

Empirical observation of negligible fairness-accuracy trade-offs in machine learning for public policy

Kit T. Rodolfa,¹ Hemank Lamba,¹ Rayid Ghani^{1*}

¹Machine Learning Department and Heinz College of Information Systems and Public Policy, Carnegie Mellon University, Pittsburgh, PA 15213, USA

*To whom correspondence should be addressed; E-mail: rayid@cmu.edu.

Abstract

Growing use of machine learning in policy and social impact settings have raised concerns for fairness implications, especially for racial minorities. These concerns have generated considerable interest among machine learning and artificial intelligence researchers, who have developed new methods and established theoretical bounds for improving fairness, focusing on the source data, regularization and model training, or post-hoc adjustments to model scores. However, little work has studied the practical trade-offs between fairness and accuracy in real-world settings to understand how these bounds and methods translate into policy choices and impact on society. Our empirical study fills this gap by investigating the impact of mitigating disparities on accuracy, focusing on the common context of using machine learning to inform benefit allocation in resource-constrained programs across education, mental health, criminal justice, and housing safety. Here we describe applied work in which we find fairness-accuracy trade-offs to be negligible in practice. In each setting studied, explicitly focusing on achieving equity and using our proposed post-hoc disparity mitigation methods, fairness was substantially improved without sacrificing accuracy. This observation was robust across policy contexts studied, scale of resources available for intervention, time, and relative size of the protected groups. These empirical results challenge a commonly held assumption that reducing disparities either requires accepting an appreciable drop in accuracy or the development of novel, complex methods, making reducing disparities in these applications more practical.

There has been a rapid growth in the use of machine learning for applications with extensive impact on society, such as informing bail determination decisions,¹⁻³ hiring,⁴ healthcare delivery,^{5,6} and social service interventions.⁷⁻⁹ These wide-reaching applications have been met

with heightened concerns about their potential for introducing or amplifying inequities, especially for racial minorities and economically disadvantaged individuals, motivating exploration of a range of potential sources and mitigation strategies for biases, including in the underlying data,¹⁰ labels,⁵ model training,¹¹⁻¹³ and post-modeling adjustments to scores.^{14,15} A common underpinning of much of this work is the assumption that trade-offs between equity and accuracy may necessitate complex methods or difficult policy choices,¹⁶⁻¹⁹ however little work to date has explicitly evaluated the magnitude (or even the existence) of these trade-offs in real-world problems. Note that in this work we use the term “accuracy” in the more colloquial sense of the correctness of a model’s predictions relative to the task at hand (in contrast to the fairness of those predictions), rather than the specific statistical property of the same name. For each policy setting we study we use a more specific “accuracy” metric based on the goal of the program.

Our study focuses on testing the assumed accuracy-fairness trade-offs in resource allocation problems across several public policy domains. Organizations with limited resources are often only able to intervene and allocate benefits to a relatively small number of individuals with need, presenting a “top k ” optimization problem where model accuracy is judged by *precision* (also known as positive predictive value) among the k highest-scoring individuals (although precision in the top k readily maps to a concept of efficiently allocating limited resources in these settings, it is worth noting that for a given value of k on a fixed dataset, knowing the value of precision in the top k will fully determine both the values of recall in the top k and accuracy in the top k such that optimizing for any one of these metrics is equivalent to optimizing for the other two as well). In such assistive intervention settings, we¹⁵ and others¹⁴ have argued that recall (also known as sensitivity or true positive rate) disparities are often an appropriate equity metric, reflecting a concept of “equality of opportunity.” In a recent case study,¹⁵ we found that **explicitly focusing on achieving equity and using subgroup-specific score thresholds as a post-hoc disparity mitigation method improved the equity of predictions with only a very modest decrease in accuracy.** While that study focused on our experiences incorporating fairness into the deployment of a machine learning system and the related policy decision-making in a single context (developing social service interventions as a means of jail diversion in Los Angeles, CA), the empirical work here extends that study’s surprising result to several new policy contexts and modeling choices. Our results suggest that trade-offs between fairness and effectiveness can in fact be negligible in practice, suggesting that improvement in may be equity easier and more practical across a wide range of applications than is often expected. We come to this conclusion using a variety of projects we undertook over the past few years with government agencies across criminal justice, mental health, housing safety, and education, finding that each of these contexts poses a counterexample to the assumption of large trade-offs between fairness and accuracy.

Policy Settings and Data

Information about the policy settings included in the present study are provided in Table 1, and a brief description of the context and problem for each is provided below. Note that the machine learning formulation (including outcome measures, train-validation splits, evaluation metrics) and details in each case were determined through careful scoping and collaboration with the partner organization. For instance, the top k list size in each case reflects the organization’s resource constraints for intervening on entities, while policymaker and stakeholder priorities informed the label definition as well as the sensitive attribute used for measuring fairness.

Inmate Mental Health Seeking to break the cycle of incarceration for individuals with untreated mental health conditions, Johnson County, KS, partnered with us to prioritize limited resources for mental health outreach on individuals at risk of a future jail booking. We developed a predictive model of risk for a booking in the next year, focusing on identifying 500 individuals for outreach in a 4-month window based on the resources available to the program. Disparities on race and ethnicity are particularly salient in the criminal justice context, and we focused on this attribute in our bias analyses (here, we look either at 2-way results between white and non-white individuals, or 3-way results between white, Black, and Hispanic individuals; individuals in other or unknown racial or ethnic groups make up only 0.4% of the dataset and were included with white individuals for these analyses).

Housing Safety The Code Enforcement Office in San Jose, CA, is tasked with protecting occupants of properties with multiple units (such as apartment buildings) by conducting safety inspections, but doesn’t have sufficient staffing to inspect all 4,500 properties every year. Using internal data supplied by the program, we developed a model for the risk that a serious violation would be found if a given property were prioritized for inspection. We focused on disparities between housing units in higher- and lower-income neighborhoods (median income above or below \$55,000) where considerable disparities were observed in our initial models favoring higher-income areas.

Student Outcomes El Salvador’s Ministry of Education seeks to support students to reduce the country’s substantial dropout rates (recently as high as 29% in some years), but the budget for these programs is insufficient to reach every student. Student-level data was provided by the Ministry to develop a model of students at risk of dropping out, and our analysis here focuses on identifying the 10,000 highest-risk students in the state of San Salvador. The Ministry of Education was concerned with potential disparities in gender, age relative to grade level, and urban-rural divide. Initial analyses found large disparities with “over-age” students (those at least 2 standard deviations above the mean of their grade level), which we focus on in the present study.

Since the projects above used confidential and sensitive data and were done under data use agreements, we are not able to make that data publicly available. For our work to be easily reproducible, we include a fourth problem in this study where the data is available publicly:

Education Crowdfunding The non-profit DonorsChoose helps alleviate school funding shortages by providing a crowdfunding platform for teachers to post requests for their classroom

needs. Here, we make use of a dataset DonorsChoose made publicly available in 2014 and posit an effort to assist projects at risk of going unfunded (for instance, providing a review and consultation) capable of helping 1,000 projects in a 2-month window. Reflecting the platform’s goal of helping schools and teachers most at need, we focus in this context on disparities across school poverty levels (65% free/reduced lunch vs others). Although unlike the other settings described above, the Education Crowdfunding analysis did not arise from a project undertaken in partnership with the organization, we sought to scope and formulate a project with the same characteristics of those we typically encounter with partners.

Details of the machine learning modeling and evaluation performed in each context are described in the Methods. In short, we explored an expansive grid of machine learning model types and hyperparameters in each problem to generate candidate models for our analysis. Because of the inherently non-stationary nature of real-world problems in general and policy applications in particular (with policies, practices, and context changing over time), training and validation sets were generated through a process of inter-temporal cross-validation. This approach creates a set of temporally sequential train and validation datasets, reflecting deployment of machine learning in these settings where models trained on past data must generalize into a changing future context. Although not the primary focus of the current work, we note that in each context, modeling was able to provide meaningful improvements over the base rate for the outcome (label) of interest in the population (shown in Table 1), which could translate into substantive efficiency gains in policy implementation. Note that in a baseline model which randomly chose a group for intervention, the expected fraction of true positives in this group (that is, the precision) would be given by this base rate of the label. In each of the Housing Safety, Student Outcomes, and Education Crowdfunding problems, model performance provided a roughly 2-fold increase over this baseline, and in the Inmate Mental Health context (where the underlying event is more rare, with a base rate of 12%), this improvement was by more than a factor of 4, making the ML models useful compared to simpler solutions and baselines.

The following results focus on using group-specific score thresholds to mitigate disparities observed in these models’ predictions. Fig. 1a illustrates the intuition behind making group-specific adjustments to a uniform score threshold in order to equalize recall (also known as sensitivity or true positive rate) across groups. Because recall increases monotonically with lower score thresholds, a unique solution can be found that equalizes recall across groups while keeping a fixed total number of entities selected for intervention. In Fig. 1b, we provide a schematic of the temporal strategy used for making these adjustments and testing their ability to generalize to unseen future data. In this figure, the dark-colored rectangles indicate points in time at which a cohort of training or validation examples is defined (this might be school years in the Student Outcomes context or quarters in the Housing Safety context) and the associated light-colored rectangles reflect a buffer during which outcomes are measured (for instance, the 4 months a project in the Education Crowdfunding context has to collect donations) in order to avoid leakage²⁰ between the training and validation sets. Note that features for each cohort will make use of all earlier information as well. For one point in our analysis, a grid of models is

trained using the blue cohort (t_{-2}), with performance on the purple cohort (t_{-1}) used to select well-performing models and find thresholds which equalize recall across groups. In order to evaluate how well these thresholds would then improve fairness (as well as any associated accuracy trade-off) when applied to new data, we step the process forward in time: training models using data up to the purple cohort and applying them along with the group-specific thresholds on the orange cohort (t_0), which reflects our true generalization performance. To understand consistency and stability of our results over time, we repeat this process several times, shifting the three cohorts together across different temporal validation splits.

Measuring Fairness

One challenge encountered by both researchers and practitioners seeking to improve the fairness of machine learning models is the nebulous nature of the concept of “fairness” itself. Much has been written about the wide range of metrics which can be used to measure or conceptualize fairness in different contexts,^{21,22} and, in particular, the mathematical incompatibility of various sets of these metrics.^{1,23} As a result, an important aspect of project scoping is understanding the fairness metric (or metrics) most appropriate to the given context. We have previously described a framework for this process.¹⁵ For instance, machine learning applications which inform benefit allocation decisions might be primarily concerned with avoiding disparities in false negatives (that is, failing to reach individuals with need), while applications that are more punitive in nature (such as bail denial decisions) are likely to be more concerned with disparities in false positives.

In the present work, we focus on resource-constrained assistive programs, a context in which we argue disparities in recall (also known as sensitivity or true positive rate) is a natural conceptualization of fairness. Because recall is the complement of the false negative rate (that is, $recall = TPR = 1 - FNR$), improving equity in terms of recall reflects a focus on errors of omission, but provides more readily-interpretable ratios than the false negative rate itself when resources are limited (such that even a highly predictive model would have a false negative rate near 1). Hardt¹⁴ provides some additional helpful intuition here, describing this metric as a reflection of “equality of opportunity.” Specifically, when only a fraction of individuals with need can receive a benefit, recall measures the fraction of those individuals a program reaches and disparities in recall therefore measure whether individuals with need across different sub-groups have an equal chance of receiving the benefit. In particular, we measure recall disparity as the ratio between the recall attained by a model for one group of interest (such as Black individuals) and the model’s recall for a reference group (such as white individuals). In contexts with more than two sub-groups of interest, we choose a reference group and consider disparities for all other groups relative to this group separately.

Evaluating Trade-Offs

We started with the assumption, based on existing theoretical work, that there is a trade-off between “accuracy” and “fairness.” As an initial experiment, we explored how model “accuracy” changes upon adjusting for disparities in the Inmate Mental Health setting using a single temporal validation split (with validation set outcomes spanning 4/2018 to 4/2019). In Figs. 1c and 1d, each pair of points is a model specification: with results obtained without adjusting for equity in blue and those with the equity adjustment in orange. The x-axis shows the precision (positive predictive value) of each model on the 500 selected individuals and the y-axis shows the recall disparity between white and non-white individuals. Fig. 1c shows all models considered (the model specifications used here are detailed in Supplementary Table 1), while Fig. 1d provides a more detailed view of the better-performing models that might reasonably be selected. All of the unadjusted models had significant disparities, indicating that merely measuring disparities to inform the process of model selection would be insufficient for achieving fairness. However, applying our proposed bias mitigation method to adjust for disparities, we find little evidence of a fairness-accuracy trade-off: overall, the mean change in precision after adjustment is -0.0006 (std: 0.0087); Fig. 1e shows the distribution of these shifts.

We also investigated creating a composite model (based on the work of Dwork and colleagues¹³) by choosing the best-performing model for each subgroup and using the recall-equalizing set of individuals (maintaining a total list size of 500) from each subgroup-optimized model. The performance of this composite on the subsequent validation set is shown as a red diamond in Figs. 1c and 1d. On both fairness and accuracy-related metrics, this composite appears to perform competitively with other models, but does not stand out as out-performing the fairness-adjusted individual models on either metric. The promise of the composite strategy is that it might improve the accuracy on each subgroup (and, hence, the overall performance across subgroups) while yielding more equitable results. In this initial experiment, however, the more complex approach of building a composite model shows no indication of providing gains beyond the simpler approach of making fairness-enhancing adjustments to a single, well-performing model.

Comparing Mitigation Strategies

These initial results suggested disparity mitigation could be integrated into the process of model selection, and we next sought to compare strategies for doing so (summarized in Table 2 and described in more detail in the Methods). In the strategy labeled “Mitigated - Single Model,” model specifications are compared based on their precision at top k after applying group-specific thresholds to mitigate disparities and the chosen model is evaluated on new data. Selecting a model without regard to fairness and applying disparity mitigation only to the chosen model showed no practical difference in performance on either fairness or accuracy metrics in any of our analyses (see Supplementary Discussion). As a baseline, the “Unmitigated” approach

performs model selection without accounting for any disparities.

Additionally, the “Mitigated - Composite Model” strategy chooses the best-performing model on each subgroup at the model selection stage and picks recall-balancing thresholds across these. We also explored composite models created from fully-decoupled models trained only on subgroup-specific examples (also suggested by Dwork¹³), but saw no difference in performance from the composite method described here in these initial experiments and did not pursue this line of investigation further. Likewise, in considering other fairness-enhancing methods that have been proposed, we found that some (such as the regularization method developed by Zafar¹²) were not well-suited to the “top k ” problem setting and others (such as the methods proposed by Celis¹¹ and Menon²⁴) could be shown to be equivalent to group-specific thresholds under certain conditions (see Supplementary Discussion).

Applying our three strategies to a wide variety of models and hyperparameter combinations including random forests, logistic regression, boosting, and decision trees built for each policy problem, **appreciable recall disparities were present in the unadjusted models in all cases**, ranging from 50% higher recall favoring white individuals in the Inmate Mental Health context to as much as a 250% disparity in the Student Outcomes setting. In the results shown in Fig. 2, strategies yielding larger disparities will have higher values along the y-axis while strategies yielding precision decreases would move left along the x-axis relative to the unmitigated models (solid squares marked with a 'U'). Note that in each context, the top k list size and sensitive attribute for measuring performance and disparities, respectively, can be found in Table 1, reflecting the resource constraints and stakeholder priorities of the policy settings discussed above. However, we see that **mitigating disparities, surprisingly, does not come with any appreciable degradation of overall model performance** on unseen data for any of the policy settings investigated: precision at top k is similar in magnitude for the unmitigated and mitigated models, with little difference between the fairness-enhancing approaches (large error bars in the unmitigated Student Outcomes result reflect the small number of temporal splits and a single cohort with low baseline disparity). Across all four problems, the precision at top k for the mitigated models was statistically indistinguishable from that of the unmitigated models (t-tests for these differences yielded p-values of 0.83 for Education Crowdfunding, 0.47 for Housing Safety, 0.84 for Inmate Mental Health, and 0.71 for Student Outcomes), as well as practically negligible with no difference in precision higher than 1 percentage point. This finding empirically demonstrates that **any inherent trade-off that might be present in adjusting for disparities appear to be practically non-existent**.

To explore these results in more detail, Fig. 3 shows the effect of these bias mitigation strategies on the Inmate Mental Health setting over temporal validation cohorts (see Supplementary Figs. 2-4 for results from the other policy settings). In this context, we also wanted to understand how these methods performed when adjusting for disparities across multiple subgroups, looking at race/ethnicity across white, Black, and Hispanic individuals.

Overall performance (precision at top 500) was similar for models selected from all three strategies over time (Fig. 3a) and recall disparities between white and Black individuals were consistently and significantly improved by the adjustments (Fig. 3b). Without accounting for

equity or fairness, recall for white individuals was 40-100% higher than Black individuals in the chosen models, while both fairness enhancing strategies yielded a ratio near one (parity) for these two subgroups. The results for disparities between Hispanic and white individuals (Fig. 3c) are broadly consistent, but show significantly more variation, which appears to be particularly acute with the composite model approach. We hypothesized that this variability might arise in part from an interaction between the relatively small size of this group in the population (about 11% of each cohort) and process by which the composite model is created: while the model selection process for the single model approach makes use of the full top k set of individuals selected for intervention, the composite approach performs model selection on each group separately, which may be less robust as the groups become smaller.

To better understand the interaction between the overall top k list size, subgroup sizes, and our results, we performed a series of sensitivity experiments. Although in practice, each policy problem comes with a specific value of k determined by the resource constraints of the organization taking action. Fig. 4a-c shows the fairness and accuracy metrics for each strategy at different levels of program resource availability (k). Consistent with the findings above, Fig. 4c shows essentially identical precision at top k performance across all selection strategies, both bias-mitigated and unmitigated, and over all values of k explored. Across all k , the fairness improving strategy choosing a single model consistently showed improvement in disparity metrics for both Black (Fig. 4a) and Hispanic (Fig. 4b) individuals. By contrast, the composite strategy was less effective at equalizing recall across race/ethnicity subgroups, particularly at lower list sizes and with the smaller Hispanic subgroup.

The under-performance of the composite model suggests a mechanism driving this result: evaluating a model’s performance on such a small group may be especially prone to overfitting, choosing a model with an unreasonably high estimate of precision on the subgroup that won’t generalize well into the future. Because constructing the composite model across subgroups involves determining the number of individuals to select from each subgroup with the goal of equalizing recall across them, a process that leads to systematically over-estimating the precision for Hispanic individuals would bias towards selecting a smaller set of individuals than is actually needed (this can be observed in Supplementary Fig. 5b, which shows the fraction of Hispanic individuals in the selected list). When the performance of the chosen model fails to generalize well on the unseen, future data, the under-estimated subgroup size results in a lower-than-expected recall and, thus, a higher disparity. By contrast, choosing a single model across all groups (rather than building a composite) seems likely to be more resilient to this issue, both by virtue of reducing variance in the model selection process (by way of the larger overall sample sizes) and to the degree that any potential reduction in generalization performance might be likely to affect all subgroups similarly.

To further investigate the impact of subgroup size on the stability of the results (e.g., through sampling variation), we performed a resampling experiment to progressively increase the Hispanic fraction in the population, ranging from 5% to 38% (focusing on $k = 500$ based on the actual program resources). Note that although this experiment focused on changing the fraction of Hispanic individuals in the population, the analysis made use of data from all subgroups,

keeping the ratio between Black and white individuals constant (see Supplementary Fig. 6 for additional results). As observed above, the composite model performs less well at reducing disparities between white and Hispanic individuals for the baseline Hispanic fraction of 11% (Fig. 4d). While this under-performance is also observed at other low fractions, this strategy becomes competitive with the other the single model strategy at Hispanic fractions above 20%, consistent with our hypothesis that the small population of Hispanic individuals in the underlying data may be creating a tendency towards overfitting in the formation of the composite models.

Discussion

Taken together, these results suggest a promising and novel conclusion: across a range of policy domains, mitigating disparities does not inherently require new and complex machine learning methods or a prohibitively large accuracy sacrifice as is often assumed. Instead, explicitly defining the fairness goal upfront in the machine learning process and making design choices to achieve that goal by taking active, practical steps, such as the post-hoc bias mitigation strategies investigated here, are important steps to achieving that goal. While it is of course important to keep in mind the empirical nature of this finding, we see an important contribution of this work in providing practical counterexamples to commonly-held assumptions about the nature of the relationship between fairness and accuracy. As such, we see the consistency of these results across policy settings and contexts here as strongly suggestive that practitioners applying machine learning methods to inform other high-stakes policy decisions should question this assumption in their own work. Additionally, this work contributes to a growing body of evidence that, in practice, straightforward approaches such as thoughtful label choice,⁵ model design,⁸ or post-modeling mitigation can effectively reduce biases in many machine learning systems.

A detailed understanding of the mechanisms driving the nominal trade-offs observed here is beyond the scope of the current empirical study, but two factors that we posit may play important roles are the resource-constrained top k setting and the relative predictive performance of the models across subgroups. With regards to the top k context, the relatively small number of individuals who can be selected for intervention will translate into a relatively high score threshold where examples with positive labels are relatively dense both just above and just below the threshold. Because reducing recall disparities involves trading true positives from one subgroup for true positives from another, operating at a region of the score distribution where positive labels are relatively dense means that the perturbations to group-specific score thresholds required to eliminate a given level of disparity may be reasonably small. Here, our hypothesis is that settings in which resources are less constrained may pose a greater challenge to disparity mitigation, because operating deeper in the score distribution (that is, at lower predicted score thresholds) means the recall gradient will typically be flatter and larger adjustments may be required to achieve equity (and, hence, resulting in greater trade-offs). Likewise, the relative performance of the model across subgroups will be a key factor in determining both

how much the score thresholds will need to change for each group to collect enough true positives to achieve recall equity as well as how many false positives are included in the process of doing so. Contexts in which the available features are much less predictive of the outcome for some subgroups relative to others could lead to both larger baseline disparities (if recall at a given threshold is higher for the more readily-predicted group) and more significant fairness-accuracy trade-offs (by necessitating larger adjustments to achieve recall equity). Although these two factors may be more salient in other settings, our focus on high-impact benefit allocation programs with limited resources reflects a setting that occurs very commonly in policy contexts, and frequently encountered in our applied work^{6,7,9,15} and the empirical results here suggest some of these concerns may not be prohibitively common in practice. Nevertheless, an important avenue for future work will be further exploring these hypothesis and beginning to develop a better theoretical understanding of the conditions under which improvements in equity will or will not come at appreciable expense in accuracy.

Much has also been written about the wide variety of fairness metrics that may be relevant depending on the context,^{1, 14, 15, 21, 22} and further exploration of the fairness-accuracy trade-offs in those contexts is certainly warranted, particularly where balancing multiple fairness metrics may be desirable. Likewise, it may be possible that there is a tension between improving fairness across different attributes (e.g., sex and race) or at the intersection of attributes. Future work should also extend these results to explore the impact not only on equity in decision making, but also equity in longer-term outcomes and implications in a legal context such as those discussed by Huq.²⁵

Social Impact of This Work

Beyond the immediate implications for the policy and social impact problems explored here, our goal in presenting these empirical findings here is to better inform the use of machine learning in high-stakes decisions by challenging the commonly held belief that there necessarily is a trade-off between “accuracy” and “fairness” in these applications. However, it is also important to note that fairness is not only a function of a model’s predictions but also how those predictions are acted on by human decision-makers and, more broadly, the context in which it operates, with historical, cultural, and structural sources of inequities that society as a whole must strive to overcome through the ongoing process of remaking itself to better reflect its highest ideals of justice and equity. Improving the fairness of the machine learning models that continue to find growing applications in critical decisions that affect many aspects of people’s lives may be only one small element of that process, but we hope this work will inspire researchers, policymakers, and data science practitioners alike to explicitly consider fairness as a goal and take steps, such as those proposed here, in their work that can, collectively, contribute to bending the long arc of history towards a more just and equitable society.

Methods

Policy Contexts and Data

A key aim of this work was to explore the fairness-accuracy trade-offs encountered in practice in the context of machine learning applications for public policy settings. As such, we drew on several projects we have worked on in partnership with government agencies across policy domains. We describe these contexts briefly in the main text and provide more details about each setting below:

Inmate Mental Health Untreated mental health conditions often result in a negative spiral, which can culminate in repeated periods of incarceration with long term consequences both for the affected individual and the community as a whole.²⁶ Surveys of inmate populations have suggested a high prevalence of multiple and complex needs, with 64% of people in local jails suffering from mental health issues and 55% meeting criteria for substance abuse or dependence.²⁷ The criminal justice system is poorly suited to address these needs, yet houses three times as many individuals with serious mental illness as hospitals.²⁸ In 2016, Johnson County, KS, partnered with our group to help them break this cycle of incarceration by identifying individuals who might benefit from outreach with mental health resources and are at risk for future incarceration. While the Johnson County Mental Health Center (JCMHC) currently provides services to the jail population, needs are generally identified reactively, for instance through screening instruments individuals fill out when entering jail. The new program being developed will supplement these existing approaches by adding a new automatic referral system for people who are at risk of being booked into jail, with the hope that they can be outreached before they return to jail. Through our partnership, the county provided administrative data from their mental health center, jail system, police arrests, and ambulance runs. Modeling was focused on a cohort of Johnson County residents with any history of mental health need who had been released from jail within the past three years. Early results from this work were described previously.⁷ A field evaluation of the predictive model is ongoing at the time of this writing, but validation on historical data demonstrated a 12% improvement over a baseline based on the number of bookings in the prior year and 4.8-fold increase over the population prevalence.

Housing Safety The Multiple Housing team in San Jose's Code Enforcement Office is tasked with protecting the occupants of properties with three or more units, such as apartment buildings, fraternities, sororities, and hotels. They do so by conducting routine inspections of these properties, looking for everything from blight and pest infestations to faulty construction and fire hazards (see work by Holtzen²⁹ and Klein³⁰ for a discussion of the importance of housing inspections to public health). Although the city of San Jose inspects all of the properties on its Multiple Housing roster over time, and expects to find minor violations at many of them, it is important that they can identify and mitigate dangerous situations early to prevent accidents. With more than 4,500 multiple housing properties in San Jose, CA – many of which comprise multiple buildings and hundreds of units – it is not possible for the city to inspect every unit every year. San Jose recently instituted a tiered approach to prioritizing inspections, inspecting riskier properties more frequently and thoroughly. Although the tier system helped focus

inspections on riskier properties, the new system has its limitations. The city evaluates tier assignments for properties infrequently (every 3 to 6 years), and these adjustments require a great deal of expertise and manual work while leaving out a rich amount of information. In order to provide a more nuanced view of properties' violation risk over time and allow for more efficient scheduling of inspections, the Code Enforcement Office partnered with us to develop a model to predict the risk that a serious violation would be found if a given property was prioritized for inspection (similar tools have been developed for allocating fire inspections in New York³¹ and health inspections in Boston³²). Evaluation of the model on historical data indicated that it could provide a 30% increase in precision relative to the current tier system and the model's predictive accuracy was confirmed during a 4-month field trial in 2017.

Student Outcomes Each year from 2010 through 2016, 15-29% of students enrolled in school in El Salvador did not return to school in the following year. This high dropout rate is cause for serious concern, with significant consequences for economic productivity, workforce skill, inclusiveness of growth, social cohesion, and increasing youth risks.^{33,34} El Salvador's Ministry of Education has programs available to support students with the goal of reducing these high dropout rates, but the budget for these programs is not large enough to reach every student and school in El Salvador. Predictive modeling has been deployed to help schools identify students at risk of dropping out in several contexts³⁵⁻³⁷ and El Salvador partnered with us in 2018 to make use of these methods to focus their limited resources on the students at highest risk of not returning each year. Student-level data was provided by the Ministry of Education, including demographics, urbanicity, school-level resources (e.g., classrooms, computers, etc), gang and drug violence, family characteristics, attendance records, and grade repetition. For the present study, we focused on the state of San Salvador and identifying the 10,000 highest-risk students, considering annual cohorts of approximately 300,000 students and drawing on 5 years' of prior examples as training data.

Education Crowdfunding Many schools in the United States, particularly in poorer communities, face funding shortages.³⁸ Often, teachers themselves are left to fill this gap, purchasing supplies for their classrooms when they have the individual resources to do so.³⁹ The non-profit DonorsChoose was founded in 2000 to help alleviate these shortages by providing a platform where teachers post project requests focused on their classroom needs and community members can make individual contributions to support these projects. Since 2000, they have facilitated \$970 million in donations to 40 million students in the United States.⁴⁰ However, approximately one third of all projects posted on the platform fail to reach their funding goal. Here, we make use of a dataset DonorsChoose made publicly available for the 2014 KDD Cup (an annual data mining competition) including information about projects, the schools posting them, and donations they received. Because the other case studies explored here focused on proprietary and often sensitive data shared with us under data use agreements that cannot be made publicly available, we included a case study surrounding this publicly-available dataset. While we have not partnered with DonorsChoose to deploy the machine learning system described, we otherwise treated this case study as we would any of our applied projects. Here, we consider a resource-constrained effort to assist projects at risk of going unfunded (for instance, providing

a review and consultation) capable of helping 1,000 projects in a 2-month window, focusing on the most recent 2 years’ of data available in the extract (earlier data had far fewer projects and instability in the baseline funding rates as the platform ramped up). This dataset is publicly available at [kaggle.com](https://www.kaggle.com).⁴¹

Machine Learning Details

All machine learning models, including feature engineering, model training, and performance evaluation were run using our open-source python ML pipeline package, `trriage`. Machine learning methods used are from `sklearn` (a python package) or `catwalk` (a component of `trriage` for baselines methods as well as `ScaledLogisticRegression`, which wraps the `sklearn` logistic regression to ensure input features/predictors as scaled between 0 and 1). The modeling grid for each project is described in Supplementary Tables 1-4, reflecting the modeling space explored by the teams working on each project. For each estimator in the tables, the grid search considered reflects the full cross-product of the hyperparameter values specified. Here we make use of a variety of state-of-the-art machine learning methods for binary classification problems. Although precision in the top k is the metric of interest in all of the settings discussed here (as described in the main text), few methods have been developed that seek to optimize for this metric directly. Instead, we are concerned with the relationship between fairness (measured here in terms of recall disparities) and accuracy (measured here in terms of precision in the top k) through real-world settings in practice. As such, we make use of well-established methods that are widely applied in practical settings, which themselves optimize for a variety of underlying target metrics (such as minimizing the regularized logistic loss in logistic regression or maximizing the information gain at each split in tree-based models). By training a large grid of estimator types and hyperparameter values, optimization for performance in terms of precision in the top k can then be performed through the process of model selection over validation sets. As illustrated in Fig. 1b, we used a strategy of inter-temporal cross-validation (as described by Roberts⁴² and Ye⁴³) to ensure that model evaluation and selection was done in a manner that reflected performance on novel data while guarding against “leakage” of information from the future affecting past results.

The method we used for mitigating disparities by post-modeling adjustment involving choosing sub-group specific thresholds (see Fig. 1a) was described in detail in our previous case study¹⁵ and draws on the idea of “equality of opportunity” discussed by Hardt.¹⁴ In brief, because the notion of fairness relevant in these policy settings relies on equalizing recall across groups, and recall monotonically increases with depth traversed in a model score, unique score thresholds that balance recall across groups can be readily found for a given combined list size. For each model, we calculate within-group recall values up to each individual in an initial validation set (purple / t_{-1} in Fig. 1b), order the combined set by within-group recall and take the top k individuals from this reordered set, calculating k_g for each group g such that $\sum k_g = k$ (the total top k list size desired) and recall is balanced across groups. To evaluate this process on novel data, models were tested on a future cohort (tan / t_0 in Fig. 1b) and the top k_g examples (ranked by score, then randomly to break ties) from each sub-group were selected to measure

precision at top k and recall disparities. In the process of model selection, we explored applying these disparity-mitigating thresholds either before choosing a model specification (“Mitigated - Single Model”) or after (“Mitigated - Unadj. Model Seln.”), finding no substantive difference in performance (Supplementary Fig. 1). For the “Mitigated - Composite Model” approach, a similar method was used, but within-group precision up to each individual is calculated for each model as well to determine the best model specification for each sub-group at each list depth (drawing on the ideas suggested by Dwork and colleagues¹³), then k_g and group-specific model specifications are chosen for evaluation on novel data.

Code for all four projects including `triage` configuration files specifying the full feature sets used as well as code used to mitigate disparities and evaluate fairness is available at github.com/dssg/peeps-chili

Data Availability

Data from the Inmate Mental Health context was shared through a partnership and data use agreement with the county government of Johnson County, KS (which collected and made available data from the county- and city-level agencies in their jurisdiction as described in the Methods above). Data from the Housing Safety context was shared through a partnership and data use agreement with the Code Enforcement Division in the city of San Jose, CA. Data from the Student Outcomes setting was shared through a partnership and data use agreement with the Ministry of Education in El Salvador. Although the sensitive nature of the data for these three contexts required that the work was performed under strict DUAs and the data cannot be made publicly available, researchers or practitioners interested in collaborating on these projects or with the agencies involved should contact the corresponding author (rayid@cmu.edu) for more information and introductions. The Education Crowdfunding dataset, however, is publicly available at: kaggle.com/c/kdd-cup-2014-predicting-excitement-at-donors-choose. Additionally, a database extract with model outputs and disparity mitigation results using this dataset is available for download (see replication instructions in the code repository noted below).

Code Availability

The code used here for modeling, disparity mitigation, and analysis for all four projects is available at github.com/dssg/peeps-chili.⁴⁴ Complete instructions for replication of the Education Crowdfunding results reported here can be found in the README of this repository, along with a step-by-step jupyter notebook to perform the analysis.

Acknowledgments

We would like to thank the Data Science for Social Good Fellowship fellows, project partners, and funders as well as our colleagues at the Center for Data Science and Public Policy at University of Chicago for the initial work on projects that were extended and used in this study.

We also thank Kasun Amarasinghe for helpful discussions on the study and drafts of this paper. Parts of this work were funded by the National Science Foundation under grant IIS-2040929 (KTR, RG) and by a grant (unnumbered) from the C3.ai Digital Transformation Institute (KTR, HL, RG).

Author Contributions Statement

KTR: Conceptualization, Methodology, Software, Investigation, Writing - Original Draft; HL: Investigation, Software, Writing - Review & Editing; RG: Conceptualization, Supervision, Funding acquisition, Writing - Review & Editing

Competing Interests Statement

The authors have no conflicts of interest to declare.

Table 1: Policy Settings and Data Details

	Inmate Mental Health	Housing Safety	Student Outcomes	Education Crowdfunding
Prediction Task	Jail booking within the next 12 months	Housing unit having a violation within the next year	Student not returning to school next year	Project not getting fully funded within 4 months
Timespan	2013-01-01 to 2019-04-01	2011-01-01 to 2017-06-01	2009-01-01 to 2018-01-01	2010-01-01 to 2014-01-01
# of Entities	61,192	4,593	801,242	210,310
# of Features	3,465	1,657	220	319
Base Rate	0.12	0.43	0.25	0.24
Evaluation Metric	Precision at top 500	Precision at top 500	Precision at top 10,000	Precision at top 1,000
Sensitive Attribute	Race	Median Income	Age Relative to Grade	Poverty Level

Table 2: Descriptions of Model Selection Strategies

Strategy	Description
Unmitigated	Baseline strategy with no equity adjustments
Mitigated - Composite Model	The highest-precision model is chosen for each subgroup and a composite model is formed combining these, with equity-balancing thresholds used when calculating test performance
Mitigated - Single Model	All models are adjusted for recall equity and then the best model is selected based on precision and equity-balancing thresholds are applied for calculating test performance
Mitigated - Unadj. Model Seln.	Model is selected based on precision then equity-balancing thresholds are applied for calculating test performance (See the Supplementary Discussion for results from this approach)

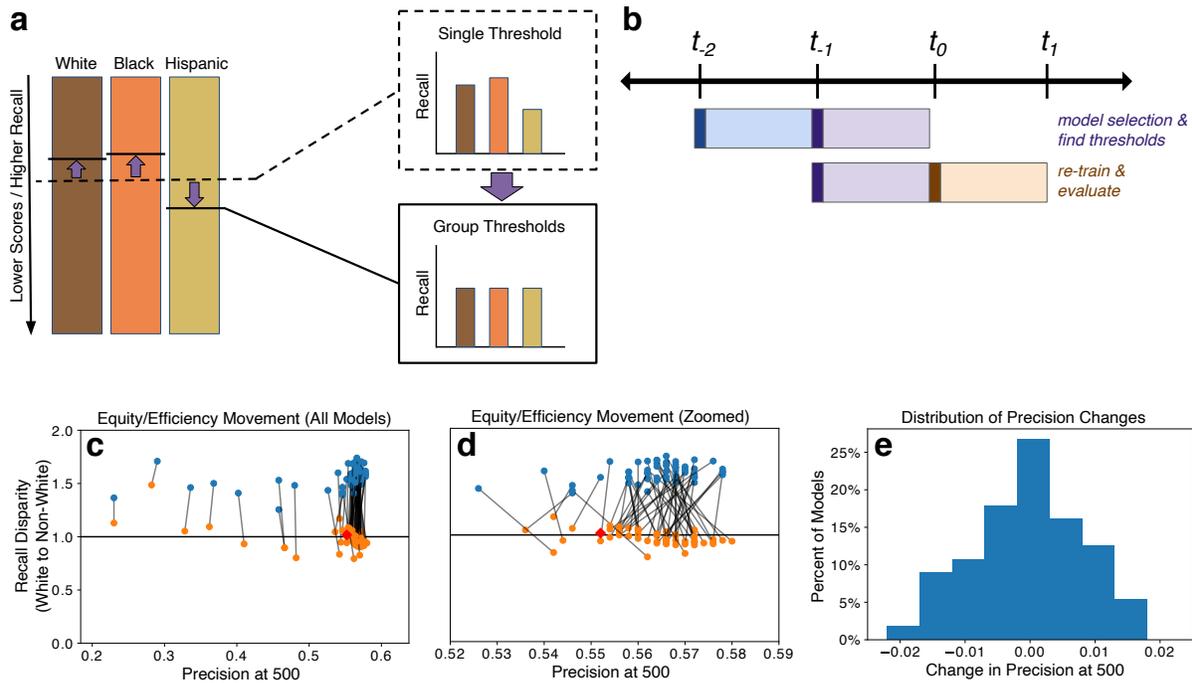


Figure 1: Illustration of the methods used and motivating results. (a) Subgroup-specific thresholds are applied to a modeled risk score to improve the recall equity among individuals chosen for intervention while maintaining a desired overall list size. (b) Temporal validation strategy: a grid of models is trained using examples as of t_{-2} (dark blue, with labels derived from the time shown in light blue) and predictions on a cohort as of t_{-1} (dark purple with labels derived from the time shown in light purple) are used to determine the equity-balancing thresholds described in (a). Models are then re-trained on this cohort for “current day” predictions as of t_0 (dark tan, with labels in light tan) used to evaluate model performance with equity adjustments. (c and d): Changes in race/ethnicity recall disparities before (blue) and after (orange) making post-hoc score adjustments for fairness in the Inmate Mental Health context. (c) shows all model specifications and (d) shows the cluster of well-performing models. The red diamond reflects the performance of a composite model combining the best-performing model for each subgroup. (e) Distribution of precision changes after adjusting for disparities for the models shown in (c), relative to the precision attained by the same model specification without adjustment (that is, the difference along the x-axis of the blue and orange dots).

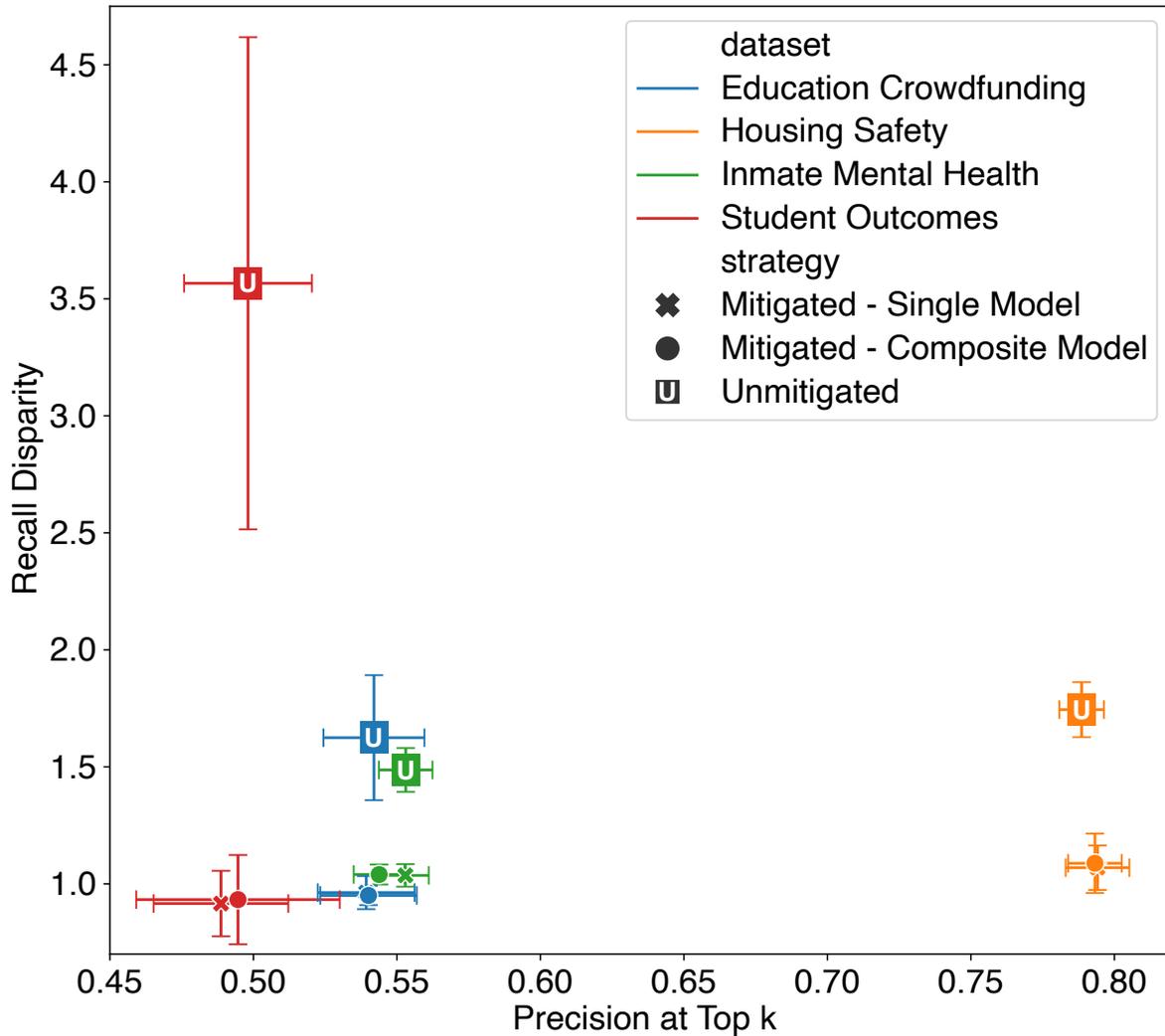


Figure 2: Comparing equity (recall disparity) and performance (precision at top k) metrics for different model selection strategies between different policy contexts. In the Education Crowdfunding context, models are evaluated at $k = 1000$ across 10 temporal validation cohorts; in the Inmate Mental Health context, $k = 500$ across 10 temporal validation cohorts; in the Student Outcomes context, $k = 10,000$ across 5 temporal validation cohorts; and in the Housing Safety context, $k = 500$ across 9 temporal validation cohorts. Unmitigated (baseline) models are shown as solid squares marked with a ‘U’. Decreases of strategies involving disparity mitigations relative to the y-axis demonstrate improvements in equity while showing little or no decrease in overall performance (that is, leftward movement on the x-axis). Error bars reflect 95% confidence intervals across temporal validation cohorts.

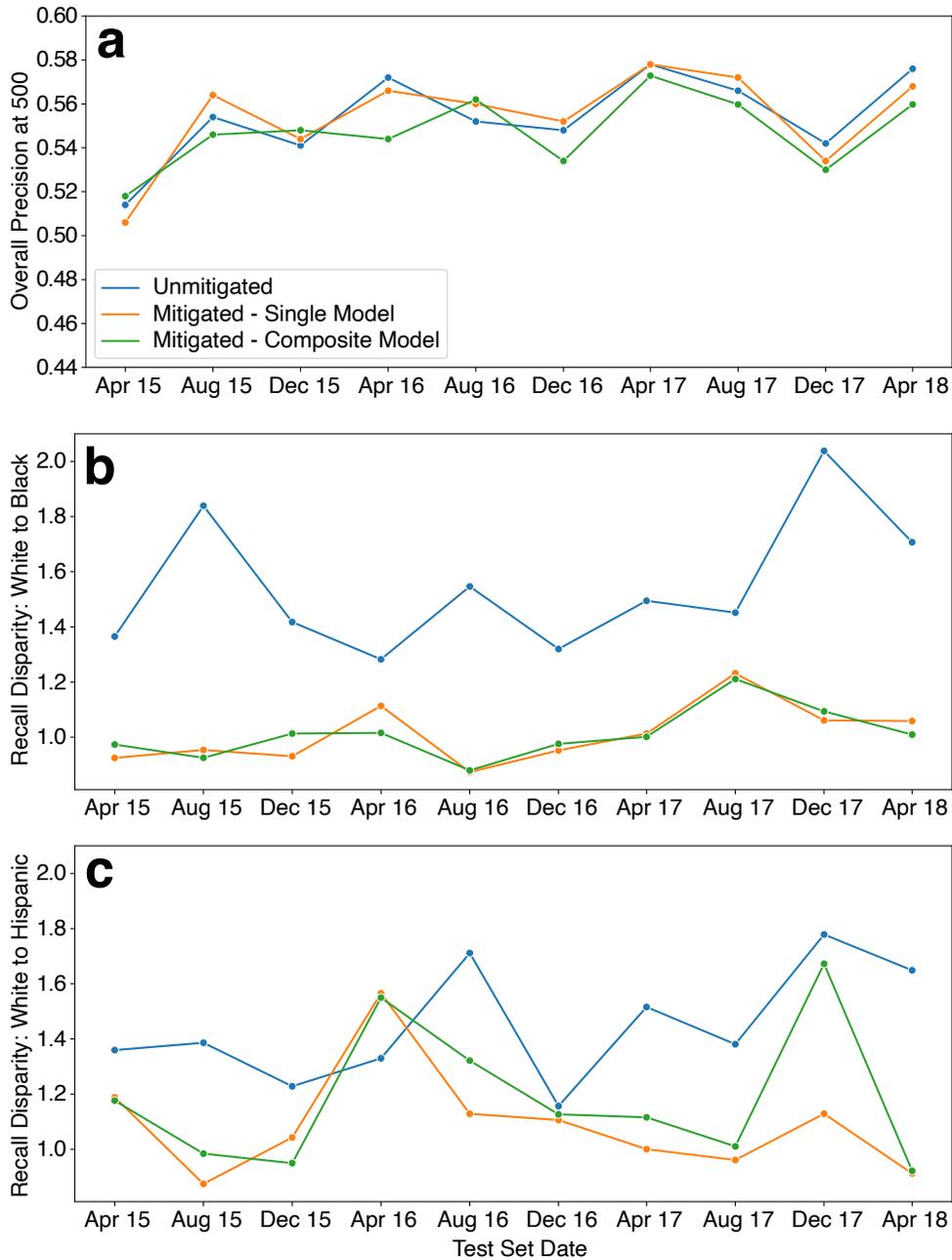


Figure 3: Comparing disparity and performance metrics over time for different model selection strategies. Results from the Inmate Mental Health policy setting using a total list size $k = 500$. (a) Model performance in terms of overall precision at top 500. (b) Recall disparities between white and Black individuals. (c) Recall disparities between white and Hispanic individuals.

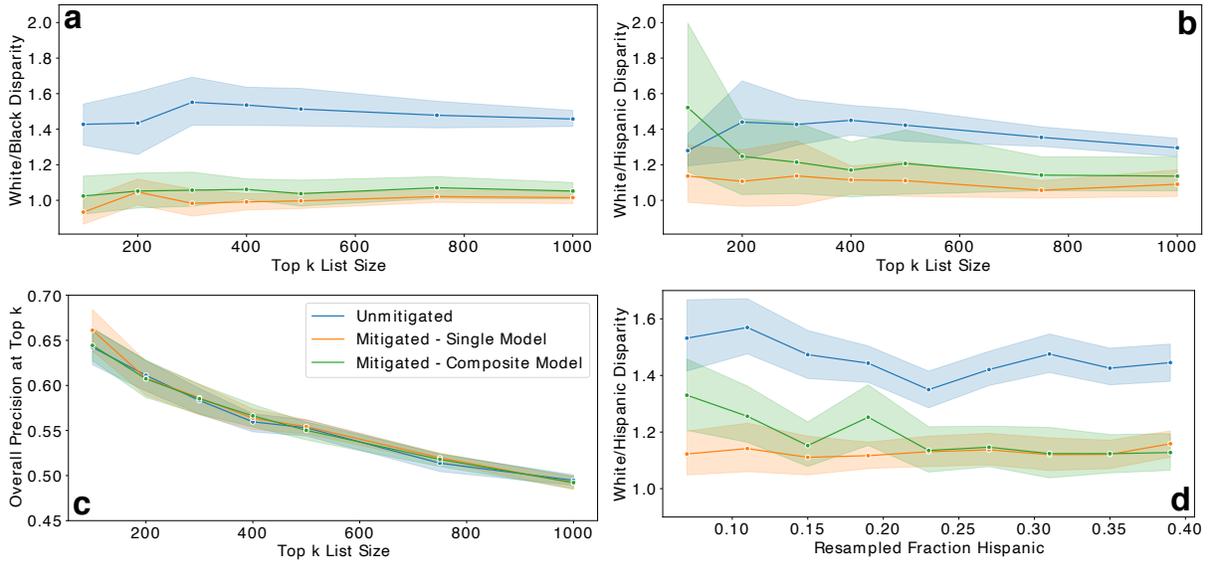


Figure 4: Comparing disparity and performance metrics across program scale and protected group size in the Inmate Mental Health policy setting. (a to c) Variation in results by list size: (a) white/Black disparities, (b) white/Hispanic disparities, and (c) overall precision at top k . (d) Disparities between white and Hispanic individuals by resampled size of the protected group in the overall population (using a list size of $k = 500$). Shaded intervals reflect 95% confidence intervals from variation across temporal validation splits (a-c) as well as bootstrap samples in (d).

References

- ¹ Chouldechova, A. Fair prediction with disparate impact: a study of bias in recidivism prediction instruments. *Big Data* **5**, 153–163 (2017).
- ² Skeem, J. L. & Lowenkamp, C. T. Risk, race, and recidivism: predictive bias and disparate impact. *Criminology* **54**, 680–712 (2016).
- ³ Angwin, J., Larson, J., Mattu, S. & Kirchner, L. “Machine bias,” *ProPublica* (23 May 2016); www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.
- ⁴ Raghavan, M., Barocas, S., Kleinberg, J. & Levy, K. Mitigating bias in algorithmic hiring: evaluating claims and practices. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2020), pp. 469–481.
- ⁵ Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* **336**, 447–453 (2019).
- ⁶ Ramachandran, A., et al. Predictive analytics for retention in care in an urban HIV clinic. *Scientific Reports* 10.1038/s41598-020-62729-x (2020).
- ⁷ Bauman, M. J., et al. Reducing incarceration through prioritized interventions. *Proceedings of the Conference on Computing and Sustainable Societies (COMPASS)* (ACM, 2018), pp. 1–8.
- ⁸ Chouldechova, A., et al. A case study of algorithm-assisted decision making in child maltreatment hotline screening decisions. *Proceedings of Machine Learning Research* **81**, 134–148 (2018).
- ⁹ Potash, E., et al. Predictive modeling for public health: preventing childhood lead poisoning. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (ACM, 2015), pp. 2039–2047.
- ¹⁰ Chen, I. Y., Johansson, F. D. & Sontag, D. Why is my classifier discriminatory? *Advances in Neural Information Processing Systems 31* (NIPS, 2018), pp. 3539–3550.
- ¹¹ Celis, L. E., Huang, L., Keswani, V. & Vishnoi, N. K. Classification with fairness constraints: a meta-algorithm with provable guarantees. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2019), pp. 319–328.
- ¹² Zafar, M. B., Valera, I., Rodriguez, M. G. & Gummadi, K. P. Fairness beyond disparate treatment and disparate impact: learning classification without disparate mistreatment. *26th International World Wide Web Conference (WWW, 2017)*, pp. 1171–1180.

- ¹³ Dwork, C., Immorlica, N., Kalai, A. T. & Leiserson, M. Decoupled classifiers for group-fair and efficient machine learning. *Proceedings of Machine Learning Research* **81**, 119–133 (2018).
- ¹⁴ Hardt, M., Price, E. & Srebro, N. Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems 29* (NIPS, 2016), pp. 3315–3323.
- ¹⁵ Rodolfa, K. T., et al. Case study: predictive fairness to reduce misdemeanor recidivism through social service interventions. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2020), pp. 142–153.
- ¹⁶ Heidari, H., Gummadi, K. P., Ferrari, C. & Krause, A. Fairness behind a veil of ignorance: a welfare analysis for automated decision making. *Advances in Neural Information Processing Systems* (NIPS, 2018), pp. 1265–1276.
- ¹⁷ Friedler, S. A., et al. A comparative study of fairness-enhancing interventions in machine learning. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2019), pp. 329–338.
- ¹⁸ Kearns, M., Roth, A., Neel, S. & Wu, Z. S. An empirical study of rich subgroup fairness for machine learning. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2019), pp. 100–109.
- ¹⁹ Zafar, M. B., Valera, I., Rogniguez, M. G. & Gummadi, K. P. Fairness constraints: mechanisms for fair classification. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics* (PMLR, 2017), pp. 962–970.
- ²⁰ Ghani, R., Walsh, J., Wang, J. Top 10 Ways Your Machine Learning Models May Have Leakage. <http://www.rayidghani.com/2020/01/24/top-10-ways-your-machine-learning-models-may-have-leakage/>.
- ²¹ Verma, S. & Rubin, J. Fairness definitions explained. *International Workshop on Software Fairness* (IEEE/ACM, 2018), pp. 1–7.
- ²² Gajane, P. & Pechenizkiy, M. On Formalizing Fairness in Prediction with Machine Learning. Preprint at <https://arxiv.org/pdf/1710.03184> (2018).
- ²³ Kleinberg, J. M., Mullainathan, S. & Raghavan, M. Inherent Trade-Offs in the Fair Determination of Risk Scores. *8th Innovations in Theoretical Computer Science Conference* (ITCS, 2017), pp. 1–43.
- ²⁴ Krishna Menon, A. & Williamson, R. C. The cost of fairness in binary classification. *Proceedings of Machine Learning Research* (PMLR, 2018), pp. 1–12.

- ²⁵ Huq, A. Racial equity in algorithmic criminal justice. *Duke Law Journal* **68**, 1043–1134 (2019).
- ²⁶ Hamilton, M. People with complex needs and the criminal justice system. *Current Issues in Criminal Justice* **22**, 307–324 (2010).
- ²⁷ James, D. J. & Glaze, L. E. “Mental health problems of prison and jail inmates” (Department of Justice, Bureau of Justice Statistics, 2006; <https://www.bjs.gov/content/pub/pdf/mhppji.pdf>)
- ²⁸ Fuller Torrey, E., Kennard, A. D., Eslinger, D., Lamb, R. & Pavle, J. “More mentally ill persons are in jails and prisons than hospitals: a survey of the states” (Treatment Advocacy Center and National Sheriffs’ Association, 2010; http://tulare.networkofcare.org/library/final_jails_v_hospitals_study1.pdf)
- ²⁹ Holtzen, H., Klein, E. G., Keller, B. & Hood, N. Perceptions of physical inspections as a tool to protect housing quality and promote health equity. *Journal of Health Care for the Poor and Underserved* **27**, 549–559 (2016).
- ³⁰ Klein, E., Keller, B., Hood, N. & Holtzen, H. Affordable housing and health: a health impact assessment on physical inspection frequency. *Journal of Public Health Management and Practice* **21**, 368–374 (2015).
- ³¹ Athey, S. Beyond prediction: using big data for policy problems. *Science* **355**, 483–485 (2017).
- ³² Glaeser, E. L., Hillis, A., Kominers, S. D. & Luca, M. Crowdsourcing city government: using tournaments to improve inspection accuracy. *Am. Econ. Rev.* **106**, 114–118 (2016).
- ³³ Levin, H. M. & Belfield, C. “The price we pay: economic and social consequences of inadequate education” (Brookings Institution Press, 2007).
- ³⁴ Atwell, M. N., Balfanz, R., Bridgeland, J. & Ingram, E. “Building a grad nation” (America’s Promise Alliance, 2019; <https://www.americaspromise.org/2019-building-grad-nation-report>)
- ³⁵ Lakkaraju, H., et al. A machine learning framework to identify students at risk of adverse academic outcomes. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (ACM, 2015), pp. 1909–1918.
- ³⁶ Aguiar, E., et al. Who, when, and why: a machine learning approach to prioritizing students at risk of not graduating high school on time. *Proceedings of the Learning Analytics and Knowledge Conference* (ACM, 2015), pp. 93–102.

- ³⁷ Bowers, A. J., Sprott, R. & Taff, S. A. Do we know who will drop out? A review of the predictors of dropping out of high school: precision, sensitivity, and specificity. *The High School Journal* **96**, 77–100 (2012).
- ³⁸ Morgan, I. & Amerikaner, A. “Funding gaps 2018” (The Education Trust, 2018; https://edtrust.org/wp-content/uploads/2014/09/FundingGapReport_2018_FINAL.pdf).
- ³⁹ Hurza, M. “What do teachers spend on supplies” (Adopt a Classroom, 2015; <https://www.adoptaclassroom.org/2015/09/15/infographic-recent-aac-survey-results-on-teacher-spending/>).
- ⁴⁰ Statistics available from DonorsChoose at <https://www.donorschoose.org/about> (accessed: 23 June 2020).
- ⁴¹ Data available at <https://www.kaggle.com/c/kdd-cup-2014-predicting-excitement-at-donors-choose/data> (accessed: 23 June 2020).
- ⁴² Roberts, D. R., et al. Cross-validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure. *Ecography* **40**, 913–929 (2017).
- ⁴³ Ye, T., et al. Using machine learning to help vulnerable tenants in New York city. *Proceedings of the Conference on Computing and Sustainable Societies (COMPASS)* (ACM, 2019), pp. 248–258.
- ⁴⁴ Rodolfa, K.T. & Lamba, H. dssg/peeps-chili: Release for trade-offs submission. (2021) doi:10.5281/zenodo.5173254.

Supplementary Information

Supplementary Discussion

The results discussed in the main text focus on post-hoc score adjustment methods, drawing on several previous lines of work³⁻⁵ and finds these straightforward methods suffice to significantly reduce disparities with little accuracy trade-off in the settings considered. However, we also did some preliminary work exploring a number of other methods that have been proposed, particularly to enhance fairness during the modeling process itself. In public policy settings with limited resources, regularization methods (such as that described by Zafar and colleagues²) provide a particular challenge. These techniques typically seek to find the best overall classifier subject to some fairness constraint. This may be well-suited to contexts where there are no hard constraints on the resources available to act on predictive positives (for instance, in deciding whether to order a relatively low-cost medical diagnostic test).

However, many applications, particularly in public policy contexts, are subject to the further constraint of limited resources — these applications are best formulated as “top-k” problems, but unfortunately this formulation introduces a non-convex optimization that is not readily integrated into current fairness-constrained methods. Likewise, naively thresholding the resulting score from these methods to yield a set number of predicted positives provides no guarantee that the fairness constraints will hold. Supplementary Fig. 7 provides an example with data from the Inmate Mental Health context: using the method described by Zafar² to generate a score with a false negative rate fairness constraint (note that $FNR = 1 - TPR$, so this constraint is equivalent to recall equity) and choosing the top 500 individuals performs no better at balancing equity than choosing the top 500 from an unconstrained score and far worse than choosing group-specific thresholds to balance equity as shown in the main text. Notably, this result isn’t a criticism of Zafar’s methods when applied in appropriate contexts, but rather an indication of the limitations of current fairness-constrained methods in the resource-constrained setting we frequently encounter in machine learning to support public policy decision making.

We also explored the methods proposed by Celis¹ and Menon and Williamson,⁶ but noticed that both methods could be shown to be equivalent to group-specific scaling or thresholding of an underlying estimate of the relevant probability distribution when applied to a single fairness metric and sub-group membership is known (notably, the method proposed by Celis¹ seems quite flexible to balancing multiple metrics or where group membership is itself being modeled as well). For a single, monotonic metric like recall/equality of opportunity, there will

be a unique balanced solution of a given total size, so any method relying on post-hoc score adjustments should yield similar results.

To see this in the context of the method described by Celis and colleagues,¹ we can start from their observation that any equity definition balancing a confusion-matrix statistic can be represented in the following form (following the notation used in their work):

$$q_{lin}^{(i)} = \frac{\alpha_0^{(i)} + \sum_{j \in [k]} \alpha_j^{(i)} P[f = 1 | G_i, A_j^{(i)}]}{\beta_0^{(i)} + \sum_{j \in [l]} \beta_j^{(i)} P[f = 1 | G_i, B_j^{(i)}]} \quad (1)$$

Here, i reflects subgroup membership, $A_j^{(i)}$ and $B_j^{(i)}$ are events such as $(Y = 1)$ and $\alpha_j^{(i)}$ and $\beta_j^{(i)}$ are parameters defined to represent generalized form of equity function (likewise, $k, l \geq 0$ are integer values allowing the function to be written as linear combination of terms). Further, some metrics (including recall (that is, TPR) as we consider in this work) can be represented with a simplified linear form in which $\beta_0^{(i)} = 1$ and $\beta_1^{(i)} = 0$, giving:

$$q_{lin}^{(i)} = \alpha_0^{(i)} + \sum_{j \in [k]} \alpha_j^{(i)} P[f = 1 | G_i, A_j^{(i)}] \quad (2)$$

In this case, they show that fairness-improved score of an instance being classified is given by:

$$s_\lambda(x) = \eta(x) - 0.5 + \sum_{i \in [p]} \lambda_i \left(\sum_{j \in [k]} \frac{\alpha_j^{(i)}}{\pi_j^{(i)}} \eta_j^{(i)}(x) \right) \quad (3)$$

Where, $\eta(x) = P[Y = 1 | X = x]$, $\eta_j^{(i)}(x) = Pr[G_i, A_j^i | X = x]$ and $\pi_j^{(i)} = Pr[G_i, A_j^{(i)}]$, while λ_i represents langrangian parameters, used to solve the optimization problem.

For recall specifically (defined as $TPR = P[f = 1 | G_i, Y = 1]$), the sum in q_{lin} has a single term ($j = 1$), with $\alpha_0^{(i)} = 0$, $\alpha_1^{(i)} = 1$, $A_1^{(i)} = (Y = 1)$. Further (and without loss of generality), for the case in which there are two sub-groups of the protected attribute (e.g., men and women), we can expand the sums and rewrite the adjusted score as:

$$s_{recall}(x) = P(Y = 1 | X = x) - 0.5 + \lambda_0 \left[\frac{1}{P[G_0, Y = 1]} P[G_0, Y = 1 | X] \right] + \lambda_1 \left[\frac{1}{P[G_1, Y = 1]} P[G_1, Y = 1 | X] \right] \quad (4)$$

This formulation can be quite useful for balancing recall equity in the case where group membership is itself being inferred. However, where we know $G = g$ for a given X , we note that this will pick out a single term as either $P[G_1, Y = 1 | X] = 0$ or $P[G_0, Y = 1 | X] = 0$

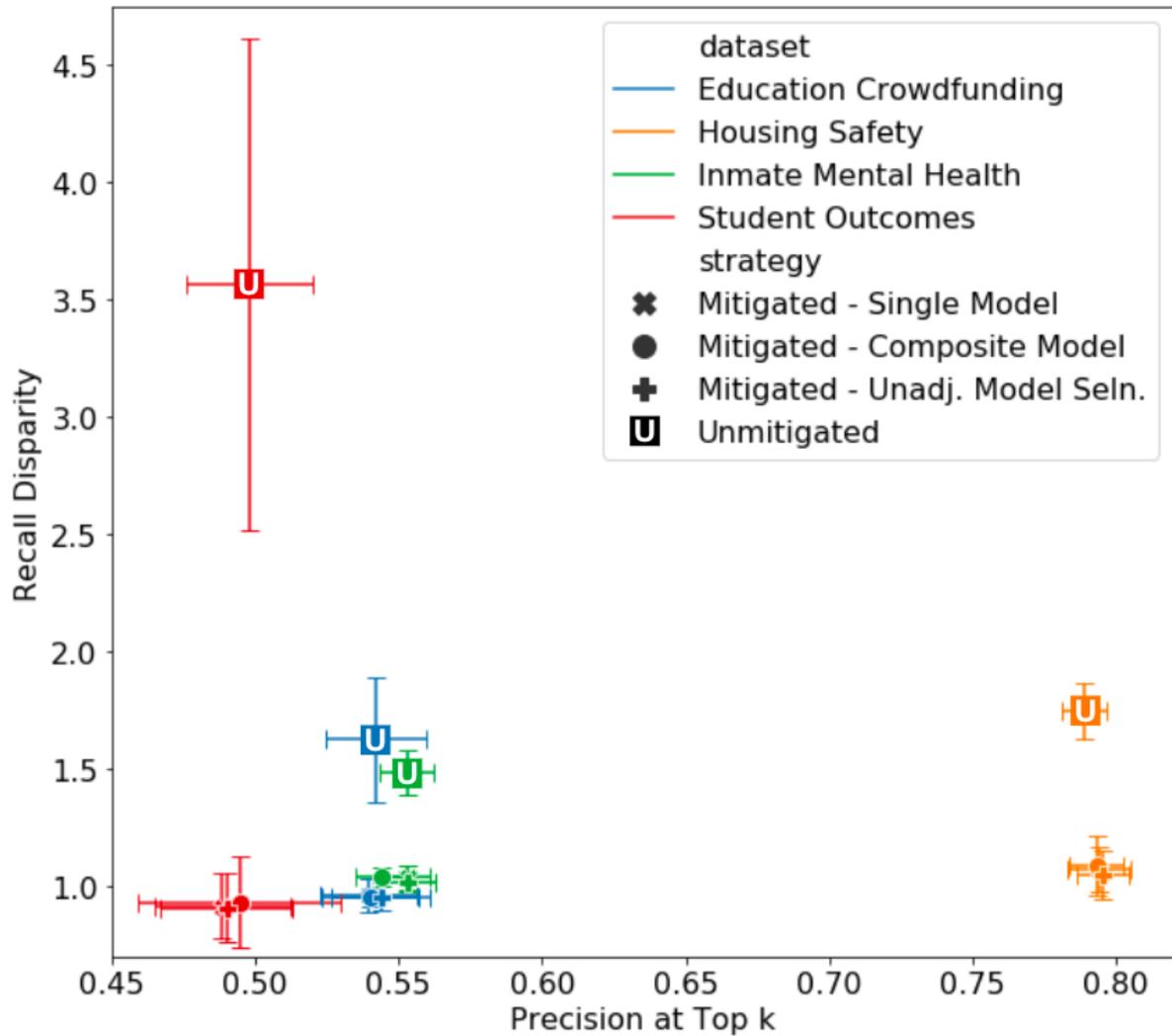
(note that this will be true if there were more sub-groups as well, simply adding additional terms). So, for a given example X , we obtain:

$$s_{recall}(x) = P(Y = 1 | X = x) - 0.5 + \lambda_g \left[\frac{1}{P[G = g, Y = 1]} P[Y = 1 | X] \right] \quad (5)$$

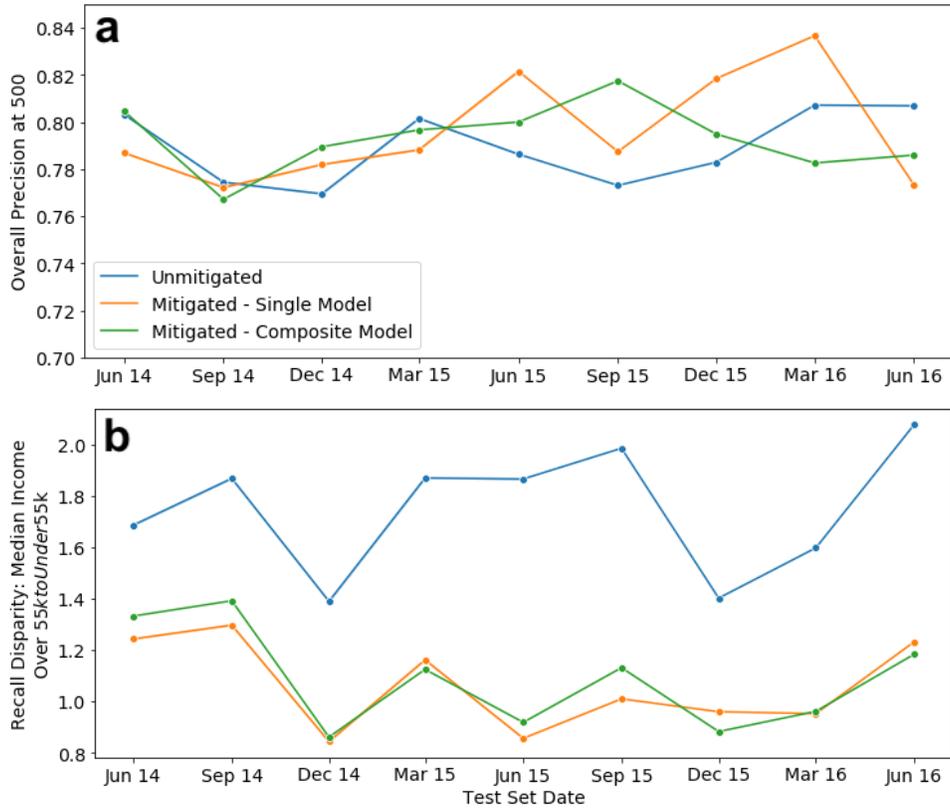
Rearranging leads to:

$$s_{recall}(x) = \left[1 + \frac{\lambda_g}{P[G = g, Y = 1]} \right] P(Y = 1 | X) - 0.5 \quad (6)$$

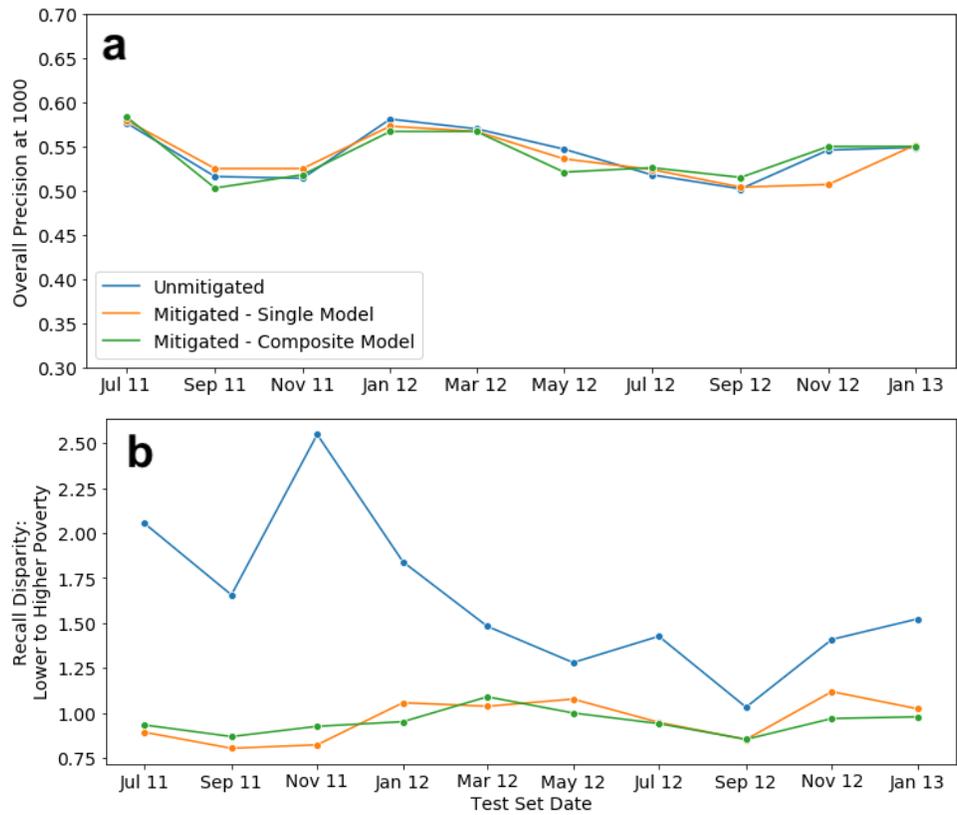
As it can be seen from the above equation, $1 + \lambda_g/P[G = g, Y = 1]$ is simply a group-specific constant with no dependence on the value of X beyond group membership. As such, this will result in scaling factor corresponding to rescaling of an underlying predicted score $P(Y = 1 | X)$ around a threshold, but without reordering the score within a given subgroup itself. Further, because of the monotonically increasing nature of recall as list depth is traversed, there will be a unique (up to exact ties in the score) recall-balancing solution of a given total list size, regardless of whether that solution is obtained by setting group-specific thresholds (as in the method we make use of in the current study) or group-specific stretching scores around a fixed threshold (as we see here), so long as the within-group ranking is not reordered.



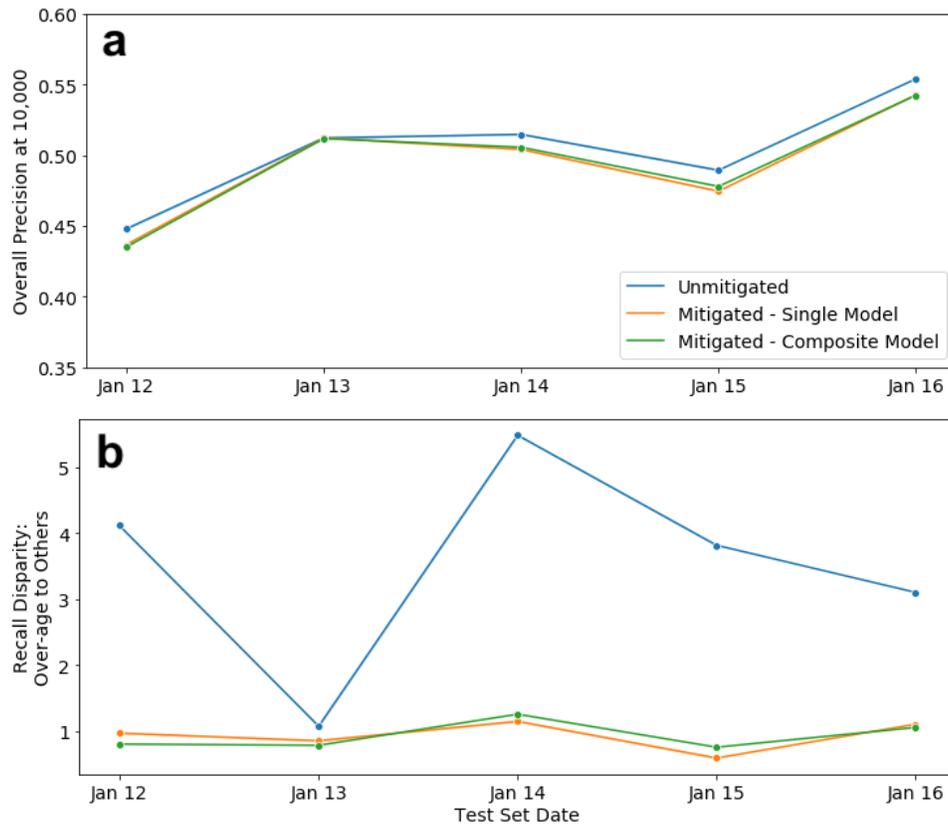
Supplementary Figure 1: Comparing equity (recall disparity) and performance (precision at top k) metrics for different model selection strategies between different policy contexts, as in Fig. 2, including an additional model selection strategy (Mitigated - Unadj. Model Seln.) in which model specification is chosen without regard for disparities then group-specific thresholds are chosen for the selected model. Error bars reflect 95% confidence intervals across temporal validation cohorts. This strategy performed similarly to strategy accounting for disparities in the model selection itself in all policy settings considered.



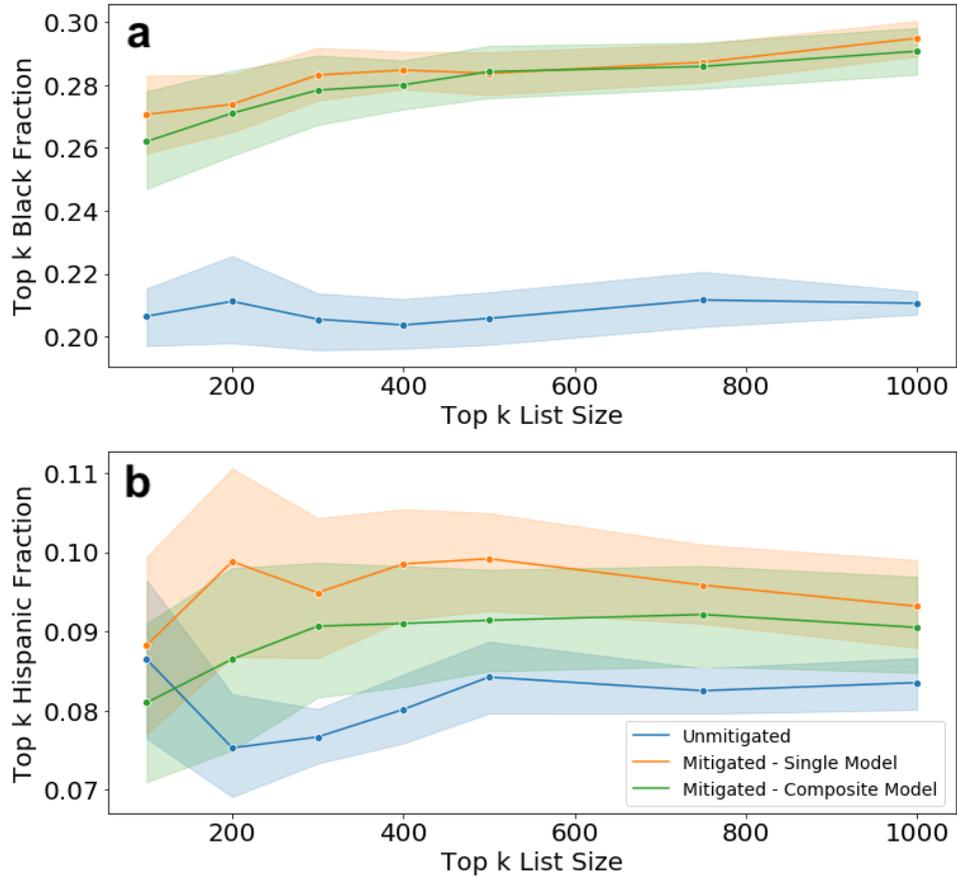
Supplementary Figure 2: Comparing model precision (a) and disparity (b) metrics over time for different model selection strategies in the Housing Safety policy context.



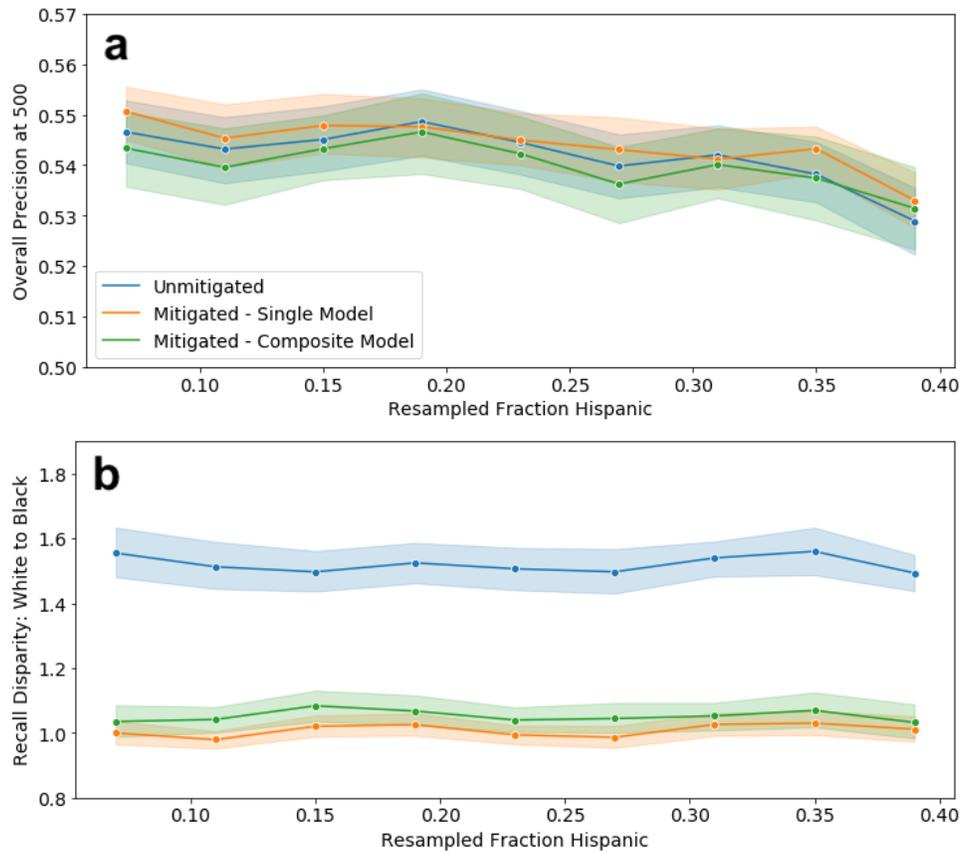
Supplementary Figure 3: Comparing model precision (a) and disparity (b) metrics over time for different model selection strategies in the Education Crowdfunding policy context.



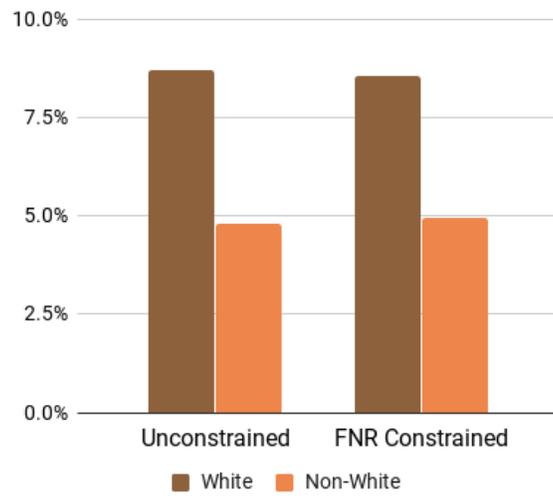
Supplementary Figure 4: Comparing model precision (a) and disparity (b) metrics over time for different model selection strategies in the Student Outcomes policy context.



Supplementary Figure 5: Fraction of the selected “top-k” list that is African American (a) or Hispanic (b) by list size in the Inmate Mental Health setting (see Fig. 4a-c). Shaded intervals reflect variation across temporal validation splits.



Supplementary Figure 6: Precision (a) and white-to-Black disparities (b) by re-sampled Hispanic fraction from the experiment shown in Fig. 4d. Shaded intervals reflect variation across bootstrap samples and temporal validation splits.



Supplementary Figure 7: Comparison of recall disparity after selecting the top 500 individuals from unconstrained and False Negative Rate (FNR)-constrained models following the methods described in Zafar et al. Results from the Inmate Mental Health policy context.

Supplementary Table 1: Model Hyperparameter Grid for Inmate Mental Health Policy Setting

Estimator	Hyperparameter	Grid Search Values
RandomForestClassifier	max_features	sqrt
	criterion	gini, entropy
	n_estimators	100, 1000, 5000
	min_samples_split	10, 25, 100
	max_depth	5, 10, 50
ScaledLogisticRegression	penalty	l1, l2
	C	0.001, 0.1, 1, 10
PercentileRankOneFeature	feature	Jail bookings last 1 or 5 years
SimpleThresholder	rules	Any mental health history, Jail releases in last: 1, 3, 6 months, or 1 year
	logical_operator	and

Supplementary Table 2: Model Hyperparameter Grid for Education Crowdfunding Policy Setting

Estimator	Hyperparameter	Grid Search Values
DecisionTreeClassifier	criterion	gini
	min_samples_split	2, 5, 10, 100, 1000
	max_depth	1, 5, 10, 20, 50, 100
RandomForestClassifier	max_features	sqrt
	criterion	gini, entropy
	n_estimators	100, 1000, 2000, 3000
	min_samples_split	10, 50
	max_depth	10, 50, 100
ExtraTreesClassifier	max_features	log2
	criterion	entropy
	n_estimators	10, 50, 1000
	min_samples_split	5, 25, 50
	max_depth	20, 50, 100
ScaledLogisticRegression	penalty	l1, l2
	C	0.0001, 0.001, 0.01, 0.1, 1, 10
AdaBoostClassifier	n_estimators	500, 1000

Supplementary Table 3: Model Hyperparameter Grid for Housing Safety Policy Setting

Estimator	Hyperparameter	Grid Search Values
DecisionTreeClassifier	criterion	gini, entropy
	min_samples_split	10, 20, 50, 100
	max_depth	1, 2, 3, 5, 10, 20, 50
RandomForestClassifier	max_features	sqrt, log2
	criterion	gini, entropy
	n_estimators	100, 1000, 5000
	min_samples_split	10, 20, 50, 100
	max_depth	2, 5, 10, 20, 50, 100
ExtraTreesClassifier	max_features	sqrt, log2
	criterion	gini, entropy
	n_estimators	100, 1000, 5000
	min_samples_split	10, 20, 50, 100
	max_depth	2, 5, 10, 50, 100
ScaledLogisticRegression	penalty	l1, l2
	C	0.001, 0.01, 0.1, 1, 10
PercentileRankOneFeature	feature	Days since last routine inspection

Supplementary Table 4: Model Hyperparameter Grid for Student Outcomes Policy Setting

Estimator	Hyperparameter	Grid Search Values
DecisionTreeClassifier	criterion	gini
	min_samples_split	2, 5, 10, 100, 1000
	max_depth	1, 5, 10, 20, 50, 100
RandomForestClassifier	max_features	sqrt
	criterion	gini
	n_estimators	100, 500
	min_samples_split	2, 10
	class_weight	~, balanced subsample, balanced
	max_depth	5, 50, 100
ExtraTreesClassifier	max_features	log2
	criterion	gini
	n_estimators	100
	max_depth	5, 50
ScaledLogisticRegression	penalty	l1, l2
	C	0.0001, 0.001, 0.01, 0.1, 1, 10

Supplementary References

- ¹ Celis, L. E., Huang, L., Keswani, V. & Vishnoi, N. K. Classification with fairness constraints: a meta-algorithm with provable guarantees. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2019), pp. 319–328.
- ² Zafar, M. B., Valera, I., Rodriguez, M. G. & Gummadi, K. P. Fairness beyond disparate treatment and disparate impact: learning classification without disparate mistreatment. *26th International World Wide Web Conference* (WWW, 2017), pp. 1171–1180.
- ³ Dwork, C., Immorlica, N., Kalai, A. T. & Leiserson, M. Decoupled classifiers for group-fair and efficient machine learning. *Proceedings of Machine Learning Research* **81**, 119–133 (2018).
- ⁴ Hardt, M., Price, E. & Srebro, N. Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems 29* (NIPS, 2016), pp. 3315–3323.
- ⁵ Rodolfa, K. T., et al. Case study: predictive fairness to reduce misdemeanor recidivism through social service interventions. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (ACM, 2020), pp. 142–153.
- ⁶ Krishna Menon, A. & Williamson, R. C. The cost of fairness in binary classification. *Proceedings of Machine Learning Research* (PMLR, 2018), pp. 1–12.