh, and Aaron Roth. Minimax group fairness: Algorithms and experiments. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2021.
- [62] Natalia Martinez, Martin Bertran, and Guillermo Sapiro. Minimax pareto fairness: A multi objective perspective. In *International Conference on Machine Learning*, pages 6755–6764. PMLR, 2020.
- [63] Robert Chen, Brendan Lucier, Yaron Singer, and Vasilis Syrgkanis. Robust optimization for non-convex objectives. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 4708–4717, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- [64] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.

- [65] Aharon Ben-Tal, Dick Den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2): 341–357, 2013.
- [66] Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037. PMLR, 2018.
- [67] Steve Yadlowsky, Rodney A Hayward, Jeremy B Sussman, Robyn L McClelland, Yuan-I Min, and Sanjay Basu. Clinical implications of revised pooled cohort equations for estimating atherosclerotic cardiovascular disease risk. *Annals of internal medicine*, 169(1):20–29, 2018.
- [68] Andrew P. DeFilippis, Rebekah Young, Christopher J. Carrubba, John W. McEvoy, Matthew J. Budoff, Roger S. Blumenthal, Richard A. Kronmal, Robyn L. McClelland, Khurram Nasir, and Michael J. Blaha. An Analysis of Calibration and Discrimination Among Multiple Cardiovascular Risk Scores in a Modern Multiethnic Cohort. *Annals of Internal Medicine*, 162(4):266, feb 2015. ISSN 0003-4819. doi: 10.7326/M14-1281.
- [69] Nancy R. Cook and Paul M. Ridker. Calibration of the Pooled Cohort Equations for Atherosclerotic Cardiovascular Disease. *Annals of Internal Medicine*, 165(11):786, dec 2016. ISSN 0003-4819. doi: 10.7326/M16-1739. URL <http://annals.org/article.aspx?doi=10.7326/M16-1739>.
- [70] Michael J Pencina, Ann Marie Navar-Boggan, Ralph B D’Agostino Sr, Ken Williams, Benjamin Neely, Allan D Sniderman, and Eric D Peterson. Application of new cholesterol guidelines to a population-based sample. *N Engl J Med*, 370:1422–1431, 2014.
- [71] Jamal S Rana, Grace H Tabada, Matthew D Solomon, Joan C Lo, Marc G Jaffe, Sue Hee Sung, Christie M Ballantyne, and Alan S Go. Accuracy of the Atherosclerotic Cardiovascular Risk Equation in a Large Contemporary, Multiethnic Population. *Journal of the American College of Cardiology*, 67(18):2118–2130, may 2016. ISSN 15583597. doi: 10.1016/j.jacc.2016.02.055.
- [72] Andrew Paul DeFilippis, Rebekah Young, John W McEvoy, Erin D Michos, Veit Sandfort, Richard A Kronmal, Robyn L McClelland, and Michael J Blaha. Risk score overestimation: the impact of individual cardiovascular risk factors and preventive therapies on the performance of the american heart association-american college of cardiology-atherosclerotic cardiovascular disease risk score in a modern multi-ethnic cohort. *European heart journal*, 38(8):598–608, 2017.
- [73] Keum Ji Jung, Yangsoo Jang, Dong Joo Oh, Byung-Hee Oh, Sang Hoon Lee, Seong-Wook Park, Ki-Bae Seung, Hong-Kyu Kim, Young Duk Yun, Sung Hee Choi, et al. The acc/aha 2013 pooled cohort equations compared to a korean risk prediction model for atherosclerotic cardiovascular disease. *Atherosclerosis*, 242(1):367–375, 2015.
- [74] Maryam Afkarian, Ronit Katz, Nisha Bansal, Adolfo Correa, Bryan Kestenbaum, Jonathan Himmelfarb, Ian H De Boer, and Bessie Young. Diabetes, kidney disease, and cardiovascular outcomes in the jackson heart study. *Clinical Journal of the American Society of Nephrology*, 11(8):1384–1391, 2016.
- [75] Samia Mora, Nanette K Wenger, Nancy R Cook, Jingmin Liu, Barbara V Howard, Marian C Limacher, Simin Liu, Karen L Margolis, Lisa W Martin, Nina P Paynter, et al. Evaluation of the pooled cohort risk equations for cardiovascular risk prediction in a multiethnic cohort from the women’s health initiative. *JAMA internal medicine*, 178(9):1231–1240, 2018.
- [76] Terry A Jacobson, Matthew K Ito, Kevin C Maki, Carl E Orringer, Harold E Bays, Peter H Jones, James M McKenney, Scott M Grundy, Edward A Gill, Robert A Wild, et al. National lipid association recommendations for patient-centered management of dyslipidemia: part 1—full report. *Journal of clinical lipidology*, 9(2):129–169, 2015.
- [77] Charles R Harper and Terry A Jacobson. Managing dyslipidemia in chronic kidney disease. *Journal of the American College of Cardiology*, 51(25):2375–2384, 2008.

- [78] Gulsen Ozen, Murat Sunbul, Pamir Atagunduz, Haner Direskeneli, Kursat Tigen, and Nevsun Inanc. The 2013 acc/aha 10-year atherosclerotic cardiovascular disease risk index is better than score and risk ii in rheumatoid arthritis: is it enough? *Rheumatology*, 55(3):513–522, 2016.
- [79] Inge A.M. van den Oever, Alper M. van Sijl, and Michael T. Nurmohamed. Management of cardiovascular risk in patients with rheumatoid arthritis: evidence and expert opinion. *Therapeutic Advances in Musculoskeletal Disease*, 5(4):166, 2013. doi: 10.1177/1759720X13491025.
- [80] Andrew Ward, Ashish Sarraju, Sukyung Chung, Jiang Li, Robert Harrington, Paul Heidenreich, Latha Palaniappan, David Scheinker, and Fatima Rodriguez. Machine learning and atherosclerotic cardiovascular disease risk prediction in a multi-ethnic population. *npj Digital Medicine*, 3(1):1–7, dec 2020. ISSN 23986352. doi: 10.1038/s41746-020-00331-1. URL <https://doi.org/10.1038/s41746-020-00331-1>.
- [81] Ioannis A Kakadiaris, Michalis Vrigkas, Albert A Yen, Tatiana Kuznetsova, Matthew Budoff, and Morteza Naghavi. Machine learning outperforms acc/aha cvd risk calculator in mesa. *Journal of the American Heart Association*, 7(22):e009476, 2018.
- [82] Yuan Zhao, Erica P. Wood, Nicholas Mirin, Stephanie H. Cook, and Rumi Chunara. Social Determinants in Machine Learning Cardiovascular Disease Prediction Models: A Systematic Review. *American Journal of Preventive Medicine*, 0(0):1–10, jul 2021. ISSN 0749-3797. doi: 10.1016/J.AMEPRE.2021.04.016.
- [83] James M. Robins and Andrea Rotnitzky. Recovery of Information and Adjustment for Dependent Censoring Using Surrogate Markers. In *AIDS Epidemiology*, pages 297–331. Birkhäuser Boston, 1992. doi: 10.1007/978-1-4757-1229-2\_14.
- [84] Annette M. Molinaro, Sandrine Dudoit, and Mark J. Van Der Laan. Tree-based multivariate regression and density estimation with right-censored data. *Journal of Multivariate Analysis*, 90(1 SPEC. ISS.): 154–177, jul 2004. ISSN 10957243. doi: 10.1016/j.jmva.2004.02.003.
- [85] Mark J Van der Laan, MJ Laan, and James M Robins. *Unified methods for censored longitudinal data and causality*. Springer Science & Business Media, 2003.
- [86] Paul Blanche, Jean-François Dartigues, and H el ene Jacqmin-Gadda. Review and comparison of roc curve estimators for a time-dependent outcome with marker-dependent censoring. *Biometrical Journal*, 55(5):687–704, 2013.
- [87] Hajime Uno, Tianxi Cai, Lu Tian, and Lee-Jen J. Wei. Evaluating prediction rules for t-year survivors with censored regression models. *Journal of the American Statistical Association*, 102(478):527–537, jun 2007. ISSN 01621459. doi: 10.1198/016214507000000149.
- [88] Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- [89] James M Robins and Dianne M Finkelstein. Correcting for noncompliance and dependent censoring in an aids clinical trial with inverse probability of censoring weighted (ipcw) log-rank tests. *Biometrics*, 56(3):779–788, 2000.
- [90] H avard Kvamme and  rnulf Borgan. Continuous and discrete-time survival prediction with neural networks. *arXiv preprint arXiv:1910.06724*, 2019.
- [91] Handrean Soran, Jonathan D Schofield, and Paul N Durrington. Cholesterol, not just cardiovascular risk, is important in deciding who should receive statin treatment. *European heart journal*, 36(43): 2975–2983, 2015.
- [92] Cholesterol Treatment Trialists et al. Efficacy and safety of cholesterol-lowering treatment: prospective meta-analysis of data from 90 056 participants in 14 randomised trials of statins. *The Lancet*, 366(9493):1267–1278, 2005.

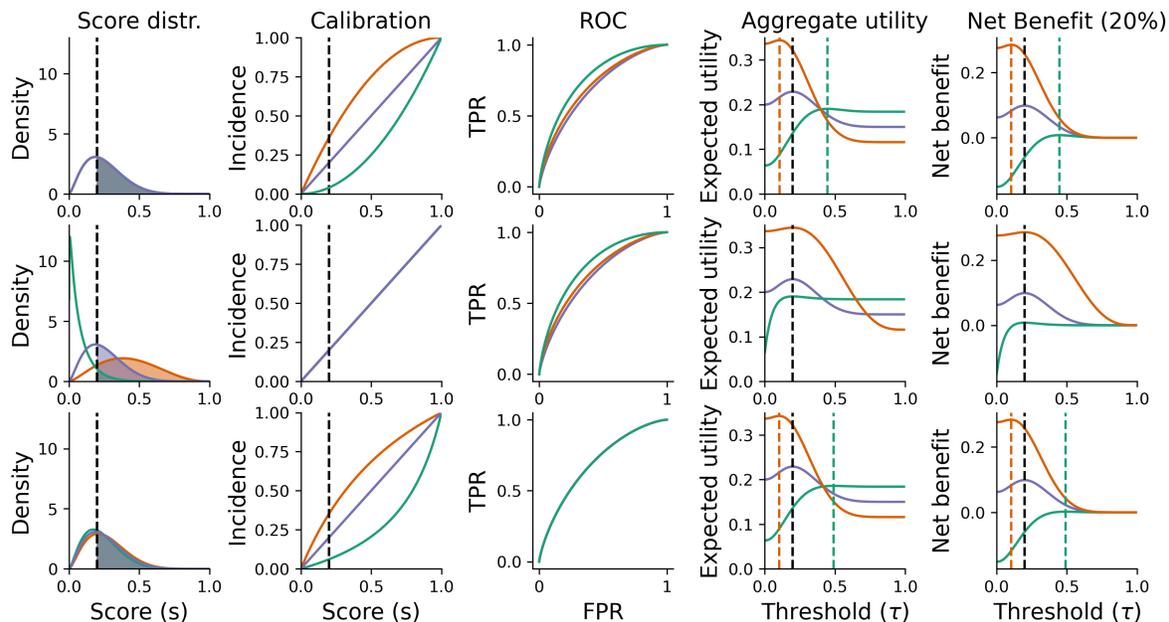
- [93] National Institute for Health and Care Excellence. Cardiovascular disease: risk assessment and reduction, including lipid modification (cg181)[online]. 2014.
- [94] Rory Collins, Christina Reith, Jonathan Emberson, Jane Armitage, Colin Baigent, Lisa Blackwell, Roger Blumenthal, John Danesh, George Davey Smith, David DeMets, et al. Interpretation of the evidence for the efficacy and safety of statin therapy. *The Lancet*, 388(10059):2532–2561, 2016.
- [95] George Hripcsak, Jon D Duke, Nigam H Shah, Christian G Reich, Vojtech Huser, Martijn J Schuemie, Marc A Suchard, Rae Woong Park, Ian Chi Kei Wong, Peter R Rijnbeek, Johan Van Der Lei, Nicole Pratt, G Niklas Norén, Yu-Chuan Chuan Li, Paul E Stang, David Madigan, and Patrick B Ryan. Observational Health Data Sciences and Informatics (OHDSI): Opportunities for Observational Researchers. In *Studies in Health Technology and Informatics*, volume 216, pages 574–578. NIH Public Access, 2015. ISBN 9781614995630. doi: 10.3233/978-1-61499-564-7-574.
- [96] J. Marc Overhage, Patrick B. Ryan, Christian G. Reich, Abraham G. Hartzema, and Paul E. Stang. Validation of a common data model for active safety surveillance research. *Journal of the American Medical Informatics Association*, 19(1):54–60, jan 2012. ISSN 10675027. doi: 10.1136/amiajnl-2011-000376.
- [97] Jenna M Reps, Martijn J Schuemie, Marc A Suchard, Patrick B Ryan, and Peter R Rijnbeek. Design and implementation of a standardized framework to generate and evaluate patient-level prediction models using observational healthcare data. *Journal of the American Medical Informatics Association*, 25(8):969–975, apr 2018. ISSN 1067-5027. doi: 10.1093/jamia/ocy032. URL <https://academic.oup.com/jamia/advance-article/doi/10.1093/jamia/ocy032/4989437>.
- [98] Jenna Reps and Peter Rijnbeek. Network study validating the Pooled Cohort Equation Model, 2020. URL <https://github.com/ohdsi-studies/PCE>.
- [99] Marc A Suchard, Martijn J Schuemie, Harlan M Krumholz, Seng Chan You, RuiJun Chen, Nicole Pratt, Christian G Reich, Jon Duke, David Madigan, George Hripcsak, et al. Comprehensive comparative effectiveness and safety of first-line antihypertensive drug classes: a systematic, multinational, large-scale analysis. *The Lancet*, 394(10211):1816–1826, 2019.
- [100] Michael F. Gensheimer and Balasubramanian Narasimhan. A scalable discrete-time survival model for neural networks, jan 2019. ISSN 21678359.
- [101] Gerhard Tutz, Matthias Schmid, et al. *Modeling discrete time-to-event data*. Springer, 2016.
- [102] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [103] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [104] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alch’e-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [105] Peter C. Austin and Ewout W. Steyerberg. The Integrated Calibration Index (ICI) and related metrics for quantifying the calibration of logistic regression models. *Statistics in Medicine*, 38(21):4051–4065, sep 2019. ISSN 10970258. doi: 10.1002/sim.8281. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/sim.8281>.
- [106] Steve Yadlowsky, Sanjay Basu, and Lu Tian. A calibration metric for risk scores with survival data. In *Machine Learning for Healthcare Conference*, pages 424–450, 2019.

- [107] Zinzi D. Bailey, Nancy Krieger, Madina Agénor, Jasmine Graves, Natalia Linos, and Mary T. Bassett. Structural racism and health inequities in the USA: evidence and interventions. *The Lancet*, 389(10077):1453–1463, apr 2017. ISSN 1474547X. doi: 10.1016/S0140-6736(17)30569-X. URL <https://www-thelancet-com/series/america-equity-equality-in-health>.
- [108] Zinzi D. Bailey, Justin M. Feldman, and Mary T. Bassett. How Structural Racism Works — Racist Policies as a Root Cause of U.S. Racial Health Inequities. <https://doi.org/10.1056/NEJMms2025396>, 384(8):768–773, dec 2020. doi: 10.1056/NEJMMS2025396. URL <https://www.nejm.org/doi/10.1056/NEJMms2025396>.
- [109] World Health Organization et al. A conceptual framework for action on the social determinants of health. 2010.
- [110] Pratyusha Kalluri. Don’t ask if artificial intelligence is good or fair, ask how it shifts power. *Nature*, 583(7815):169–169, jul 2020. ISSN 0028-0836. doi: 10.1038/d41586-020-02003-2. URL <http://www.nature.com/articles/d41586-020-02003-2>.
- [111] Steven N. Goodman, Sharad Goel, and Mark R. Cullen. Machine Learning, Health Disparities, and Causal Reasoning. *Annals of Internal Medicine*, 169(12):883, dec 2018. ISSN 0003-4819. doi: 10.7326/M18-3297. URL <http://annals.org/article.aspx?doi=10.7326/M18-3297>.
- [112] Úrsula Hébert-Johnson, Michael P. Kim, Omer Reingold, and Guy N. Rothblum. Calibration for the (Computationally-Identifiable) Masses. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1939–1948, Stockholmsmässan, Stockholm Sweden, 2017. PMLR.
- [113] Diana Cagliero, Natalie Deutch, Nigam Shah, and Danton Char. Evaluating ethical concerns with machine learning to guide advance care planning. In *2021 Western Medical Research Conference*, volume 69, pages 103–296. BMJ Publishing Group Limited, January 2021.
- [114] Elad Eban, Mariano Schain, Alan Mackey, Ariel Gordon, Ryan Rifkin, and Gal Elidan. Scalable learning of non-decomposable objectives. In *Artificial intelligence and statistics*, pages 832–840. PMLR, 2017.

# Supplementary material

## A Supplementary methods

### A.1 Simulation study



**Supplementary Figure A1:** The relationship between score distributions, calibration, receiver operating characteristic curves, aggregate expected utility, and net benefit in simulation when demographic parity (first row), group calibration (second row), and equalized odds (third row) are satisfied. Dashed lines indicate optimal decision thresholds.

We conduct a simulation study in order to help clarify the relationship between the key concepts considered in this work. We consider evaluation of a predictive model of a binary outcome without censoring using the notion of utility and net benefit discussed in section 2.3. For evaluation in all settings, we evaluate aggregate expected utility using the fixed-cost utility function with  $u_{TP} = 0.8$ ,  $u_{TN} = 0.2$ , and  $u_{FP} = u_{FN} = 0$ , corresponding to an optimal threshold of 0.2, as well as the net benefit parameterized by  $\tau^* = 0.2$ .

We evaluate over three subgroups in three different settings (Supplementary Figure A1). In the first setting, we assume the model satisfies demographic parity and the score distribution is given by  $S \sim \text{Beta}(2.5, 7.5)$ , corresponding to an outcome incidence of 25%. We assume that the model is perfectly calibrated for one subgroup (purple) and that risk is systematically underestimated (orange) or overestimated (green) for the other two subgroups. We encode the miscalibration by setting  $c(s) = -(s - 1)^2 + 1$  for underestimation and  $c(s) = 2s + (s - 1)^2 - 1$  for overestimation. For the second setting, we consider a hypothetical recalibration procedure for each subgroup that constructs a new set of scores by setting the value of the score to be that of the subgroup calibration curve, propagating the changes to the score distributions using change of variables for probability density functions. For the third context, we consider the effect of enforcing an equalized odds constraint by setting the conditional distributions  $P(S | Y)$  for each subgroup to be equal to that of the subgroup for which the model was perfectly calibrated in the first setting, where demographic parity was satisfied, as the ROC curve for that subgroup defines the convex hull of the ROC curves across all subgroups. In this case, the resulting changes to the score distributions and calibration curves are computed using simple conditional probability rules.

The results of the simulation study are consistent with the presentation of section 2.4. In particular, we note that when the model is calibrated for each subgroup, the utility and net benefit maximizing decision

threshold is the same for each subgroup. However, when demographic parity or equalized odds are satisfied and sufficiency is violated, the optimal threshold differs from 0.2 for the subgroups for which the model is miscalibrated. We note that satisfying equalized odds is similar to that of demographic parity in terms of the effect on the score distributions, aggregate utility, and net benefit.

## A.2 Assessing net benefit in terms of risk reduction

This section expands upon the presentation in section 3.2.2 and includes full derivations for equations (13) and (14). We present a conceptual model that reasons about the utility of an intervention allocated on the basis of a decision rule applied to a predictive model in terms of a downstream reduction in the risk of the predicted outcome as a result of the intervention. We show that the use of this conceptual model results in similar conclusions as the fixed expected cost/utility setting.

We define  $u_1^y$  be the utility associated with the presence of the outcome  $Y$  and  $u_0^y$  be the utility associated with its absence. The probability of the outcome in the absence of the intervention is given by  $p_y^0(s) = c(s)$  where  $c(s)$  is the calibration curve. The probability of the outcome in the presence of intervention is given by  $p_y^1(s)$ , where the precise form of  $p_y^1(s)$  is governed by the effectiveness of the intervention. We further assume that there is some harm  $Z$ , representing all costs, harms, or side effects, that occurs with probability  $p_z^1(s)$  and utility  $u_1^z$  following the intervention and with  $p_z^0(s)$  and utility  $u_0^z$  in the absence of the intervention.

This formulation implies the conditional utilities

$$U_{\text{cond}}^0(s) = p_y^0(s)(u_1^y - u_0^y) + u_0^y + p_z^0(s)(u_1^z - u_0^z) + u_0^z, \quad (15)$$

$$U_{\text{cond}}^1(s) = p_y^1(s)(u_1^y - u_0^y) + u_0^y + p_z^1(s)(u_1^z - u_0^z) + u_0^z, \quad (16)$$

and

$$U_{\text{cond}}(s) = (u_0^y - u_1^y)(p_y^0(s) - p_y^1(s)) + (u_0^z - u_1^z)(p_z^0(s) - p_z^1(s)). \quad (17)$$

Setting equation (17) to zero shows that the value of the optimal threshold  $\tau_y^*$  is governed by the following relationship

$$(u_0^y - u_1^y)(p_y^0(\tau_y^*) - p_y^1(\tau_y^*)) = (u_0^z - u_1^z)(p_z^1(\tau_y^*) - p_z^0(\tau_y^*)) \rightarrow \frac{p_y^0(\tau_y^*) - p_y^1(\tau_y^*)}{p_z^1(\tau_y^*) - p_z^0(\tau_y^*)} = \frac{u_0^z - u_1^z}{u_0^y - u_1^y}, \quad (18)$$

when  $U_{\text{cond}}(s)$  is monotonically increasing in  $s$ . To interpret this expression, consider that  $\text{ARR}(s) = p_y^0(s) - p_y^1(s)$  is the absolute reduction in risk as a result of the intervention,  $p_z^0(s) - p_z^1(s)$  indicates a corresponding increase in the risk of harm, and  $u_0^y - u_1^y$  and  $u_0^z - u_1^z$  indicate the utilities associated with avoiding  $Y$  and  $Z$ , respectively. It follows that the optimal threshold is the one where the benefits of the intervention are balanced against its harms.

To simplify the model, we now assume that  $p_z^1(s) - p_z^0(s)$  do not depend on the risk score, indicating that the expected harm  $k_{\text{harm}} = (u_0^z - u_1^z)(p_z^1 - p_z^0)$  is a constant.

With this assumption, the conditional utility may be represented as

$$U_{\text{cond}}(s) = (u_0^y - u_1^y)(p_y^0(s) - p_y^1(s)) - k_{\text{harm}}. \quad (19)$$

Setting equation (19) to zero shows that the value of the optimal threshold  $\tau_y^*$  is governed by the following relationship, consistent with Vickers et al. [35]:

$$p_y^0(\tau_y^*) - p_y^1(\tau_y^*) = \frac{k_{\text{harm}}}{u_0^y - u_1^y}. \quad (20)$$

This expression relates the absolute risk reduction  $\text{ARR}(s) = p_y^0(s) - p_y^1(s)$  evaluated at the optimal threshold  $\tau_y^*$  to both the expected harm of intervention  $k_{\text{harm}}$  and the utility of avoiding the outcome  $u_0^y - u_1^y$ . It follows that the optimal threshold  $\tau_y^*$  is given by

$$\tau_y^* = \text{ARR}^{-1}\left(\frac{k_{\text{harm}}}{u_0^y - u_1^y}\right). \quad (21)$$

Furthermore, the aggregate utility over the population when treating at a threshold  $\tau_y$  can be derived as

$$\begin{aligned} U_{\text{agg}}(\tau_y) &= \left(u_1^y - u_0^y\right) \left( \int_0^{\tau_y} p_y^0(s)P(s)ds + \int_{\tau_y}^1 p_y^1(s)P(s)ds \right) - k_{\text{harm}} + p_z^0(u_1^z - u_0^z) + u_0^y + u_0^z \\ &= \left(u_1^y - u_0^y\right) \left( \mathbb{E}[p_y^0(s) | S < \tau_y]P(S < \tau_y) + \mathbb{E}[p_y^1(s) | S \geq \tau_y]P(S \geq \tau_y) \right) \dots \\ &\quad - k_{\text{harm}} + p_z^0(u_1^z - u_0^z) + u_0^y + u_0^z. \end{aligned} \quad (22)$$

To construct a net benefit measure that represents the aggregate utility given that  $\tau_y$  is the optimal threshold, we divide equation (22) by  $u_0^y - u_1^y$ , perform a substitution following equation (20), and define a constant  $k$  such that the net benefit of the treat-none strategy is zero:

$$\text{NB}(\tau_y; \tau_y^*) = -\mathbb{E}[p_y^0(s) | S < \tau_y]P(S < \tau_y) - \mathbb{E}[p_y^1(s) | S \geq \tau_y]P(S \geq \tau_y) - \text{ARR}(\tau_y^*)P(S \geq \tau_y) + k. \quad (23)$$

From this expression, it follows that the appropriate value of  $k$  is given by  $P(Y = 1) = \mathbb{E}[p_y^0(s) | S < 1]P(S < 1)$ , giving the following expression for the net benefit:

$$\begin{aligned} \text{NB}(\tau_y; \tau_y^*) &= -\mathbb{E}[p_y^0(s) | S < \tau_y]P(S < \tau_y) - \mathbb{E}[p_y^1(s) | S \geq \tau_y]P(S \geq \tau_y) \dots \\ &\quad - \text{ARR}(\tau_y^*)P(S \geq \tau_y) + P(Y = 1). \end{aligned} \quad (24)$$

The expression for the calibrated variant of the net benefit is given by:

$$\begin{aligned} \text{cNB}(\tau_y; \tau_y^*) &= -\mathbb{E}[p_y^0(s) | S < c^{-1}(\tau_y)]P(S < c^{-1}(\tau_y)) - \mathbb{E}[p_y^1(s) | S \geq c^{-1}(\tau_y)]P(S \geq c^{-1}(\tau_y)) \dots \\ &\quad - \text{ARR}(\tau_y^*)P(S \geq c^{-1}(\tau_y)) + P(Y = 1). \end{aligned} \quad (25)$$

This formulation differs from that of Vickers et al. [35] in that that work defines the treat-all strategy as having a net benefit of zero whereas we do so for the treat-none strategy in order to maintain consistency with the net benefit defined for the fixed-cost utility function.

### A.2.1 Constant relative risk reduction

Given only observational data that corresponds to an untreated population, it is necessary to provide assumptions on the form of  $p_y^1(s)$  in order to assess the net benefit of a model using a utility function defined in terms of the risk reduction induced by the intervention. A simple choice for the relationship between  $p_y^0(s)$  and  $p_y^1(s)$  is one where the intervention reduces the risk of the outcome by a constant multiplicative factor  $r \in (0, 1)$ , such that  $p_y^1(s) = (1 - r)p_y^0(s)$ . Given this assumption,  $p_y^0(s) = c(s)$ ,  $p_y^1(s) = (1 - r)c(s)$ , and  $\text{ARR}(s) = rc(s)$ .

With these assumptions, it follows that  $U_{\text{cond}}(s)$  is a linear transformation of the calibration curve, just as was the case for the fixed-cost utility function:

$$U_{\text{cond}}(s) = (u_0^y - u_1^y)rc(s) - k_{\text{harm}}. \quad (26)$$

Furthermore, the optimal threshold is given by

$$\tau_y^* = c^{-1}\left(\frac{k_{\text{harm}}}{r(u_0^y - u_1^y)}\right), \quad (27)$$

which can be simplified to

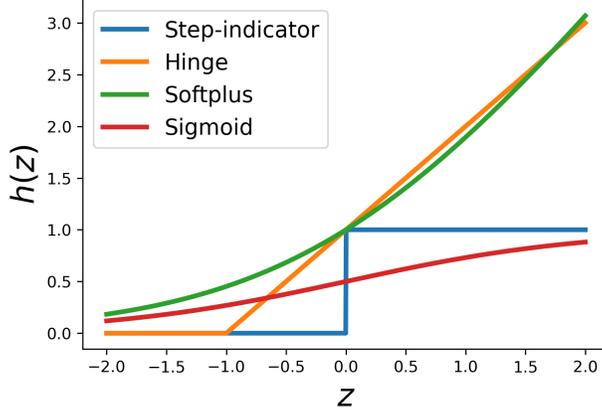
$$\tau_y^* = \frac{k_{\text{harm}}}{r(u_0^y - u_1^y)} \quad (28)$$

when the model is assumed to be calibrated.

Making the appropriate substitutions into equation (24), and noting that  $\mathbb{E}[c(s) | S < \tau_y] = \mathbb{E}[Y | S < \tau_y]$  and  $\mathbb{E}[c(s) | S \geq \tau_y] = \mathbb{E}[Y | S \geq \tau_y]$ , gives an expression for the net benefit:

$$\begin{aligned} \text{NB}(\tau_y; \tau_y^*) &= -\mathbb{E}[Y | S < \tau_y]P(S < \tau_y) - P(S \geq \tau_y) \left( (1 - r) \mathbb{E}[Y | S \geq \tau_y] + \text{ARR}(\tau_y^*) \right) + P(Y = 1) \\ &= -(1 - \text{NPV}(\tau_y))P(S < \tau_y) - P(S \geq \tau_y) \left( (1 - r)\text{PPV}(\tau_y) + r\tau_y^* \right) + P(Y = 1) \end{aligned} \quad (29)$$

where  $\text{NPV}(\tau_y)$  and  $\text{PPV}(\tau_y)$  designate the negative and positive predictive values that result from operating at the decision threshold  $\tau_y$ . Note that this matches equation (13) and that the calibrated variant in equation (14) is analogous.



Supplementary Figure A2: Surrogates to the indicator function.

### A.3 Regularized training objectives for fairness

Here, we present the regularized training objective that allows for the flexible specification of penalties on the violation of fairness criteria. To begin, we consider the MMD-based penalty for equalized odds presented in equation (8). The MMD uses the distance between the mean embedding of samples from two distributions in a kernel space to define a statistic that takes a value of zero in a population setting if and only if two distributions are the same [64]. To construct a regularizer, we use an empirical estimate of the squared population MMD [64]

$$\begin{aligned} \hat{D}_{\text{MMD}}(\mathcal{D}_0 \parallel \mathcal{D}_1) &= \mathbb{E}_{(z, z') \sim \mathcal{D}_0, \mathcal{D}_0} [k(z, z')] - \\ &2 \mathbb{E}_{(z, z') \sim \mathcal{D}_0, \mathcal{D}_1} [k(z, z')] + \\ &\mathbb{E}_{(z, z') \sim \mathcal{D}_1, \mathcal{D}_1} [k(z, z')], \end{aligned} \quad (30)$$

where  $k(z, z') = \exp(-\gamma \|z - z'\|)$  is the Gaussian Radial Basis Function kernel defined for a positive scalar hyperparameter  $\gamma$ , and  $\mathbb{E}_{(z, z') \sim \mathcal{D}_0, \mathcal{D}_1}$  indicates an empirical mean following sampling a pair of data  $(z, z')$  from  $\mathcal{D}_0$  and  $\mathcal{D}_1$ , respectively. In our experiments, we set  $\gamma = 1$ . As noted in Gretton et al. [64], this statistic is non-negative, but has a small upward bias.

To account for censoring, we use a weighted extension of the maximum mean discrepancy, as a modification to each of the expectations over the pairwise evaluation of the kernel function (equation 30). As an example, the term  $\mathbb{E}_{(z, z') \sim \mathcal{D}_0, \mathcal{D}_1} [k(z, z')]$  can be replaced with  $\sum_{z_i, z_j \in \{\mathcal{D}_0, \mathcal{D}_1\}} w_{ij} k(z_i, z_j)$  for weights defined as

$$w_{ij} = \frac{\delta_i^y}{G(u_i^y, x_i)} \frac{\delta_j^y}{G(u_j^y, x_j)} \left( \sum_{z_i \in \mathcal{D}_0} \sum_{z_j \in \mathcal{D}_1} \frac{\delta_i^y}{G(u_i^y, x_i)} \frac{\delta_j^y}{G(u_j^y, x_j)} \right)^{-1}. \quad (31)$$

To define the full MMD, each of the three expectations in equation (30) are replaced with weighted variants analogous to equation (31).

Now, we consider the operationalization of equation (9) to penalize differences in model performance metrics. Recall that equation (9) is well-suited to penalties that rely on smooth and differentiable  $g_j$ . Unfortunately, naively plugging-in threshold-based metrics, including the classification rate, true positive rate, false positive rate, PPV, as well as those defined as ranking performance, including the AUC, does not produce a practical regularized objective due to the presence of the indicator function embedded in the definition of each of those metrics. As an example, consider that the classification rate  $\mathbb{E}[\mathbb{1}[f_\theta(X) > \tau_y]]$  can be represented as  $\mathbb{E}[\mathbb{1}[f_\theta(X) - \tau_y > 0]]$  or  $\mathbb{E}[h_{\text{step}}(f_\theta(X) - \tau_y)]$  when  $h_{\text{step}}$  is the step-indicator function  $h_{\text{step}}(z) = \mathbb{1}[z > 0]$ . The shape of this function is such that it does not provide a useful signal for stochastic gradient descent, given that its derivative is zero everywhere that its derivative is defined.

One approach to addressing this issue is to use a smooth and differentiable surrogate [45, 114] to the step-indicator that either upper bounds or approximates it. A visual depiction of several options is

provided in Figure A2. The use of either the hinge,  $h_{\text{hinge}}(z) = \max(0, 1 + z)$ , or the scaled softplus,  $h_{\text{softplus}}(z) = \log(1 + \exp(z))/\log(2)$ , provides a smooth and differentiable upper bound to the indicator. Because  $h_{\text{step}}(z) \leq h_{\text{hinge}}(z)$  and  $h_{\text{step}}(z) \leq h_{\text{softplus}}(z)$ , a metric  $g$  defined as a sum over evaluations of  $h_{\text{step}}(z)$  can be upper bounded by a metric  $\hat{g}$  defined as a sum over evaluations of  $h_{\text{hinge}}(z)$  or  $h_{\text{softplus}}(z)$ . Furthermore, the use of the sigmoid function,  $h_{\text{sigmoid}}(z) = \frac{1}{(1 + \exp(-z))}$  does not directly bound the indicator function, but rather provides a smooth approximation to it (Figure A2), that can be similarly incorporated into a relaxed, approximate metric  $\hat{g}$ .

Given a relaxed metric  $\hat{g}$ , the corresponding training objective is given by

$$\min_{\theta \in \Theta} \sum_{i=1}^N w_i \ell(y, f_{\theta}(x)) + \lambda \sum_{j=1}^J \sum_{A_k \in \mathcal{A}} \left( \hat{g}_j(f_{\theta}, \mathcal{D}_{A_k}) - \hat{g}_j(f_{\theta}, \mathcal{D}) \right)^2. \quad (32)$$

With this relaxed objective, it is straightforward to penalize differences in threshold-based performance metrics, such as the true positive rate and false positive rates, or to penalize differences in AUC measures. The true positive and false positive rates can each be written as  $\sum_{i=1}^N w_i h(f_{\theta}(x_i) - \tau_y)$  for the following weights [87]:

$$w_i = \frac{\mathbb{1}[y_i = 1] \delta_i^y}{G(u_i^y, x_i)} \left( \sum_{i=1}^N \frac{\mathbb{1}[y_i = 1] \delta_i^y}{G(u_i^y, x_i)} \right)^{-1} \quad (33)$$

for the true positive rate, and

$$w_i = \frac{\mathbb{1}[y_i = 0] \delta_i^y}{G(u_i^y, x_i)} \left( \sum_{i=1}^N \frac{\mathbb{1}[y_i = 0] \delta_i^y}{G(u_i^y, x_i)} \right)^{-1} \quad (34)$$

for the false positive rate. The corresponding censoring-adjusted definition of the AUC [86] that incorporates IPCW is given by  $\sum_{i=1}^N \sum_{j=1}^N w_{ij} h(f_{\theta}(x_i) - f_{\theta}(x_j))$  for

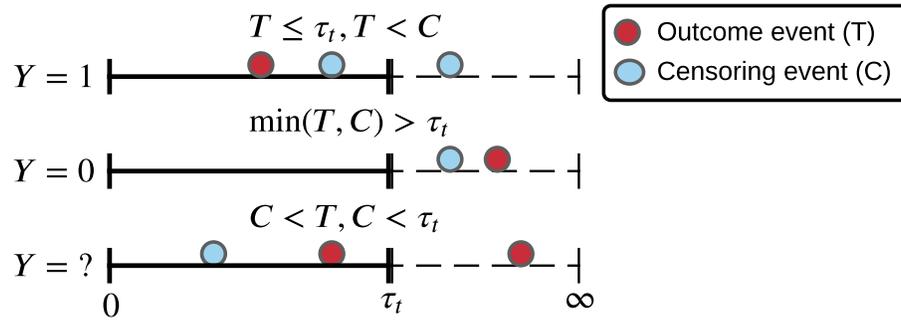
$$w_{ij} = \frac{\delta_i^y \mathbb{1}[y_i = 1]}{G(u_i^y, x_i)} \frac{\delta_j^y \mathbb{1}[y_j = 0]}{G(u_j^y, x_j)} \left( \sum_{i=1}^N \sum_{j=1}^N \frac{\delta_i^y \mathbb{1}[y_i = 1]}{G(u_i^y, x_i)} \frac{\delta_j^y \mathbb{1}[y_j = 0]}{G(u_j^y, x_j)} \right)^{-1}. \quad (35)$$

## B Supplementary tables

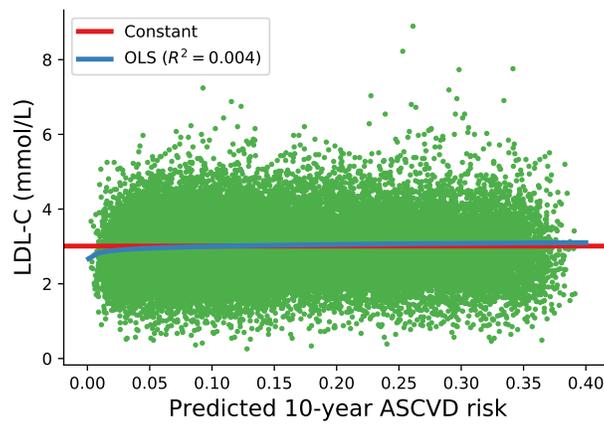
**Table B1:** Code and concept identifiers used to construct the cohort. Parentheses indicate the source vocabulary for the listed identifiers. We use the International Classification of Diseases version 9 (ICD-9), the Anatomical Therapeutic Chemical Classification System (ATC), Logical Observation Identifiers Names and Codes (LOINC), and OMOP CDM concept identifiers. Asterisks indicate the union of all possible suffixes and brackets indicate a range of included suffixes. For identifiers that are not OMOP CDM concept identifiers, we map the listed identifiers to standard OMOP CDM concepts using the mappings provided by the OMOP CDM vocabulary. Each set of OMOP CDM concepts used in the cohort definition is defined by the union of the mapped standard OMOP CDM concepts and their descendants in the OMOP CDM vocabulary followed by the exclusion of any excluded concepts and their descendants from the set.

Concept	Code or concept identifiers
Stroke (ICD-9)	430*, 431*, 432*, 433* (except 433.*0), 434* (except 434.*0), 436*
Myocardial Infarction (ICD-9)	410*
Coronary Heart Disease (ICD-9)	411*, 413*, 414*
Cardiovascular Disease (ICD-9)	410*, 411*, 413*, 414*, 430*, 431*, 432*, 433*, 434*, 436*, 427.31, 428*
Statin (ATC)	C10AA0[1-8]
Type 1 diabetes (OMOP)	201254, 40484648, 201254, 435216
Gestational diabetes (OMOP)	4058243
Type 2 diabetes (OMOP)	443238, 201820, 442793
Chronic kidney disease (OMOP)	(exclude all type 1 and gestational concepts) 192279, 192359, 193253, 194385, 195314, 201313, 261071, 4103224, 4263367, 46271022 (exclude 195014, 195289, 195737, 197320, 197930, 4066005, 37116834, 43530912, 45769152)
Rheumatoid Arthritis (OMOP)	80809
Low-density lipoprotein cholesterol (LOINC)	18262-6, 13457-7, 2089-1

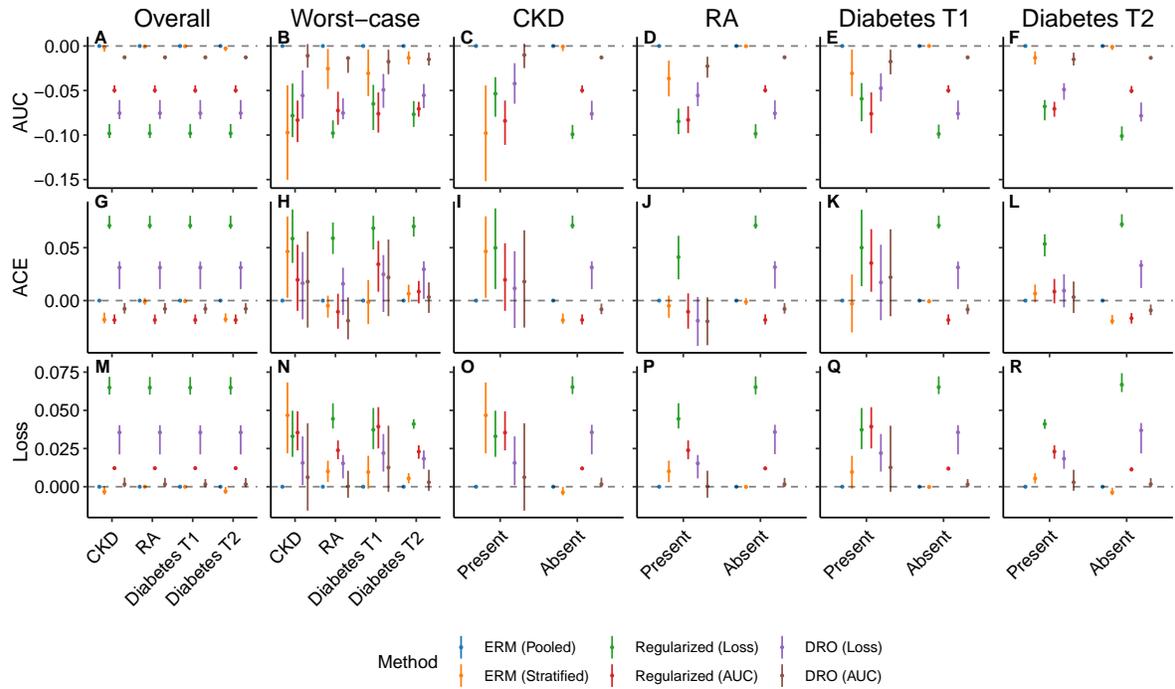
## C Supplementary figures



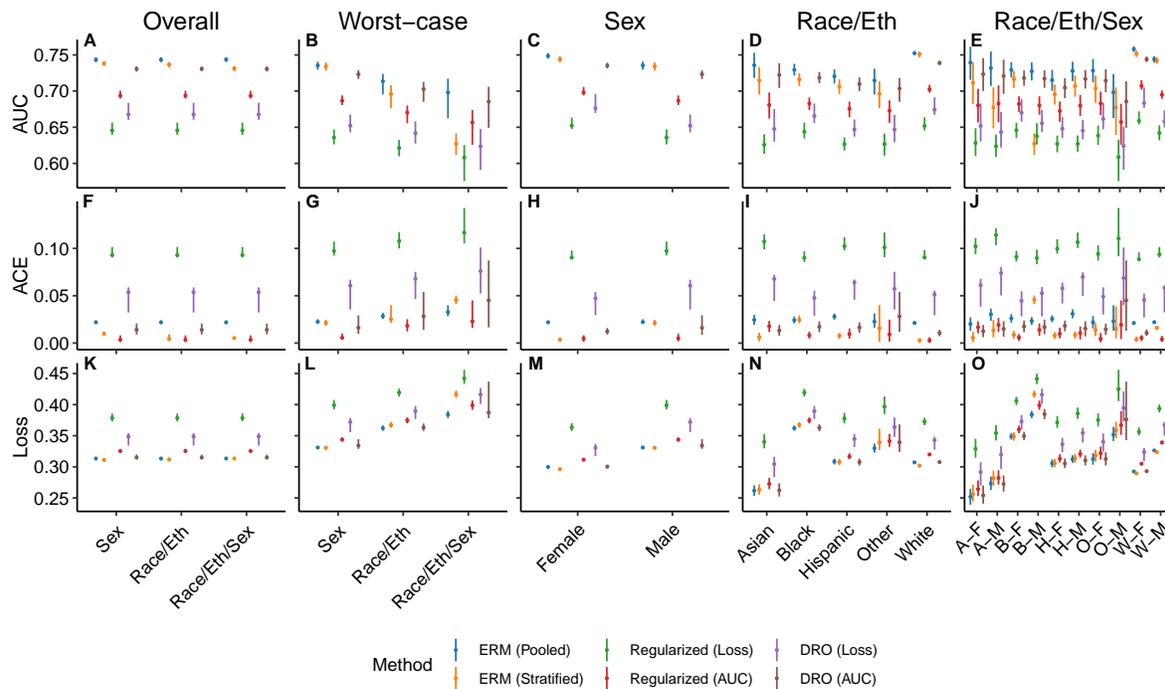
**Supplementary Figure C1:** The effect of censoring on the observation of binary outcomes



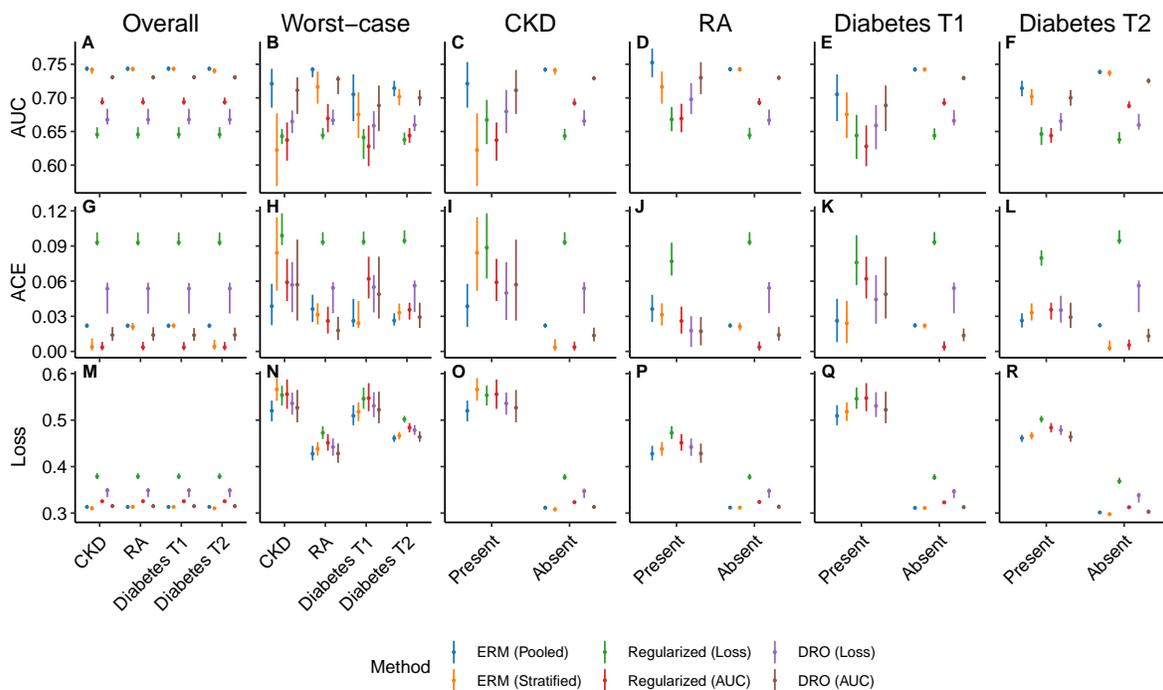
**Supplementary Figure C2:** The result of the most recent low density lipoprotein cholesterol (LDL-C) measurement versus the estimated risk of ASCVD within ten years.



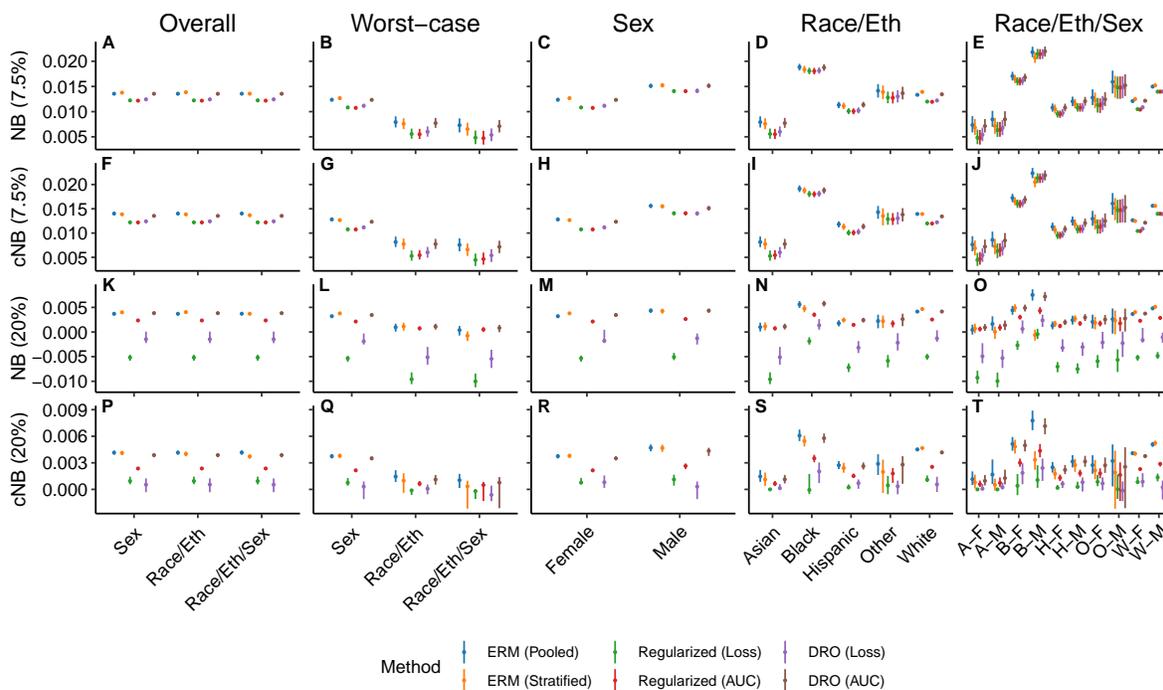
**Supplementary Figure C3:** The performance of models that estimate ten-year ASCVD risk, for subgroups defined by the presence or absence of chronic kidney disease (CKD), rheumatoid arthritis (RA), or type 1 (T1) or type 2 (T2) diabetes, relative to the results attained by the application of unpenalized ERM to the overall population. Results shown are the relative AUC, absolute calibration error (ACE), and log-loss assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unpenalized ERM, regularized objectives that penalize differences in the log-loss or AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss or AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



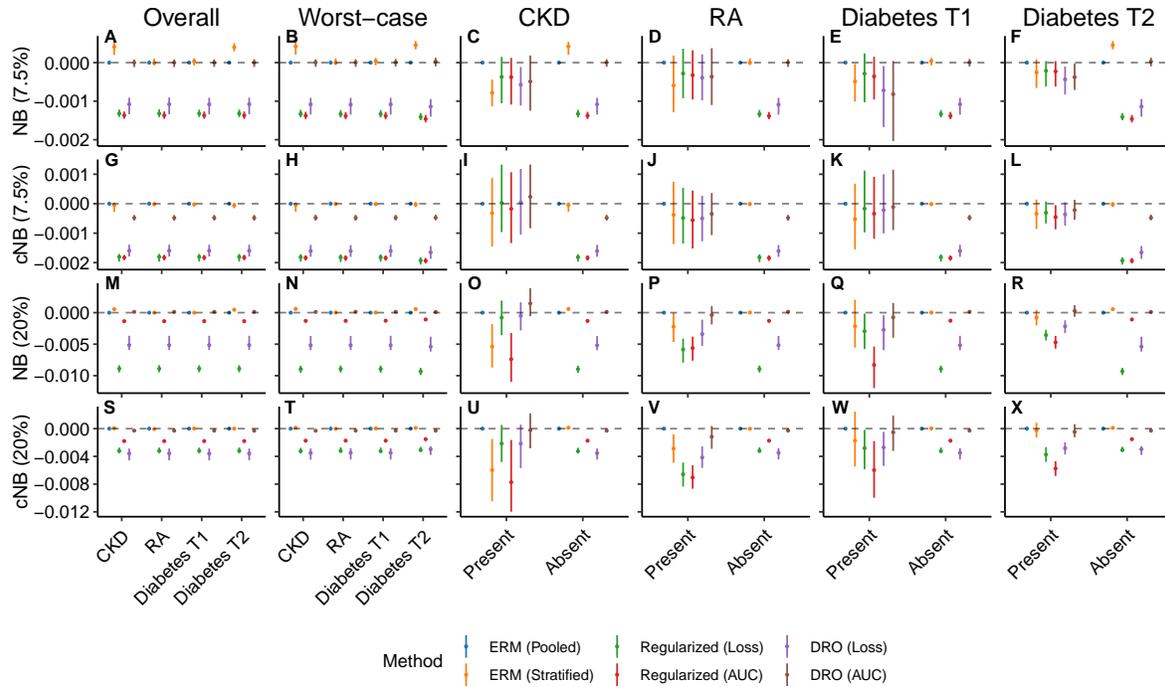
**Supplementary Figure C4:** The performance of models that estimate ten-year ASCVD risk, stratified by race, ethnicity, and sex. Results shown are the AUC, absolute calibration error (ACE), and log-loss assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unconstrained pooled or stratified ERM, regularized objectives that penalize differences in the log-loss of AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss of AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



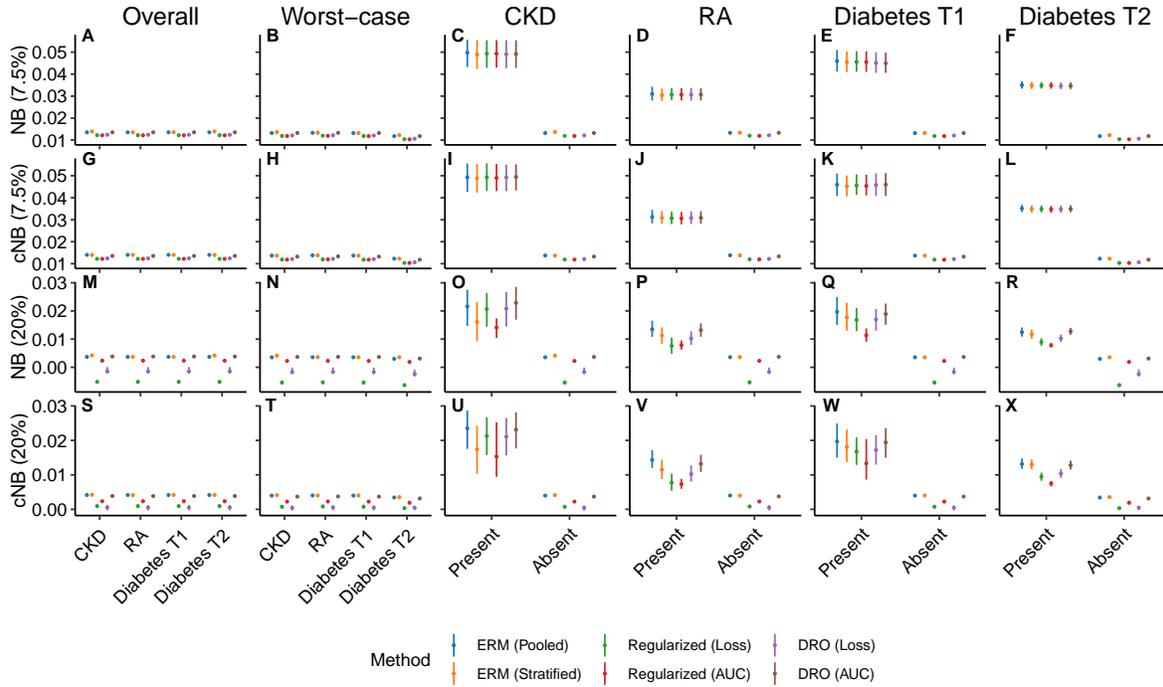
**Supplementary Figure C5:** The performance of models that estimate ten-year ASCVD risk for subgroups defined by the presence or absence of chronic kidney disease (CKD), rheumatoid arthritis (RA), or type 1 (T1) or type 2 (T2) diabetes. Results shown are the AUC, absolute calibration error (ACE), and log-loss assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unpenalized ERM, regularized objectives that penalize differences in the log-loss of AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss of AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



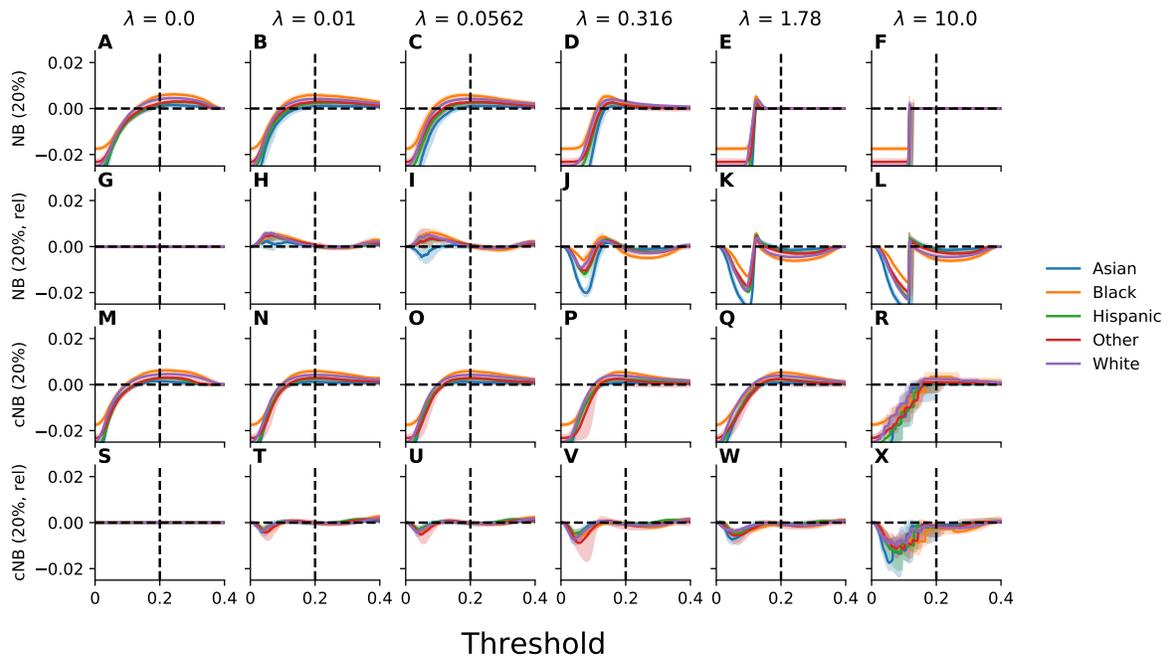
**Supplementary Figure C6:** The net benefit of models that estimate ten-year ASCVD risk, stratified by race, ethnicity, and sex. Results shown are the net benefit (NB) and calibrated net benefit (cNB), parameterized by the choice of a decision threshold of 7.5% or 20%, assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unconstrained pooled or stratified ERM, regularized objectives that penalize differences in the log-loss of AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss of AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



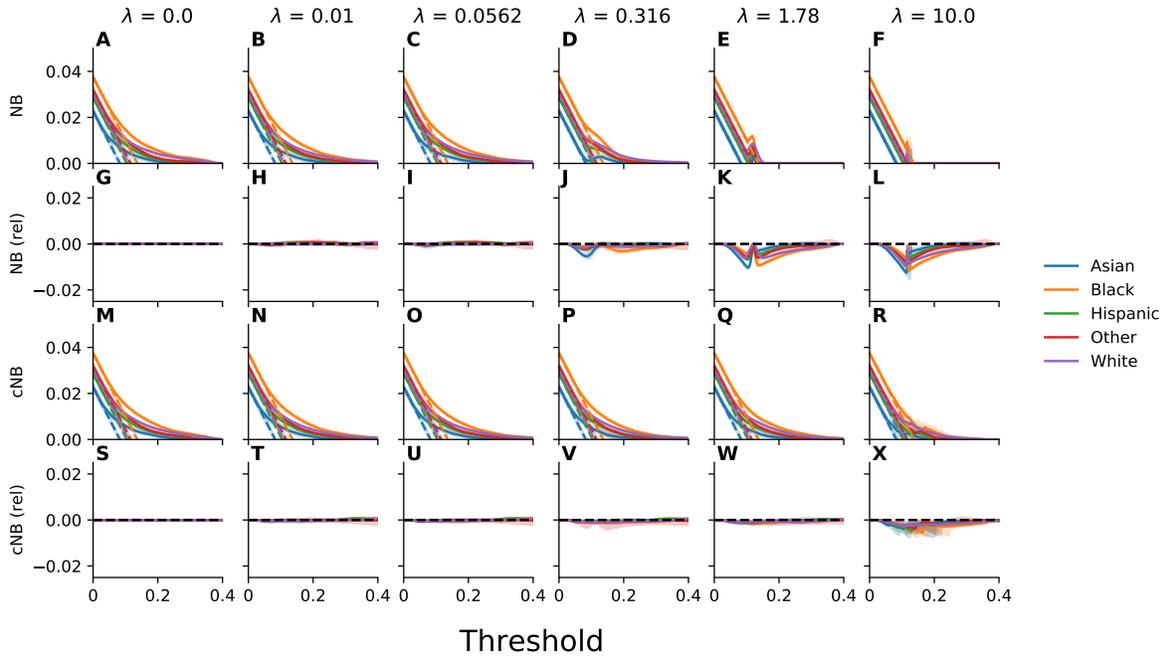
**Supplementary Figure C7:** The net benefit of models that estimate ten-year ASCVD risk, for subgroups defined by the presence or absence of chronic kidney disease (CKD), rheumatoid arthritis (RA), or type 1 (T1) or type 2 (T2) diabetes, relative to the results attained by the application of unpenalized ERM to the overall population. Results shown are the net benefit (NB) and calibrated net benefit (cNB), parameterized by the choice of a decision threshold of 7.5% or 20%, assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unconstrained pooled or stratified ERM, regularized objectives that penalize differences in the log-loss or AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss or AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



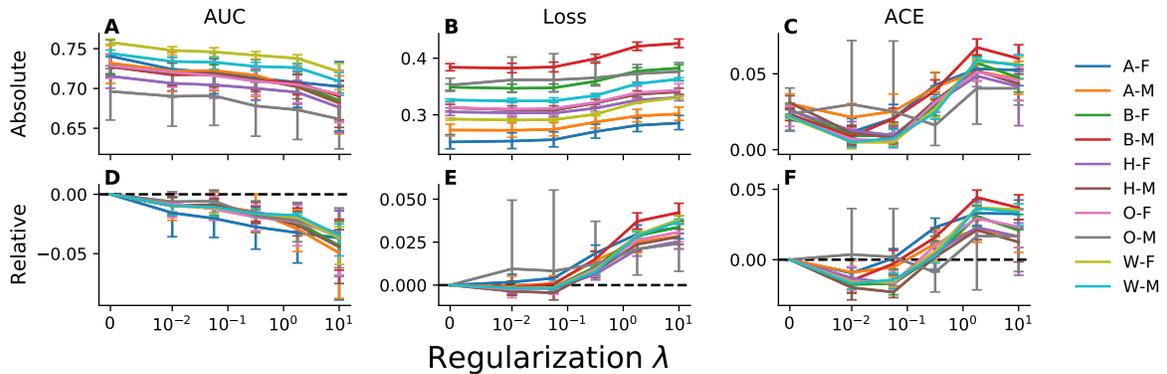
**Supplementary Figure C8:** The net benefit of models that estimate ten-year ASCVD risk for subgroups defined by the presence or absence of chronic kidney disease (CKD), rheumatoid arthritis (RA), or type 1 (T1) or type 2 (T2) diabetes. Results shown are the net benefit (NB) and calibrated net benefit (cNB), parameterized by the choice of a decision threshold of 7.5% or 20%, assessed in the overall population, on each subgroup, and in the worst-case over subgroups following the application of unconstrained pooled or stratified ERM, regularized objectives that penalize differences in the log-loss of AUC across subgroups, or DRO objectives that optimize for the worst-case log-loss of AUC across subgroups. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



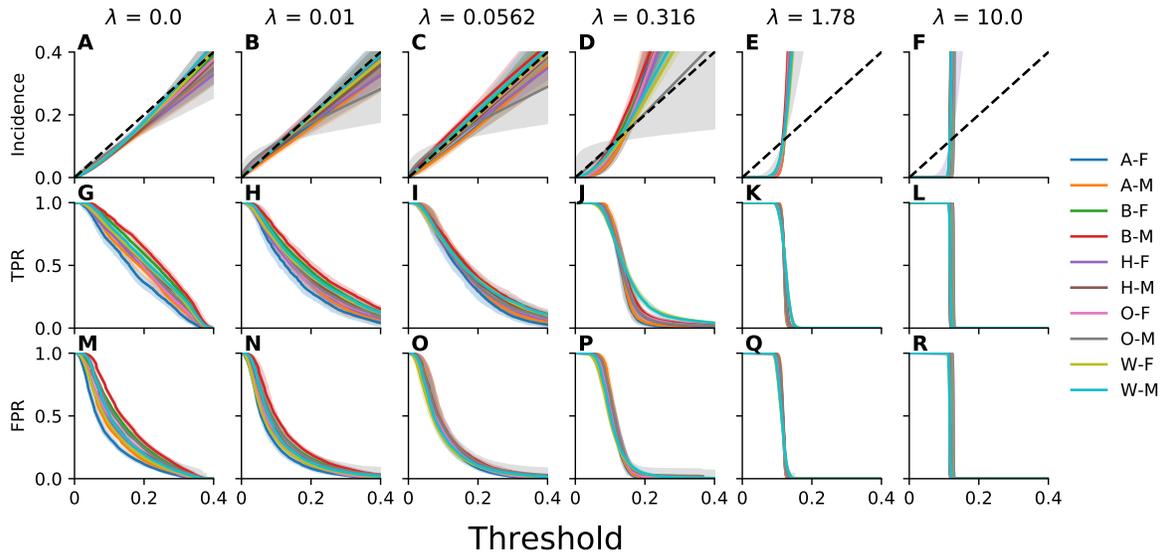
**Supplementary Figure C9:** The net benefit evaluated for a range of thresholds across racial and ethnic subgroups, parameterized by the choice of a decision threshold of 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



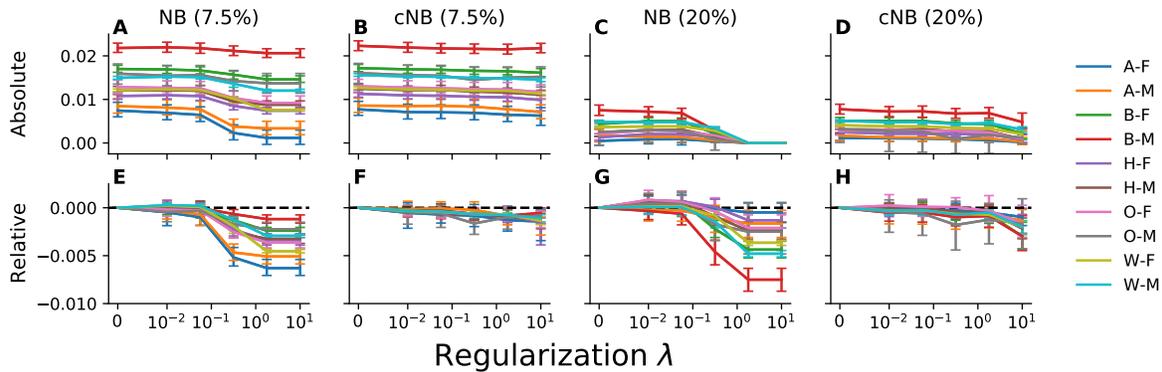
**Supplementary Figure C10:** Decision curve analysis to assess net benefit of models across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



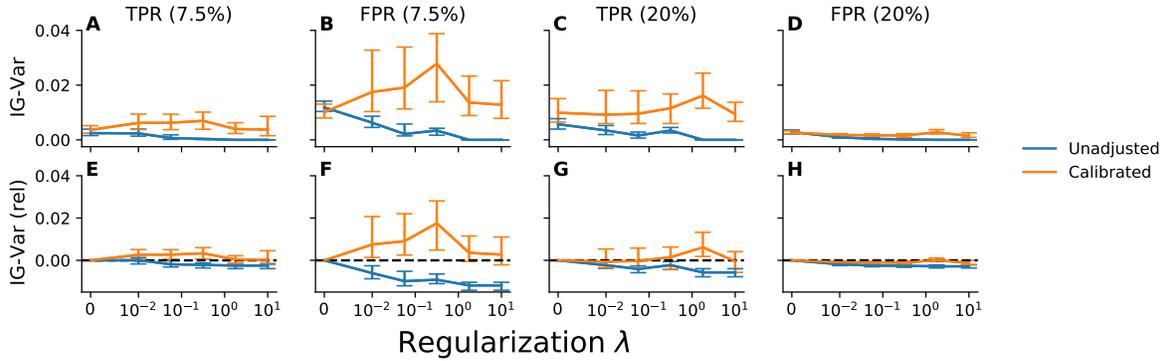
**Supplementary Figure C11:** The performance of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the area under the receiver operating characteristic curve (AUC), log-loss, and absolute calibration error (ACE). Relative results are reported relative to those attained for unconstrained empirical risk minimization. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



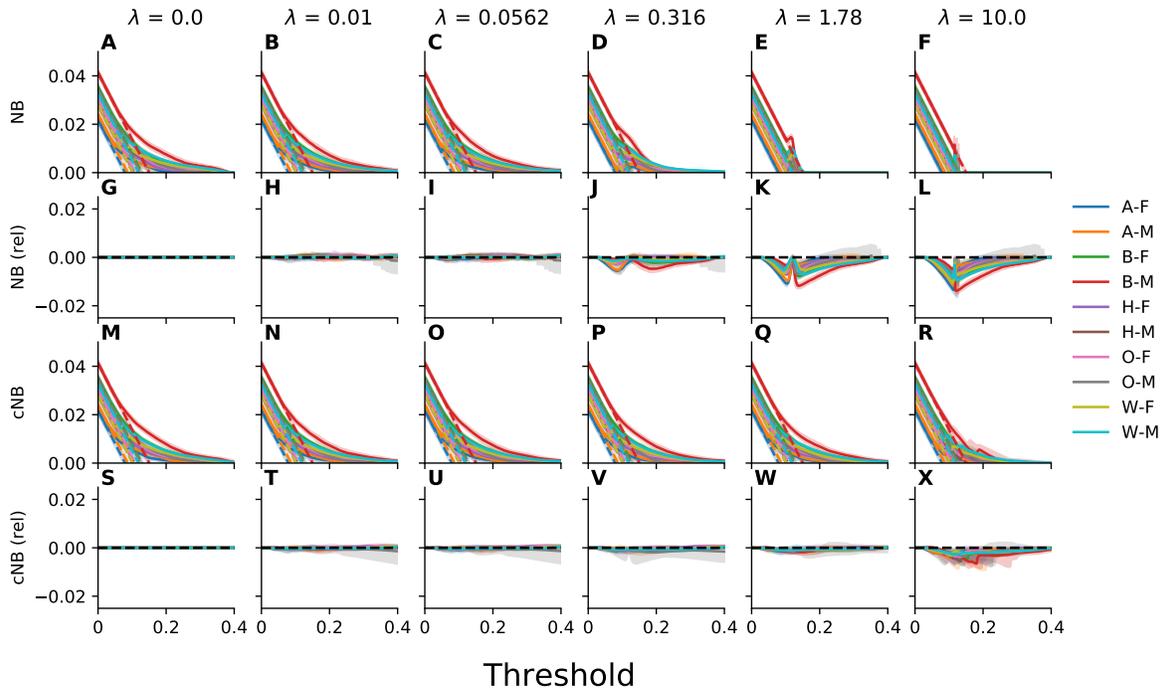
**Supplementary Figure C12:** Calibration curves, true positive rates, and false positive rates evaluated for a range of thresholds for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , are the calibration curve (incidence), true positive rate (TPR), and false positive rate (FPR) as a function of the decision threshold. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



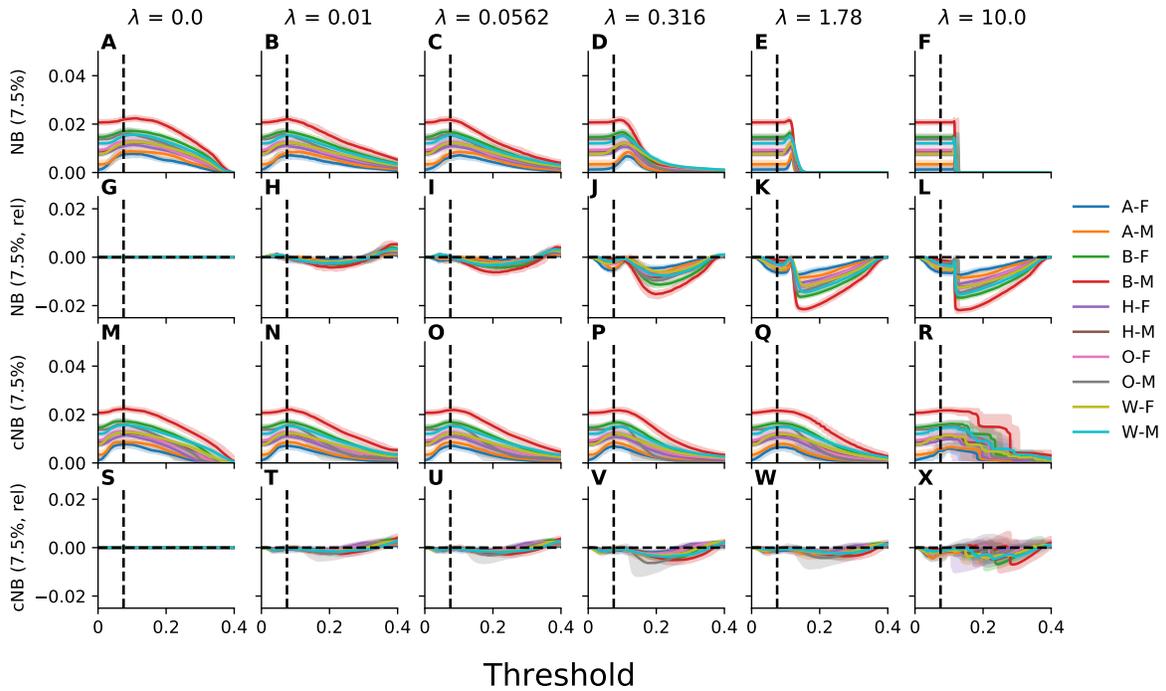
**Supplementary Figure C13:** The net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty, parameterized by the choice of a decision threshold of 7.5% or 20%. Plotted, for each subgroup is the net benefit (NB) and calibrated net benefit (rNB) as a function of the value of the regularization parameter  $\lambda$ . Relative results are reported relative to those attained for unconstrained empirical risk minimization. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



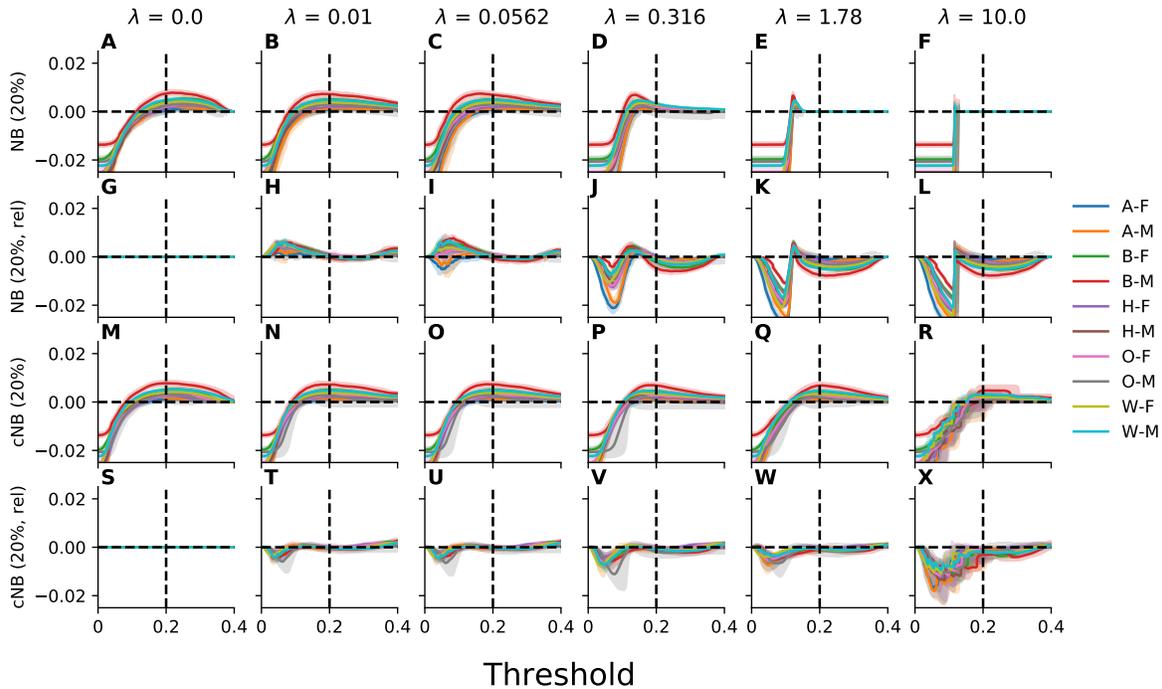
**Supplementary Figure C14:** Satisfaction of equalized odds for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted is the intergroup variance (IG-Var) in the true positive and false positive rates at decision thresholds of 7.5% and 20%. Recalibrated results correspond to those attained for models for which the threshold has been adjusted to account for the observed miscalibration. Relative results are reported relative to those attained for unconstrained empirical risk minimization. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



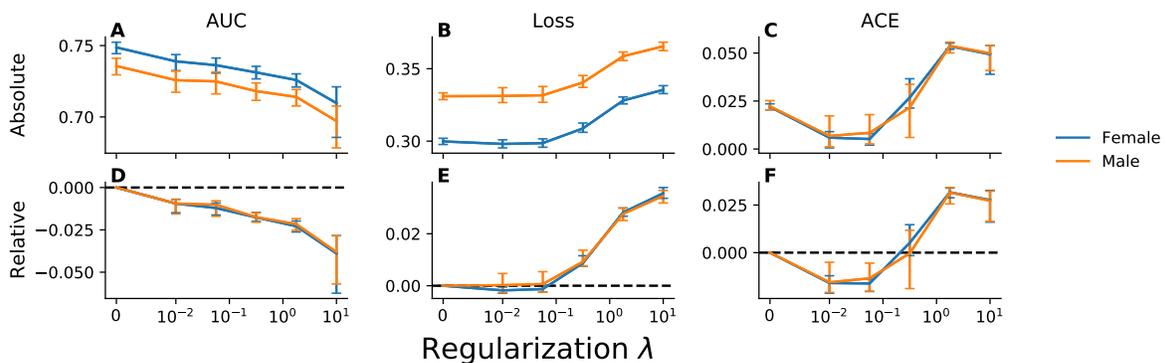
**Supplementary Figure C15:** Decision curve analysis to assess net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (rNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



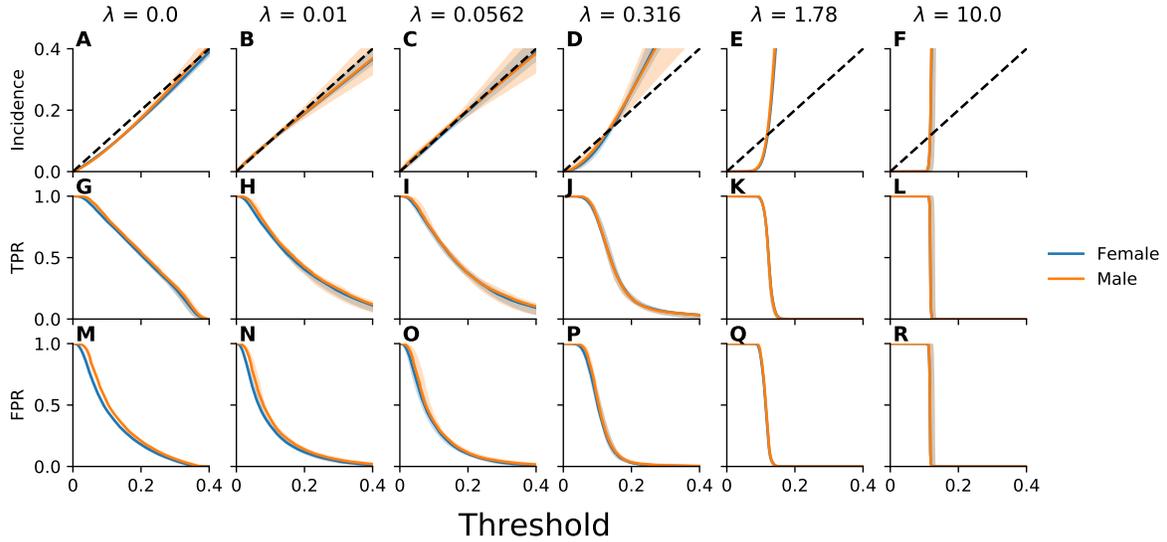
**Supplementary Figure C16:** The net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty, for the net benefit parameterized by the choice of a decision threshold of 7.5%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (rNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



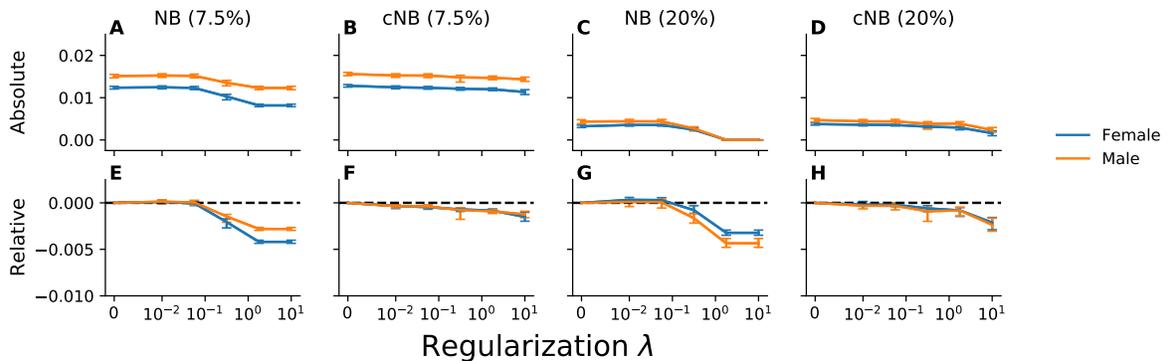
**Supplementary Figure C17:** The net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty, for the net benefit parameterized by the choice of a decision threshold of 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (rNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported to the results for unconstrained empirical risk minimization are indicated by “rel”. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



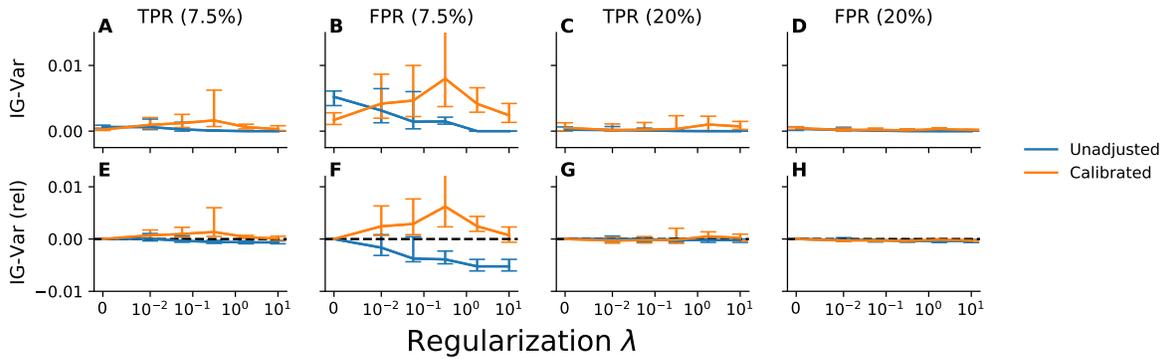
**Supplementary Figure C18:** Model performance evaluated across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the area under the receiver operating characteristic curve (AUC), log-loss, and absolute calibration error (ACE). Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



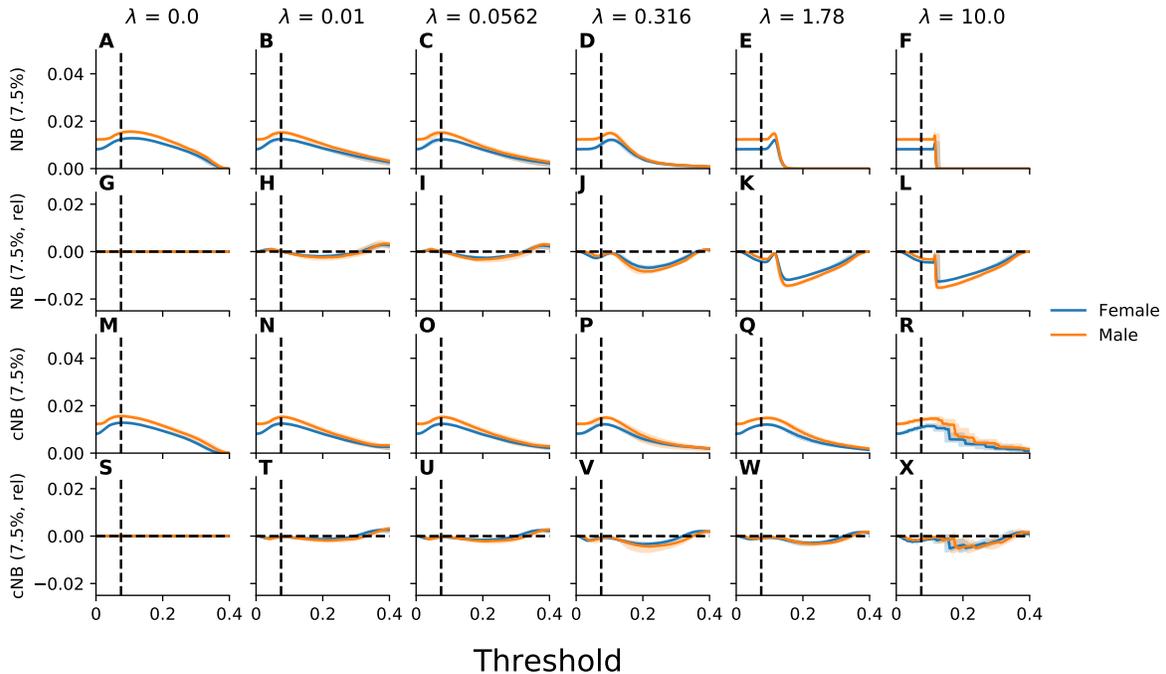
**Supplementary Figure C19:** Calibration curves, true positive rates, and false positive rates evaluated for a range of thresholds across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , are the calibration curve (incidence), true positive rate (TPR), and false positive rate (FPR) as a function of the decision threshold. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



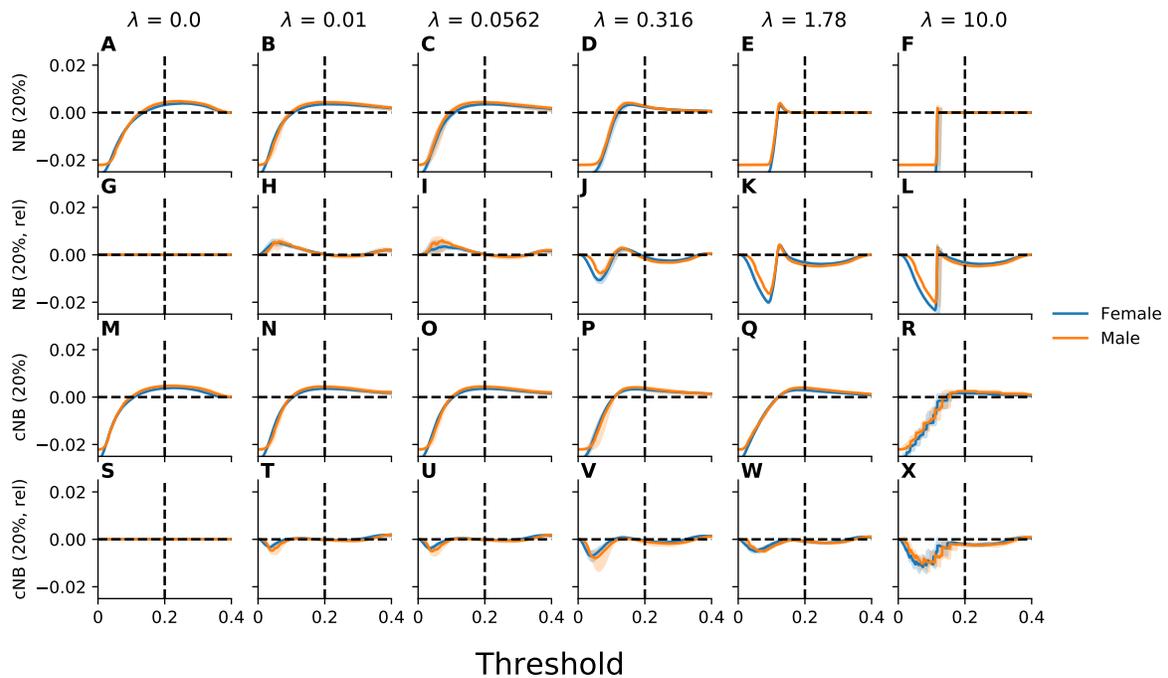
**Supplementary Figure C20:** The net benefit evaluated across subgroups defined by sex, parameterized by the choice of a decision threshold of 7.5% or 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup is the net benefit (NB) and calibrated net benefit (cNB) as a function of the value of the regularization parameter  $\lambda$ . Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



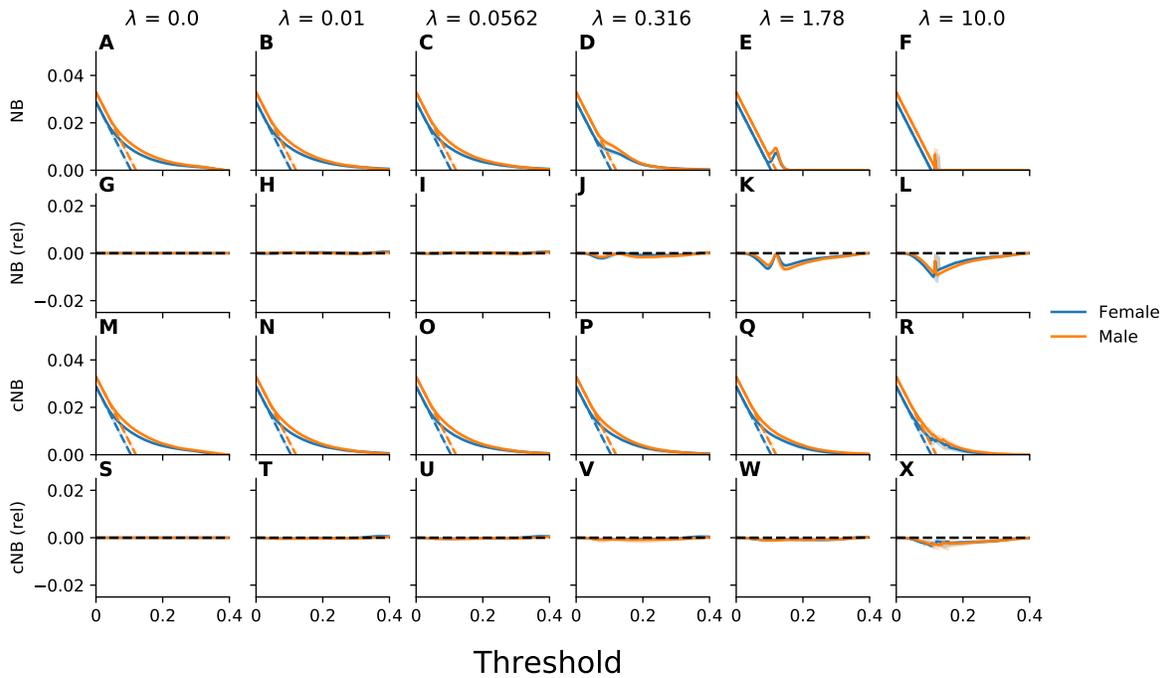
**Supplementary Figure C21:** Satisfaction of equalized odds evaluated across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted is the intergroup variance (IG-Var) in the true positive and false positive rates at decision thresholds of 7.5% and 20%. Recalibrated results correspond to those attained for models for which the threshold has been adjusted to account for the observed miscalibration. Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



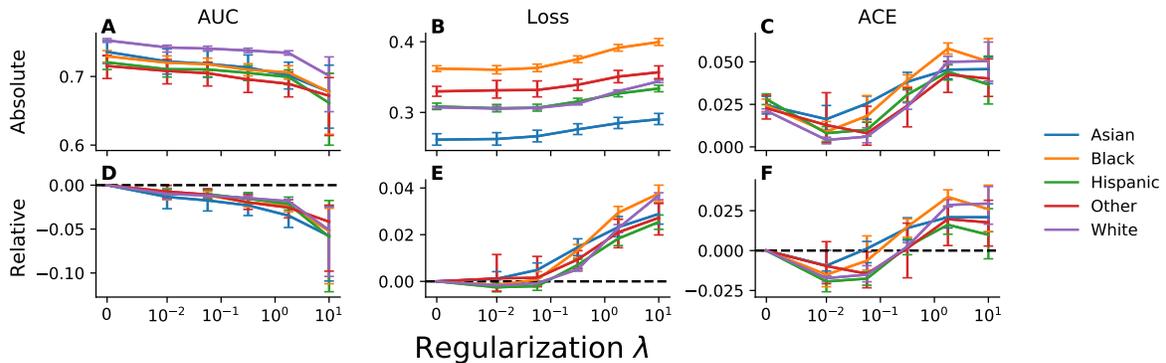
**Supplementary Figure C22:** The net benefit evaluated for a range of thresholds across subgroups defined by sex, parameterized by the choice of a decision threshold of 7.5%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



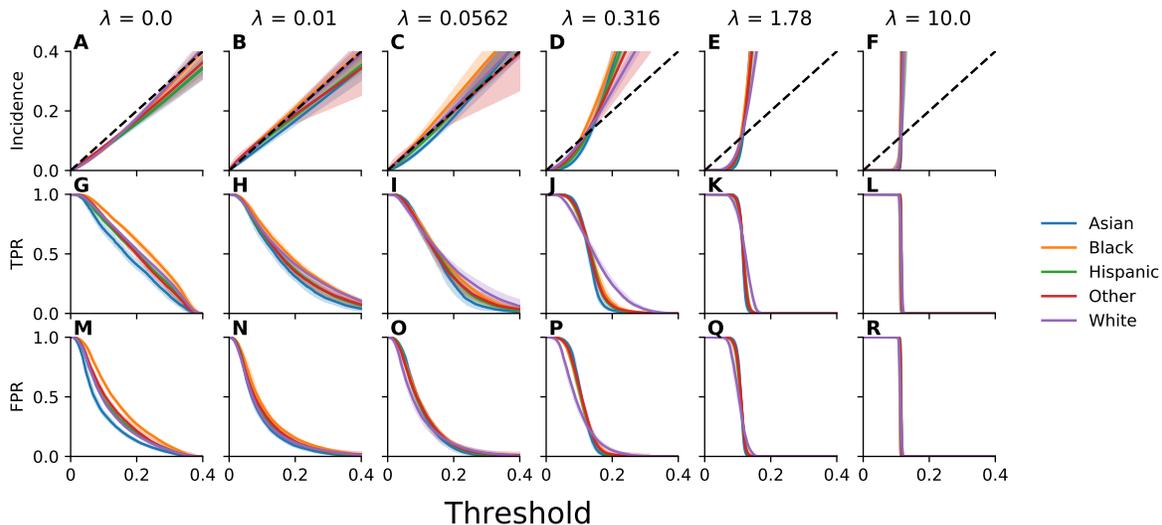
**Supplementary Figure C23:** The net benefit evaluated for a range of thresholds across subgroups defined by sex, parameterized by the choice of a decision threshold of 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



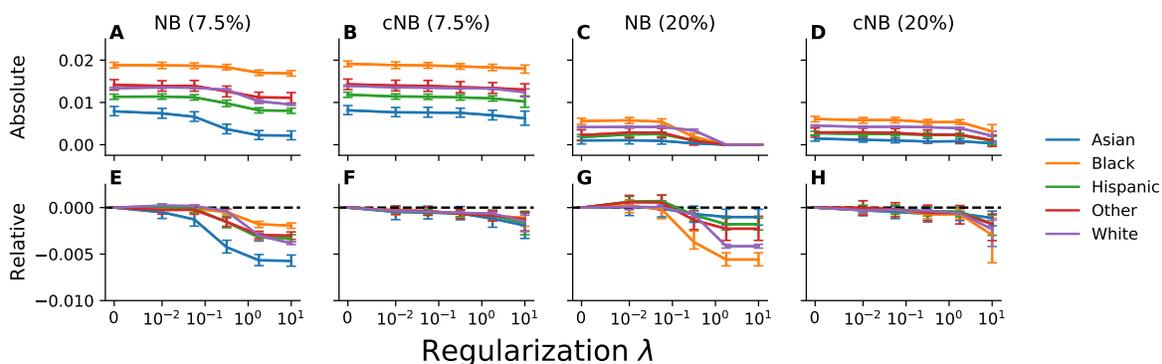
**Supplementary Figure C24:** Decision curve analysis to assess net benefit of models across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a MMD-based penalty. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



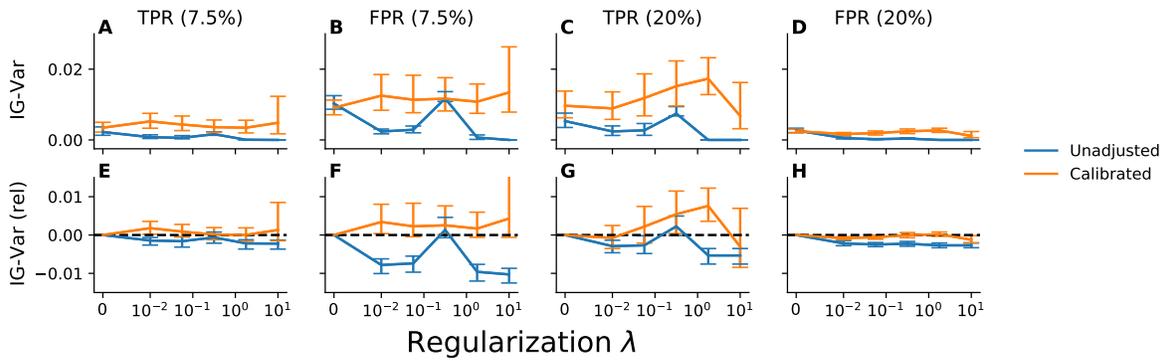
**Supplementary Figure C25:** Model performance evaluated across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the area under the receiver operating characteristic curve (AUC), log-loss, and absolute calibration error (ACE). Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



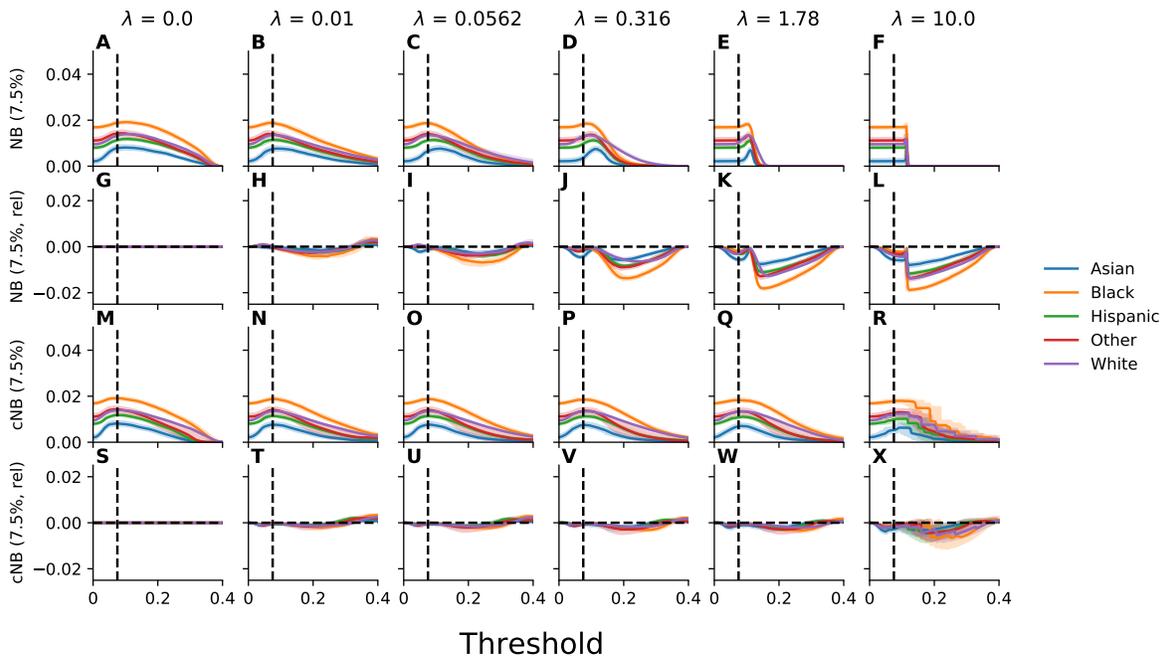
**Supplementary Figure C26:** Calibration curves, true positive rates, and false positive rates evaluated for a range of thresholds across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , are the calibration curve (incidence), true positive rate (TPR), and false positive rate (FPR) as a function of the decision threshold. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



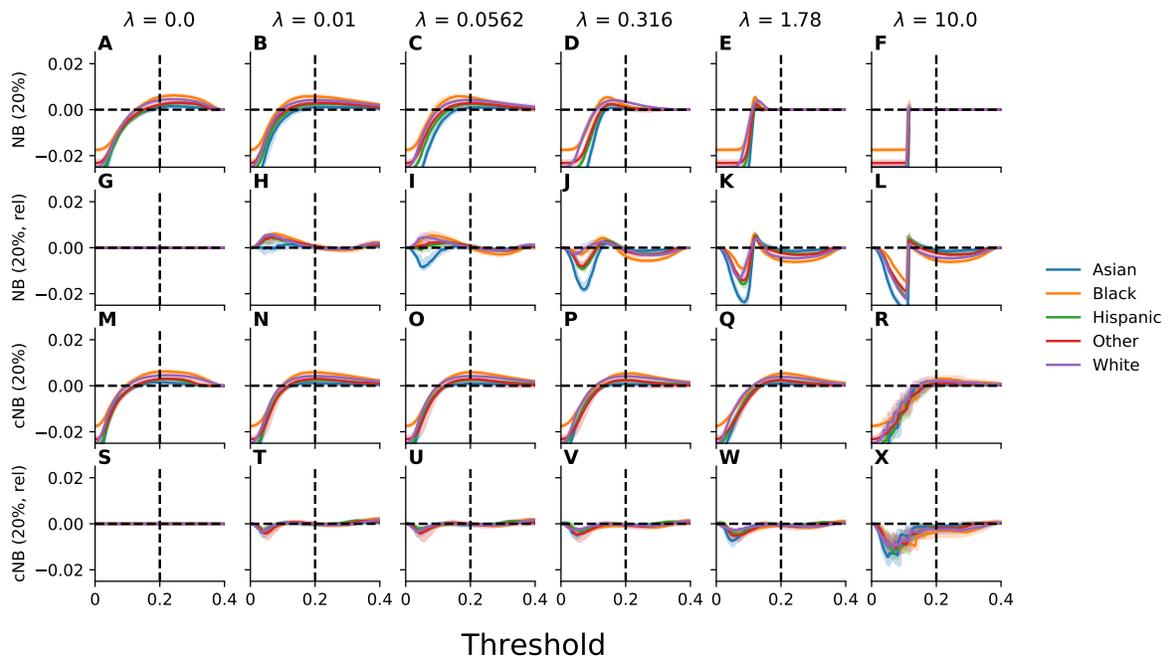
**Supplementary Figure C27:** The net benefit evaluated across racial and ethnic subgroups, parameterized by the choice of a decision threshold of 7.5% or 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup is the net benefit (NB) and calibrated net benefit (cNB) as a function of the value of the regularization parameter  $\lambda$ . Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



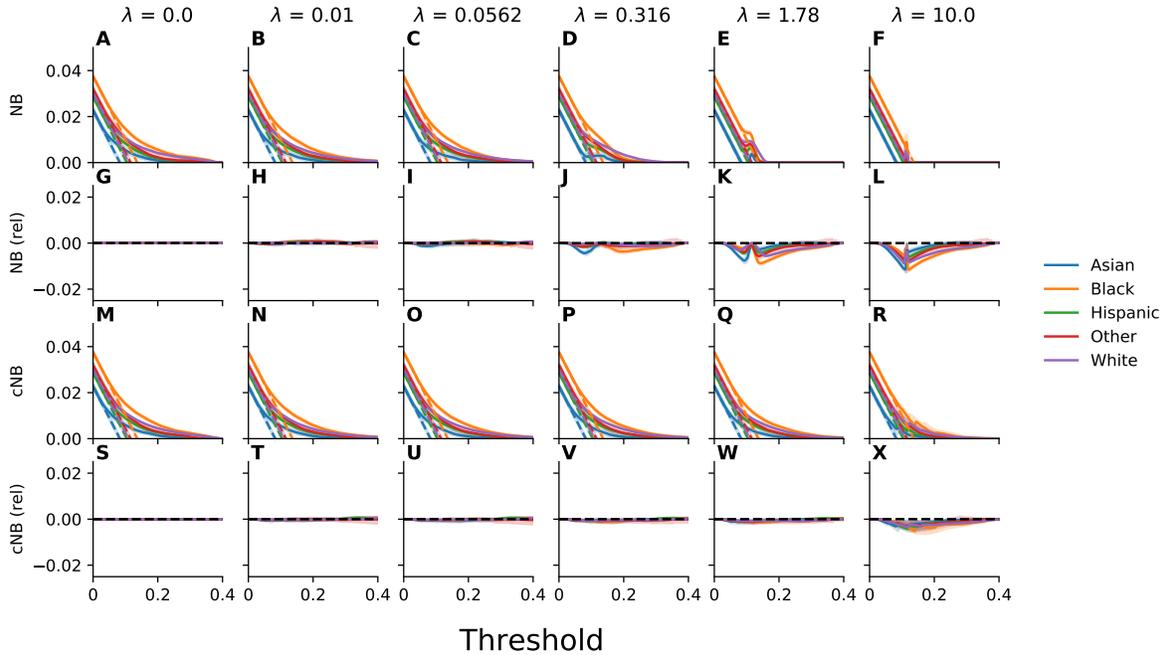
**Supplementary Figure C28:** Satisfaction of equalized odds evaluated across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted is the intergroup variance (IG-Var) in the true positive and false positive rates at decision thresholds of 7.5% and 20%. Recalibrated results correspond to those attained for models for which the threshold has been adjusted to account for the observed miscalibration. Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



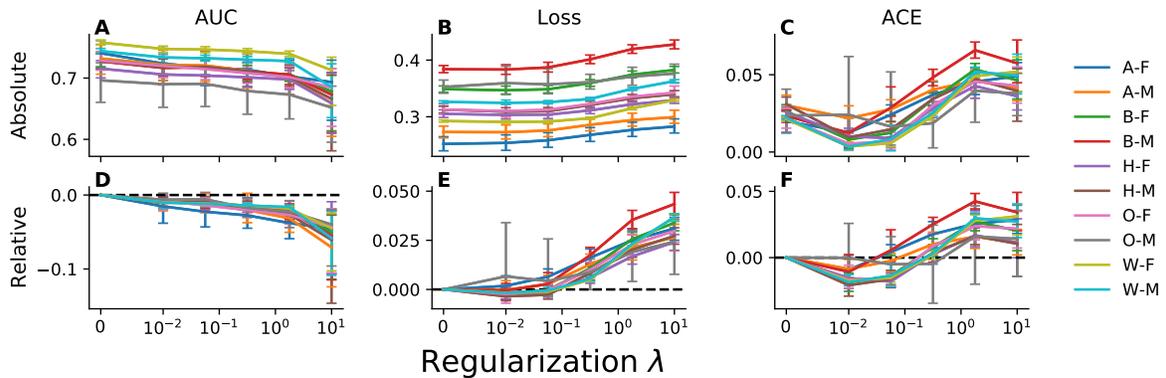
**Supplementary Figure C29:** The net benefit evaluated for a range of thresholds across racial and ethnic subgroups, parameterized by the choice of a decision threshold of 7.5%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



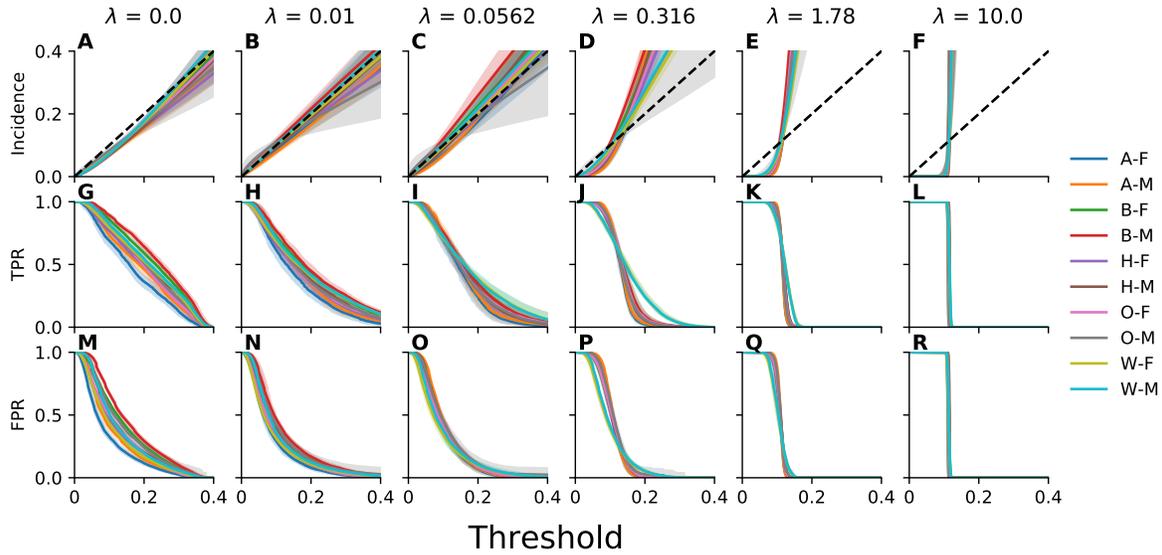
**Supplementary Figure C30:** The net benefit evaluated for a range of thresholds across racial and ethnic subgroups, parameterized by the choice of a decision threshold of 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



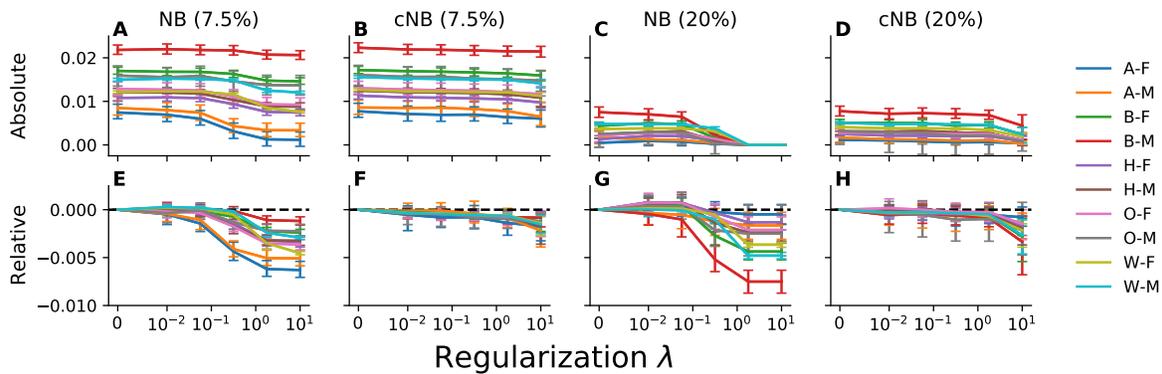
**Supplementary Figure C31:** Decision curve analysis to assess net benefit of models across racial and ethnic subgroups for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



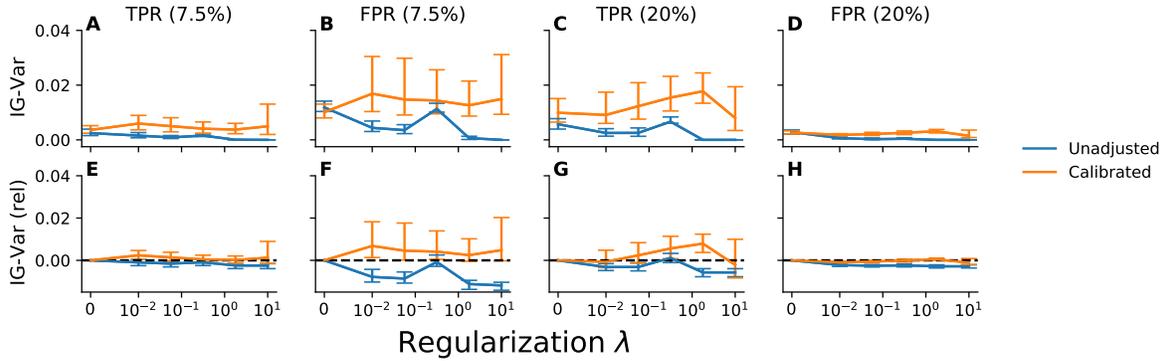
**Supplementary Figure C32:** The performance of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the area under the receiver operating characteristic curve (AUC), log-loss, and absolute calibration error (ACE). Relative results are reported relative to those attained for unconstrained empirical risk minimization. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



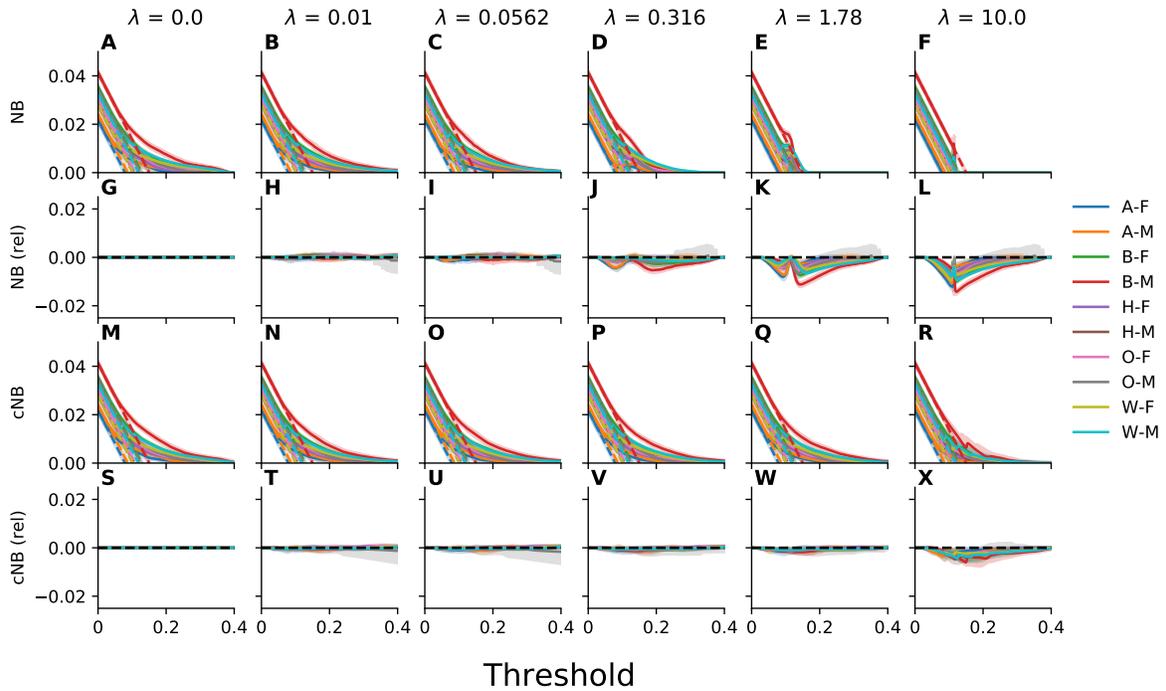
**Supplementary Figure C33:** Calibration curves, true positive rates, and false positive rates evaluated for a range of thresholds for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , are the calibration curve (incidence), true positive rate (TPR), and false positive rate (FPR) as a function of the decision threshold. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



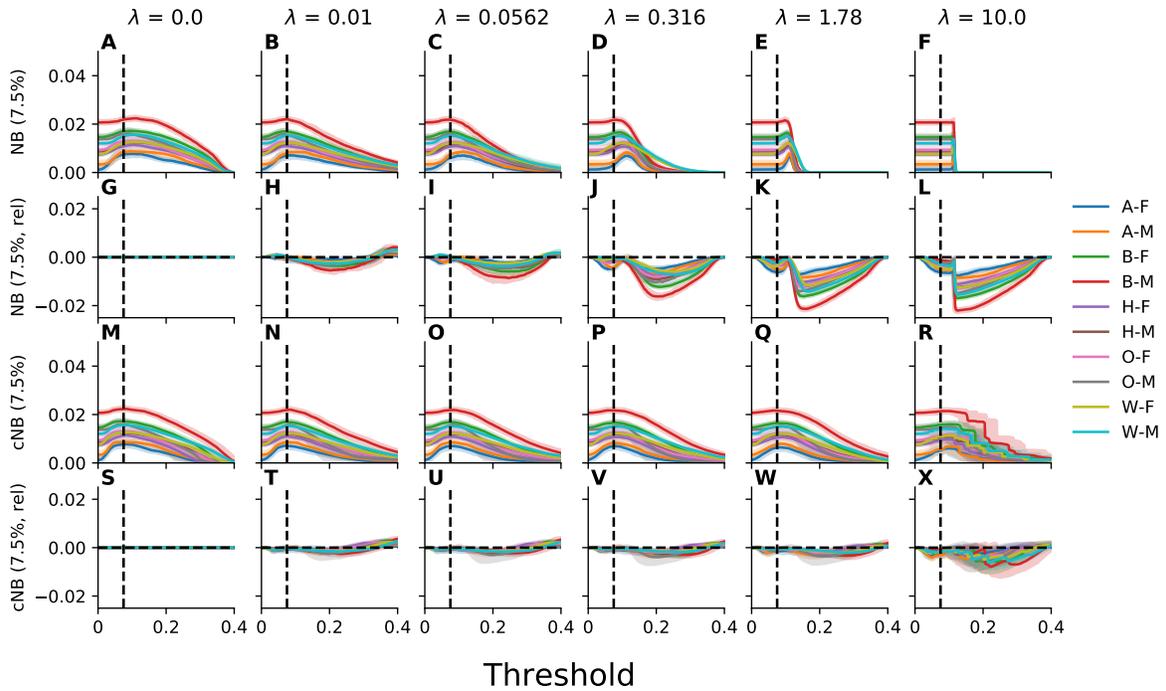
**Supplementary Figure C34:** The net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%, parameterized by the choice of a decision threshold of 7.5% or 20%. Plotted, for each subgroup is the net benefit (NB) and calibrated net benefit (rNB) as a function of the value of the regularization parameter  $\lambda$ . Relative results are reported relative to those attained for unconstrained empirical risk minimization. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



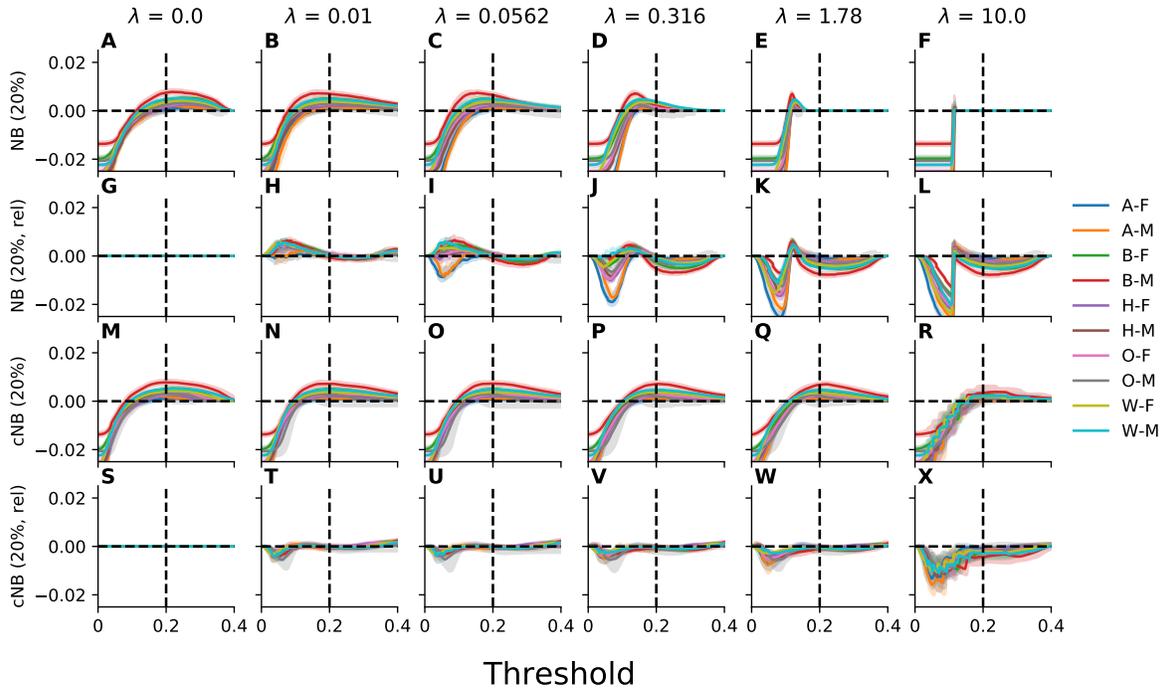
**Supplementary Figure C35:** Satisfaction of equalized odds for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted is the intergroup variance (IG-Var) in the true positive and false positive rates at decision thresholds of 7.5% and 20%. Recalibrated results correspond to those attained for models for which the threshold has been adjusted to account for the observed miscalibration. Relative results are reported relative to those attained for unconstrained empirical risk minimization. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



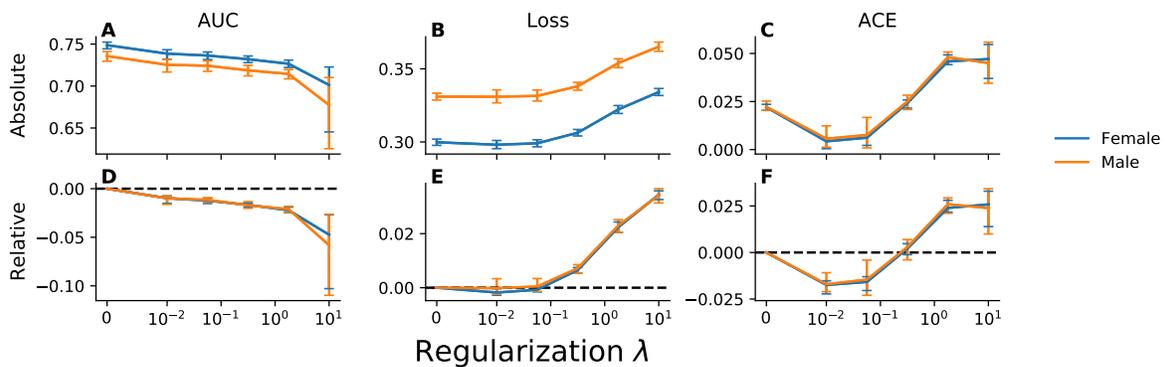
**Supplementary Figure C36:** Decision curve analysis to assess net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (rNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



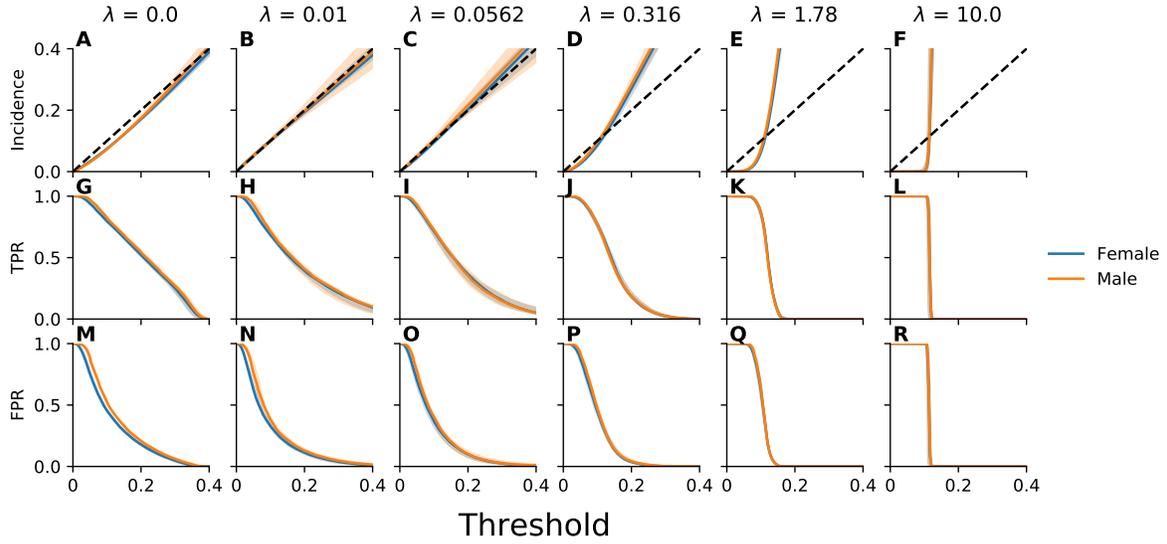
**Supplementary Figure C37:** The net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%, parameterized by the choice of a decision threshold of 7.5%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (rNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



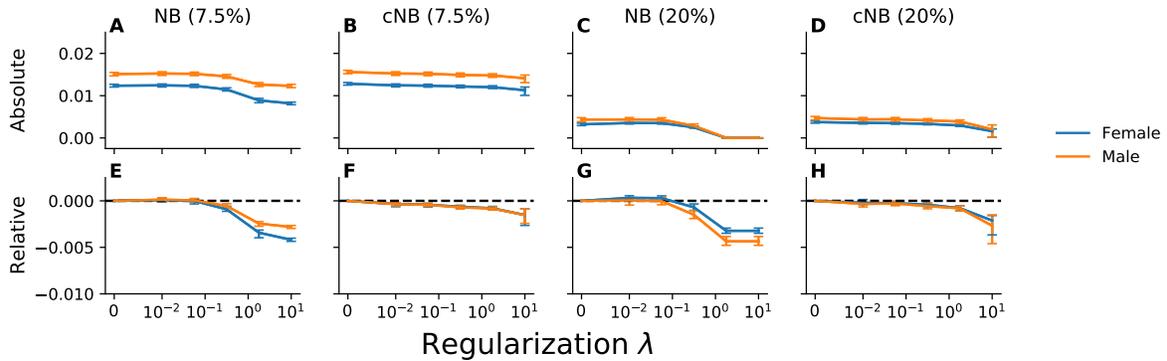
**Supplementary Figure C38:** The net benefit of models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%, parameterized by the choice of a decision threshold of 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (rNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Labels correspond to Asian (A), Black (B), Hispanic (H), Other (O), White (W), Male (M), and Female (F) patients. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



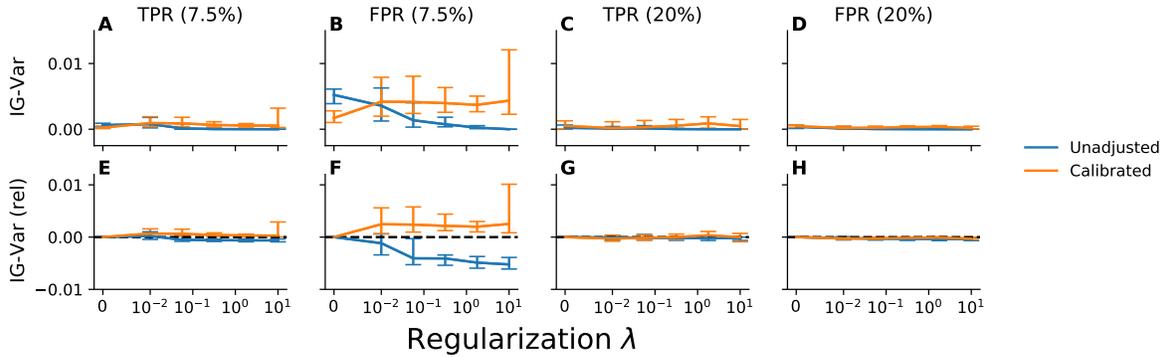
**Supplementary Figure C39:** Model performance evaluated across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the area under the receiver operating characteristic curve (AUC), log-loss, and absolute calibration error (ACE). Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



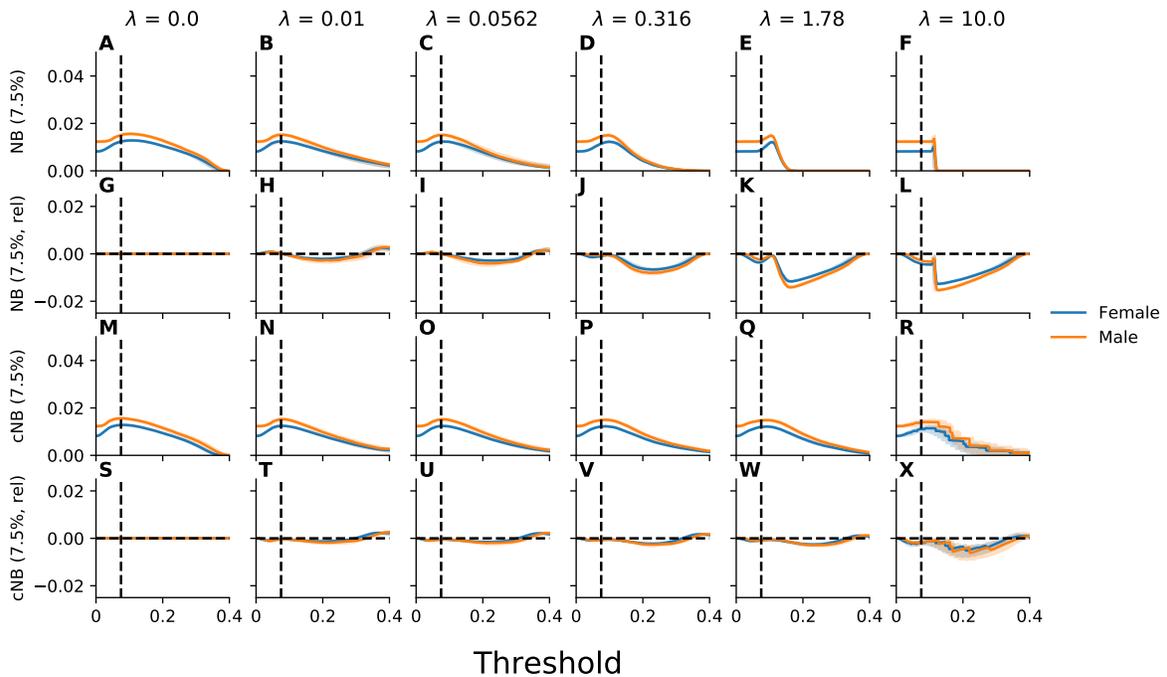
**Supplementary Figure C40:** Calibration curves, true positive rates, and false positive rates evaluated for a range of thresholds across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , are the calibration curve (incidence), true positive rate (TPR), and false positive rate (FPR) as a function of the decision threshold. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



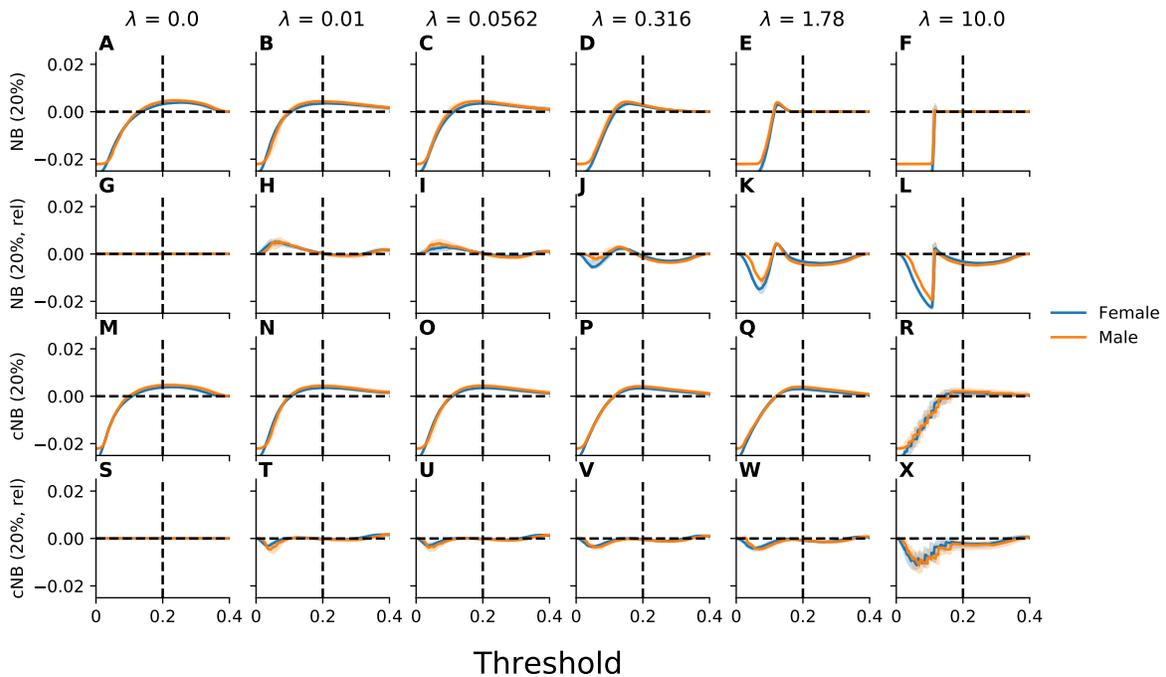
**Supplementary Figure C41:** The net benefit evaluated across subgroups defined by sex, parameterized by the choice of a decision threshold of 7.5% or 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup is the net benefit (NB) and calibrated net benefit (cNB) as a function of the value of the regularization parameter  $\lambda$ . Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



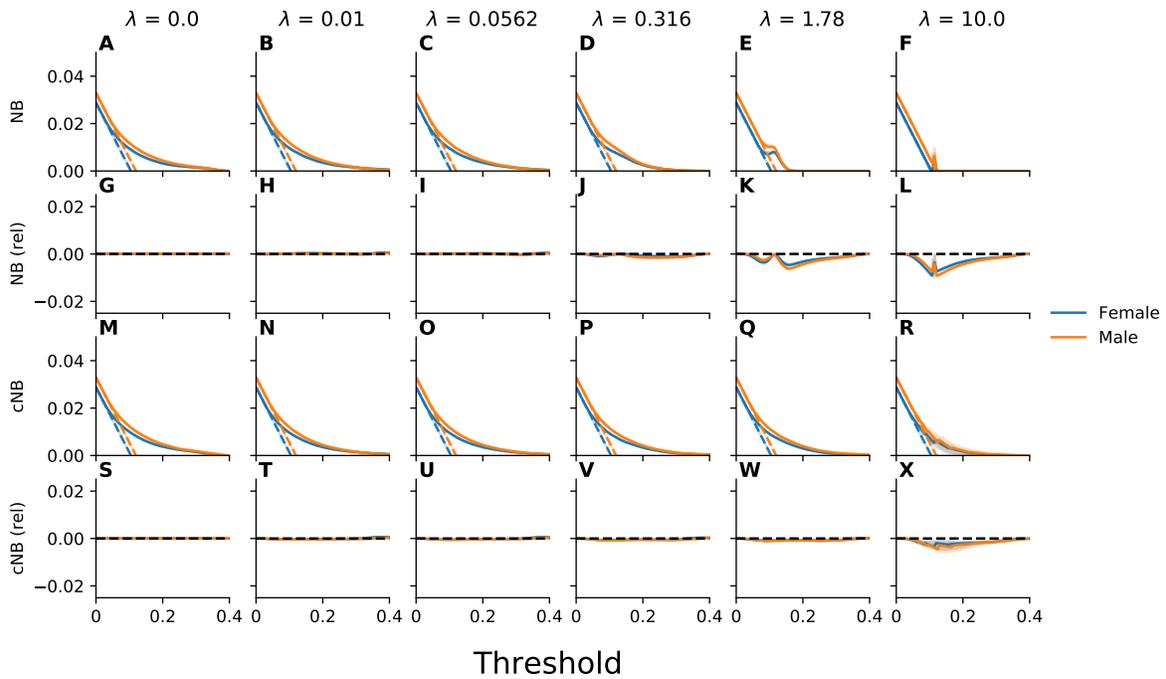
**Supplementary Figure C42:** Satisfaction of equalized odds evaluated across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted is the intergroup variance (IG-Var) in the true positive and false positive rates at decision thresholds of 7.5% and 20%. Recalibrated results correspond to those attained for models for which the threshold has been adjusted to account for the observed miscalibration. Relative results are reported relative to those attained for unconstrained empirical risk minimization. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



**Supplementary Figure C43:** The net benefit evaluated for a range of thresholds across subgroups defined by sex, parameterized by the choice of a decision threshold of 7.5%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



**Supplementary Figure C44:** The net benefit evaluated for a range of thresholds across subgroups defined by sex, parameterized by the choice of a decision threshold of 20%, for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.



**Supplementary Figure C45:** Decision curve analysis to assess net benefit of models across subgroups defined by sex for models trained with an objective that penalizes violation of equalized odds across intersectional subgroups defined on the basis of race, ethnicity, and sex using a threshold-based penalty at 7.5% and 20%. Plotted, for each subgroup and value of the regularization parameter  $\lambda$ , is the net benefit (NB) and calibrated net benefit (cNB) as a function of the decision threshold. The net benefit of treating all patients is designated by dashed lines. Results reported relative to the results for unconstrained empirical risk minimization are indicated by “rel”. Error bars indicate 95% confidence intervals derived with the percentile bootstrap with 1,000 iterations.