

Notes on “Delayed Impact of Fair Machine Learning”^[1]

Dhruv Medarametla*
dhruvm2@stanford.edu

Anav Sood*
anavsood@stanford.edu

December 8, 2018

1 Problem Setting

The paper we study considers two *groups*, A and B, which make up $g_A \in [0, 1]$ and $g_B := 1 - g_A$ proportion of the population. Throughout our review, we consider A to be disadvantaged compared to B. Each individual has a *score* $\mathcal{X} := [C]$, an abstract quantity summarizing how well an individual is suited to being *accepted* (or not *rejected*) by some *institution*. For group $j \in \{A, B\}$, scores are distributed according to $\pi_j \in \text{Simplex}^{C-1}$. For each group, the institution adopts a *policy* $\tau_j : \mathcal{X} \rightarrow [0, 1]$, where $\tau_j(x)$ represents the probability of accepting a person in group j with score x . The policy τ_j can be seen as a vector in $[0, 1]^C$ with $\tau_j(x)$ as the x th component. The institution then adopts an overall policy $\tau = (\tau_A, \tau_B)$ while attempting to maximize its *utility*

$$\mathcal{U}(\tau) := \sum_{j \in \{A, B\}} g_j \sum_{x \in \mathcal{X}} \tau_j(x) \pi_j(x) \mathbf{u}(x)$$

where the authors have assumed the existence of such a reasonable $\mathbf{u} : \mathcal{X} \rightarrow \mathbb{R}$ which models the utility contribution of a individual with score x . We also let $\mathcal{U}_j(\tau) := \sum_{x \in \mathcal{X}} \tau_j(x) \pi_j(x) \mathbf{u}(x)$.

The crux of the paper’s interest is in the socially motivated *outcome* for groups A and B. Namely, the authors assume existence of $\Delta : \mathcal{X} \rightarrow \mathbb{R}$, where $\Delta(x)$ represents the abstract change of well-being for an individual of score x should they be accepted. They then define change in the mean outcome of group j due to a policy τ

$$\Delta \mu_j(\tau) := \sum_{x \in \mathcal{X}} \pi_j(x) \tau_j(x) \Delta(x)$$

Note that this metric of social impact does not take into account those who have not been accepted; it is possible to modify $\Delta(x)$ in order to do so.

Lastly, the authors define a *success* function $\rho : \mathcal{X} \rightarrow [0, 1]$ where $\rho(x)$ is the probability that an individual with score x succeeds given that they were accepted. It is often intuitive and useful to define $\Delta(x)$ and $\mathbf{u}(x)$ using ρ . Also, note that in an initial effort to maintain fairness, $\mathbf{u}(x)$, $\Delta(x)$, and $\rho(x)$, all take values independent of the group the individual is from.

1.1 Examples

Middle Level Management. Consider a large tech company looking to promote current employees to middle management. Scores $x \in [C]$ encode a current employee’s managerial ability. The company decides to promote or not promote each individual according to their policy τ . For an employee with score x , company and personal utility are naturally dependent on $\rho : x \mapsto \rho(x) \in [0, 1]$, the probability the employee will be a successful manager. Given $\rho(x)$, we could define $\mathbf{u}(x) = 50\rho(x) - 150(1 - \rho(x))$, representing the idea that a successful manager increases the company’s utility by 50, while an unsuccessful one decreases it by 150. The individual impact of acceptance can be directly represented by change in income, in which case we can define $\Delta(x) = 25000$, representing promotion incurring a raise of 25000.

Now let the two groups of people be A, women, and B, men, with $g_A = 0.2$ and $g_B = 0.8$. Note that $\mathbb{E}_A[x]$ might be slightly less than $\mathbb{E}_B[x]$ due to years of unconscious bias and societal pushback, something to consider when determining a fair τ .

American Immigration Policies. Consider America as an institution deciding to accept or reject immigrants. The USA may accept immigrants according to a score $x \in [C]$, encoding an “American preparedness” measure. Immigrants may succeed (finding a job above a certain income and paying taxes) with some probability $\rho(x)$. The country may experience a utility $u(x)$ corresponding to a positive or negative change in GDP depending on the immigrants success. We can consider $\Delta(x)$ to be an immigrant’s change in income minus a fixed positive factor accounting for the turmoil and uncertainty of emigration. In this scenario, let the two groups of people be A, refugees, and B, non-refugees, with $g_A = 0.15$ and $g_B = 0.85$. Then, this model effectively represents the immigration policy when selecting for new residents among these two groups; most likely, $\mathbb{E}_A[x] < \mathbb{E}_B[x]$ due to the situations that refugee immigrants come from, something to consider when determining a fair τ .

1.2 Possible Constraints

In an attempt to force the chosen policy $\tau \in [0, 1]^{2C}$ to be more fair to the disadvantaged group, humans have historically constrained τ by demanding $\tau \in \mathcal{C} \subset [0, 1]^{2C}$. Below, we present the examples studied in the paper.

1. **MaxUtil**: The **MaxUtil** policy simply corresponds to $\mathcal{C} := [0, 1]^{2C}$. The only imposed fairness constraint is that $u(x)$ is independent of group.
2. **DemParity**: The **DemParity** policy demands that an equal proportion of each group be accepted by the institution. Namely, $\mathcal{C} := \{(\tau_A, \tau_B) : \sum_{x \in \mathcal{X}} \pi_A(x) \tau_A(x) = \sum_{x \in \mathcal{X}} \pi_B(x) \tau_B(x)\}$
3. **EqOpt**: The **EqOpt** policy demands that an equal proportion of people who are likely to be successful must be chosen from each group. Namely, if we define $\text{TPR}_j(\tau) := \frac{\sum_{x \in \mathcal{X}} \pi_j(x) \rho(x) \tau(x)}{\sum_{x \in \mathcal{X}} \pi_j(x) \rho(x)}$, then $\mathcal{C} := \{(\tau_A, \tau_B) : \text{TPR}_A(\tau_A) = \text{TPR}_B(\tau_B)\}$. To gain a little more intuition for this expression, we note that $\text{TPR}_A(\tau)$ is equal to $\mathcal{P}_{a \in A}[a \text{ selected} | a \text{ successful}]$, and correspondingly for B.

Considering the policy τ^{MaxUtil} achieved as the optimal solution to the **MaxUtil** constraint, for a policy τ the authors say that τ causes *relative harm* to group j if $\Delta\mu_j(\tau_j) < \Delta\mu_j(\tau^{\text{MaxUtil}})$ and *relative improvement* if $\Delta\mu_j(\tau_j) > \Delta\mu_j(\tau^{\text{MaxUtil}})$. In the same vein, the authors say that τ causes *active harm* if $\Delta\mu_j(\tau_j) < 0$, *improvement* if $\Delta\mu_j(\tau_j) > 0$, and *stagnation* if $\Delta\mu_j(\tau_j) = 0$.

2 Reduction to Threshold Policies

One of the paper’s large technical insights was a methodology for reducing the high parameter problem of picking a policy τ to a one parameter problem given by our *selection rate function* $r_{\tau_j}(\tau_j) := \sum_{x \in \mathcal{X}} \pi_j(x) \tau_j(x)$. Note when picking a policy τ_j for group j, the institution only cares about the values τ_j takes on scores for which π_j is non-zero, inducing an equivalence under π_j . Two equivalent policies provide identical utility for the institution, outcome (so our analysis of a policy is the same up to equivalence), true positive rate, and selection rate.

Now, as a thought experiment, let’s consider picking a policy under the **MaxUtil** constraint. It is easy to see from $\mathcal{U}(\tau)$ that an optimal policy would accept any individual with score $u(x) > 0$. Then, the reasonable assumption that $u(x)$ is strictly increasing in x would imply that anyone with score above some threshold score should be accepted. Under a **DemParity** constraint, we must select an equal proportion β of each group. Again if $u(x)$ is strictly increasing in x , it makes sense to select the β proportion of individuals in each population with the highest scores, yielding a similar threshold. This motivates the definition of a *threshold policy*, a policy $\tau_{c,\gamma}$ of the form

$$\tau_{c,\gamma} = \begin{cases} 1, & \text{if } x > c \\ \gamma, & \text{if } x = c \\ 0, & \text{if } x < c \end{cases}$$

for some $c \in \mathcal{X}$ and $\gamma \in [0, 1]$. The authors were able to then show these intuitive, but important results:

Result One. (Lemma 5.1) Let τ_j and τ'_j be threshold policies. Then $\tau_j \simeq_{\pi_j} \tau'_j$ if and only if $r_{\pi_j}(\tau_j) = r_{\pi_j}(\tau'_j)$. Further, $r_{\pi_j}(\tau_j)$ is a bijection from $\mathcal{T}_{\text{thresh}}(\pi_j)$ to $[0, 1]$ where $\mathcal{T}_{\text{thresh}}(\pi_j)$ is the set of equivalences classes between threshold policies under \simeq_{π_j} .

Result Two. (Proposition 5.2) Suppose $u(x)$ is strictly increasing in x . Then all optimal **MaxUtil** and **DemParity** policies (τ_A, τ_B) satisfy $\tau_j \simeq_{\pi_j} r_{\pi_j}^{-1}(r_{\pi_j}(\tau_j))$ for $j \in \{A, B\}$. If in addition $u(x)/\rho(x)$ is increasing, the same is true for all optimal **EqOpt** policies.

Thus, all optimal policies are effectively threshold policies. Combining these results, the authors have reduced from optimizing over all policies to optimizing over equivalences classes of threshold policies, which can be uniquely characterized by their selection rate.

3 Concavity of Utility Functions

The authors have re-parameterized the problem from $2C$ parameters to 2: the proportion of each group the institution accepts. The authors now show the optimal parameters β_j for utility are unique under the **MaxUtil** constraint, the **DemParity** constraint, the **EqOpt** constraint, and also for the social outcome function by proving that all of these functions are concave in β_j or a similar parameter. The first important result is the following:

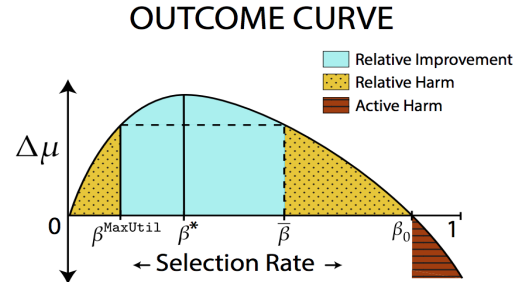
Result Three (Proposition 5.3). Let $\pi \in \text{Simplex}^{C-1}$ be a distribution over C states. Then if $w(x)$ is any non-decreasing map from $\mathcal{X} \rightarrow \mathbb{R}$, $\beta \rightarrow \langle w, \pi \circ r_{\pi}^{-1}(\beta) \rangle$ is concave.

Setting $w(x) = \Delta(x)$, we get that the function $\beta_A \rightarrow \Delta\mu_A(r_{\pi_A}^{-1}(\beta_A))$ [social outcome function] is concave. Setting $w(x) = u(x)$, we get that $\beta_j \rightarrow \mathcal{U}_j(r_{\pi_j}^{-1}(\beta_j))$ for a fixed j [**MaxUtil** function] is concave, and $\beta \rightarrow g_A \mathcal{U}_A(r_{\pi_A}^{-1}(\beta)) + g_B \mathcal{U}_B(r_{\pi_B}^{-1}(\beta)) = \mathcal{U}((r_{\pi_A}^{-1}(\beta), r_{\pi_B}^{-1}(\beta)))$ [**DemParity** function] is concave as well using properties of concave functions. The authors then prove a similar result for the **EqOpt** constraint.

Result Four (Lemma 6.1, B.1). For $j \in \{A, B\}$, the function $T_j(\beta) = \langle r_{\pi_j}^{-1}(\beta), \pi_j \circ u \rangle$ is a bijection from $[0, 1]$ to $[0, \langle \pi_j \circ u \rangle]$. Additionally, the function $t \rightarrow \mathcal{U}_j(r_{\pi_j}^{-1}(T_j^{-1}(t)))$ is concave in t .

Recall that the **EqOpt** policy forces $\text{TPR}_A(\tau_A) = \text{TPR}_B(\tau_B)$; this is equivalent to having $T_A(\beta_A) = T_B(\beta_B)$. Then, the first part of the result 4 implies that it is enough to prove that $t \rightarrow g_A \mathcal{U}_A(r_{\pi_A}^{-1}(T_A^{-1}(t))) + g_B \mathcal{U}_B(r_{\pi_B}^{-1}(T_B^{-1}(t))) = \mathcal{U}((r_{\pi_A}^{-1}(T_A^{-1}(t)), r_{\pi_B}^{-1}(T_B^{-1}(t))))$ [**EqOpt** function] is concave, which follows from the second part of result 4. The authors have thus shown that the social outcome function and the three constrained utility functions are concave, either as a function of β_j or a bijection to β_j . Since each function is concave, we have the existence of a unique interval of optimal β_j . Thus we can compare each constraint's social impact by examining $\Delta\mu_A(\beta_A)$ for β_A in the constraint's respective optimal interval.

Additionally, we can better understand our social outcome function. Under the reasonable assumption that $u(x) > 0 \implies \Delta(x) > 0$ (the institution takes on more risk than the individual) it is clear that $\beta^{\text{MaxUtil}} \leq \beta^*$. Then, we can define $\bar{\beta} \geq \beta^*$ such that $\mathcal{U}_A(r_{\pi_A}^{-1}(\bar{\beta})) = \mathcal{U}_A(r_{\pi_A}^{-1}(\beta^{\text{MaxUtil}}))$, and $\beta_0 \geq \bar{\beta}$ such that $\mathcal{U}_A(r_{\pi_A}^{-1}(\beta_0)) = 0$. The image on the right (from the paper) summarizes this paragraph pictorially by considering a graph of the social outcome function.



4 Main Results

In Section 6, the authors provide two theorems which give a robust characterization where the optimal selection rate parameter β lies for group A under **DemParity** and **EqOpt** constraints. These Theorems imply the

Corollaries in Section 3, which provide important insight into how such social constraints impact group A. Letting $G^{(A \rightarrow B)}(\beta_A) := T_B^{-1}(T_A(\beta_A))$ be the transfer function as described in the paper, we have that

Result Five (Corollary 3.2). (a) Under the assumption that $\beta_A^{\text{MaxUtil}} < \bar{\beta}$ and $\beta_B^{\text{MaxUtil}} > \beta_A^{\text{MaxUtil}}$, there exist population proportions $g_0 < g_1 < 1$ such that, for all $g_A \in [g_0, g_1]$, $\beta_A^{\text{MaxUtil}} < \beta_A^{\text{DemParity}} < \bar{\beta}$. That is, **DemParity** causes relative improvement.

(b) Under the assumption that there exist $\beta_A^{\text{MaxUtil}} < \beta < \beta' < \bar{\beta}$ such that $\beta_B^{\text{MaxUtil}} > G^{(A \rightarrow B)}(\beta), G^{(A \rightarrow B)}(\beta')$, there exist population proportions $g_2 < g_3 < 1$ such that, for all $g_A \in [g_2, g_3]$, $\beta_A^{\text{MaxUtil}} < \beta_A^{\text{EqOpt}} < \bar{\beta}$. That is, **EqOpt** causes relative improvement.

Result Six (Corollary 3.3, 3.4). (a) Fix a selection rate β . Assume that $\beta_B^{\text{MaxUtil}} > \beta > \beta_A^{\text{MaxUtil}}$. Then there exists a population proportion g_0 such that, for all $g_A \in [0, g_0]$, $\beta_A^{\text{DemParity}} > \beta$. In particular, when $\beta = \beta_0$, **DemParity** causes active harm, and when $\beta = \bar{\beta}$, **DemParity** causes relative harm.

(b) Suppose that $\beta_B^{\text{MaxUtil}} > G^{(A \rightarrow B)}(\beta)$ and $\beta > \beta_A^{\text{MaxUtil}}$. Then there exists a population proportion g_0 such that, for all $g_A \in [0, g_0]$, $\beta_A^{\text{EqOpt}} > \beta$. In particular, when $\beta = \beta_0$, **EqOpt** causes active harm, and when $\beta = \bar{\beta}$, **EqOpt** causes relative harm.

Corollary 3.5 shows that **DemParity** can cause active harm while **EqOpt** causes improvement. We see this is intuitive as, under **DemParity**, if scores of group B are particularly high in comparison to low scores in group A, maximizing utility still gives large β due to the massive increase in utility given by accepting much of group B. The **DemParity** constraint then forces many members of group A with negative $\Delta(x)$ to also be accepted. The **EqOpt** constraint avoids this by not forcing equal acceptance rates. We instead present and examine Corollary 3.6, which is much more non-intuitive.

Result Seven (Corollary 3.6) Suppose it is the case that $\beta_A^{\text{MaxUtil}} < \beta_B^{\text{MaxUtil}}$ and $\text{TPR}_A(\tau^{\text{MaxUtil}}) > \text{TPR}_B(\tau^{\text{MaxUtil}})$. Then $\beta_A^{\text{EqOpt}} < \beta_A^{\text{MaxUtil}} < \beta_A^{\text{DemParity}}$. That is, **EqOpt** causes relative harm.

The above Corollary shows it is possible for **EqOpt** to select less of group A than **MaxUtil**, meaning **EqOpt** causes relative harm. This is a surprising result, and has interesting social implications surrounding how inequality in a disadvantaged group can reverse the intended effect of the **EqOpt** constraint. See Application 2 below for a concrete example of this result.

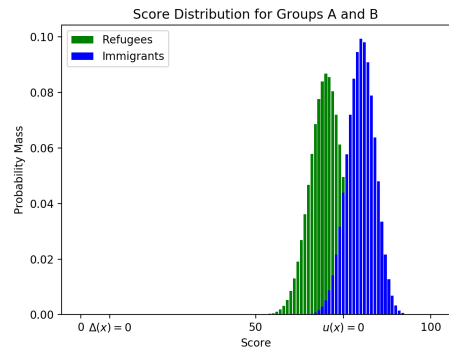
5 Applications

We go over two case studies that study how $\beta_j^{\text{MaxUtil}}, \beta_j^{\text{DemParity}}, \beta_j^{\text{EqOpt}}$, and β_j^* differ in concrete examples.

5.1 American Immigration Policies: An Example of Corollary 3.2

Recall again our example of immigrants coming to America. Say that every immigrant can be assigned a score $x \in \mathcal{X}$, where $\mathcal{X} = \{0, 1, \dots, 99, 100\}$, that measures how prepared they are to become a citizen of America. Assume that we can calculate $\rho(x) = 0.35 + 0.006x$, $u(x) = 20\rho(x) - 80(1 - \rho(x))$, and $\Delta(x) = 30\rho(x) - 20(1 - \rho(x))$.

In other words, we are saying that America gains utility of 20 for every successful immigrant and loses utility of 80 for every unsuccessful immigrant, and that an immigrant gains 30 utility if they are successful and loses 20 utility if they are unsuccessful in America (these parameters are estimated quantities; the true values may differ significantly).



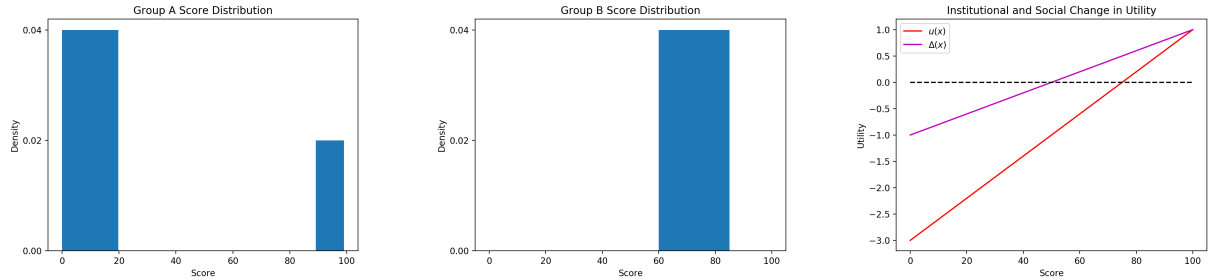
Now, say we have two subpopulations: A, refugees and B, non-refugees, with $g_A = 0.15$ and $g_B = 0.85$. Say that the scores of the refugees are distributed according to $\mathcal{B}(100, 0.7)$, where $\mathcal{B}(n, p)$ is the discrete binomial distribution with. Similarly, say that the scores of the non-refugees are distributed according to $\mathcal{B}(100, 0.8)$. We display this in the figure above, marking where $\mathbf{u}(x)$ and $\Delta(x)$ are 0. Given these set of parameters, we can calculate the optimal β_j for each policy, as well as institutional utilization and social outcome. Our results are summarized below:

Application 1: American Immigration Policies				
Parameter	MaxUtil	DemParity	EqOpt	Social Max
β_A	0.163	0.811	0.809	> 0.999
β_B	0.913	0.811	0.811	> 0.999
$\mathcal{U}(\tau)$	2.690	2.381	2.392	2.100
$\Delta\mu_A(\tau_A)$	3.35	15.38	15.35	18.50

In this case, we can see that group B dominates due to the high value of g_B ; as a result, β_B is relatively similar for the three constraints, but β_A varies noticeably. This application displays that for many reasonable modeling scenarios, parts (a) and (b) of Corollary 3.2 are at work, as both DemParity and EqOpt cause relative improvement for group A.

5.2 Application 2: Understanding Corollary 3.6

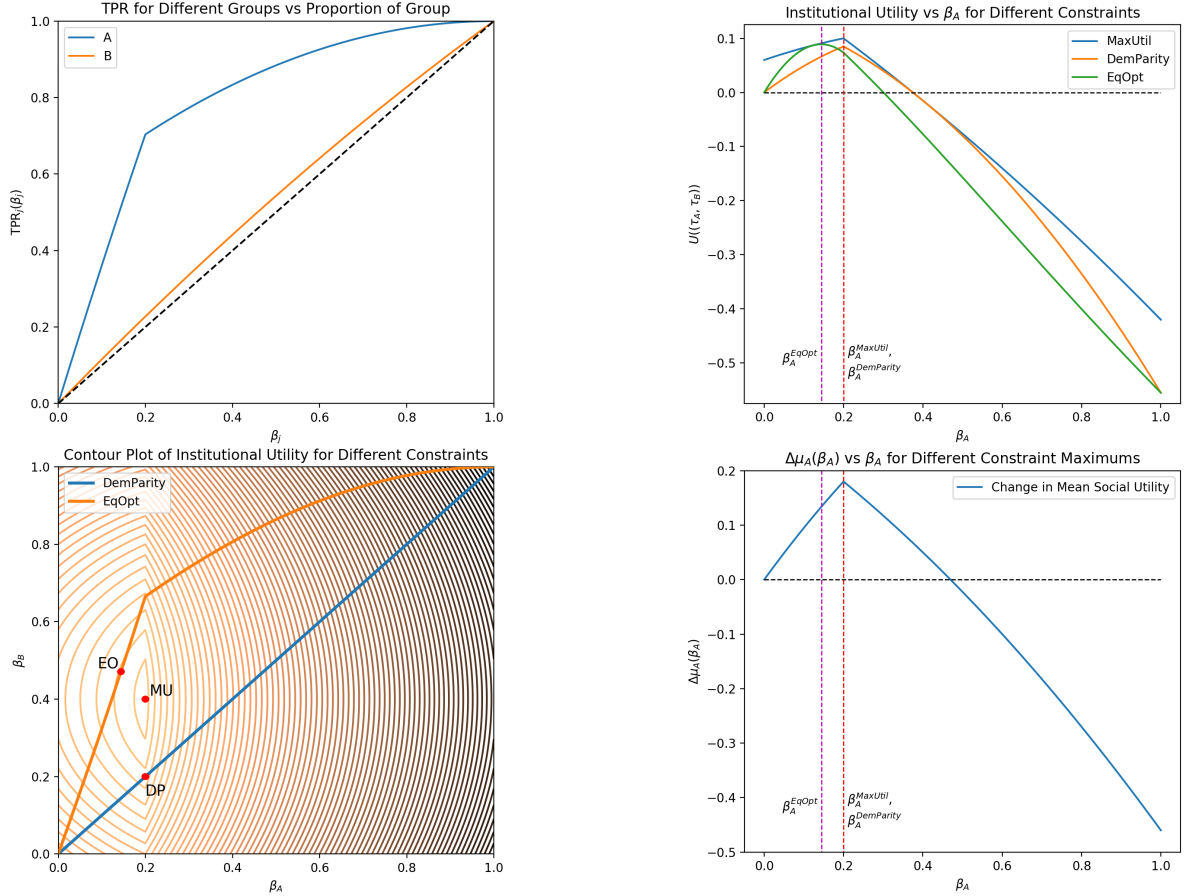
Consider the following abstract example, which is set in the natural continuous extension of the problem statement.¹ We let $g_A = 0.25$, $g_B = 0.75$. We let $\mathcal{X} = [0, 100]$ with the scores of A distributed with 80% of data uniform in $[0, 20]$ and 20% of data uniform in $[90, 100]$ and B distributed with all data uniformly in $[60, 85]$. We let $\rho(x) = x/100$, $\mathbf{u}(x) = \rho(x) - 3(1 - \rho(x)) = -3 + x/25$, and $\Delta(x) = \rho(x) - (1 - \rho(x)) = -1 + x/50$. We provide both plots visually depicting this and tabulated results below:



Application 2: High Inequality in Group A				
Parameter	MaxUtil	DemParity	EqOpt	Social Max
β_A	0.2	0.2	0.144	0.2
β_B	0.4	0.2	0.471	1.0
$\mathcal{U}(\tau)$	0.1	0.085	0.089	-0.035
TPR _A	0.704	0.704	0.514	0.704
TPR _B	0.441	0.228	0.514	1.0
$\Delta\mu_A(\tau_A)$	0.18	0.18	0.133	0.18

Due to inequality in group A, a large amount of group A's potential success is concentrated among high-achievers. The success of group B, however, is more evenly distributed. Thus, under the constraint that $\text{TPR}_A = \text{TPR}_B$, even a small β_A naturally results in a large β_B , causing the policy to accept individuals from group B who cause a decrease in \mathcal{U} . Thus we see $\beta_A^{\text{EqOpt}} < \beta_A^{\text{MaxUtil}}$, meaning that EqOpt is causing relative harm to group A. This can be very clearly visually seen in our plots below:

¹Treating scores and such continuously lends for easier computation, while still allowing us to explain the impact and importance of Corollary 3.6. The intuition transfers to the discrete case directly.



6 Limitations and Future Work

We first discuss some possibilities for future work.

1. Consider Corollary 3.6. In the scenario outlined by our second example, can we prevent **EqOpt** from relatively harming group A by selecting $\beta'_B = \min\{(1 + \epsilon)\beta_A, \beta_B\}$ for group B?
2. There is some study of an outcome-based alternative, where the institution solves the optimization problem $\max_{\tau_A} \Delta_A \mu_A(\tau_A)$ s.t. $\mathcal{U}^{\text{MaxUtil}} - \mathcal{U}(\tau) < \delta$. Essentially, the institution prioritizes the outcome of group A so long as they stay within a δ budget. We wonder if $\max_{\tau} \mathcal{U}(\tau) + \lambda \Delta \mu_A(\tau_A)$ is a more natural way to phrase the problem (in fact we believe the best way), where we have given the institution incentive (public image, tax returns etc.) to help the disadvantaged group. How similar are these two problem statements? Do they offer differing results?
3. The topic of non-selection is only touched on in one footnote. The authors let $\Delta_a(x)$ and $\Delta_n(x)$ represent the effect of being accepted and rejected respectively, then let $\Delta \mu_j(\tau) := \sum_{x \in \mathcal{X}} \pi_j(x) (\tau_j(x) \Delta_a(x) + (1 - \tau_j(x)) \Delta_n(x))$. Our above results hold if $\Delta_a(x) - \Delta_n(x)$ is increasing in x ; however, this is a rather nuanced assumption. For example, consider the case of immigrants, where a low-scoring immigrant will likely have a larger value for $\Delta_a(x) - \Delta_n(x)$ due to a very negative $\Delta_n(x)$. We wonder what can be said about non-selection without this assumption.

With regard to the work itself, it is a rather complete study of the problem setup. That being said, it applies to one epoch, or one cycle of being accepted and denied. Helping disadvantaged groups requires the proper constraints to be employed over long periods of time and monitored throughout multiple epochs. This likely requires a more restrictive and focused study per particular issue. Additionally it requires an expert to be incredibly confident in setting u , Δ , and ρ functions, which offers a lot of modeling freedom, much of which we exploited in making our examples. Currently, at best, the results provide very solid intuition to guide more focused and less general work for specific issues.

References

- [1] Liu, L. T., Dean, S., Rolf, E., Simchowitz, M., and Hardt, M. (2018). Delayed impact of fair machine learning. In International Conference on Machine Learning.
- [2] Bertrand, M., Chugh, D., and Mullainathan, S. (2005). Implicit discrimination. *The American Economic Review*, 95(2):94–98.