# A Dynamic Decision-Making Framework Promoting Long-Term Fairness

Bhagyashree Puranik bpuranik@ucsb.edu University of California Santa Barbara Upamanyu Madhow madhow@ucsb.edu University of California Santa Barbara Ramtin Pedarsani ramtin@ucsb.edu University of California Santa Barbara

# **ABSTRACT**

With AI-based decisions playing an increasingly consequential role in our society, for example, in our financial and criminal justice systems, there is a great deal of interest in designing algorithms conforming to application-specific notions of fairness. In this work, we ask a complementary question: can AI-based decisions be designed to dynamically influence the evolution of fairness in our society over the long term? To explore this question, we propose a framework for sequential decision-making aimed at dynamically influencing long-term societal fairness, illustrated via the problem of selecting applicants from a pool consisting of two groups, one of which is under-represented. We consider a dynamic model for the composition of the applicant pool, in which admission of more applicants from a group in a given selection round positively reinforces more candidates from the group to participate in future selection rounds. Under such a model, we show the efficacy of the proposed Fair-Greedy selection policy which systematically trades the sum of the scores of the selected applicants ("greedy") against the deviation of the proportion of selected applicants belonging to a given group from a target proportion ("fair"). In addition to experimenting on synthetic data, we adapt static real-world datasets on law school candidates and credit lending to simulate the dynamics of the composition of the applicant pool. We prove that the applicant pool composition converges to a target proportion set by the decision-maker when score distributions across the groups are identical.

## **CCS CONCEPTS**

• Computing methodologies  $\rightarrow$  Markov decision processes; Machine learning approaches.

# **KEYWORDS**

Long-term fairness, Fair selection, Positive reinforcement, Sequential decision-making, AI for social equity

### **ACM Reference Format:**

Bhagyashree Puranik, Upamanyu Madhow, and Ramtin Pedarsani. 2022. A Dynamic Decision-Making Framework Promoting Long-Term Fairness. In Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society (AIES'22), August 1–3, 2022, Oxford, United Kingdom., 10 pages. https://doi.org/10.1145/3514094.3534127



This work is licensed under a Creative Commons Attribution International 4.0 License.

AIES'22, August 1–3, 2022, Oxford, United Kingdom © 2022 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-9247-1/22/08. https://doi.org/10.1145/3514094.3534127

# 1 INTRODUCTION

In this paper, we seek to develop a framework for sequential decision making aimed at influencing long-term societal fairness. Machine learning models are being increasingly applied in making critical decisions that affect humans, such as recidivism prediction [8], mortgage lending [3], and recommendation systems [25]. While the algorithms offer increased efficiency, speed, and scalability in the decision-making process, they could introduce bias leading to the decisions being unfair towards certain groups of the population. There is a rich and rapidly growing literature on "fair" strategies that mitigate bias in algorithmic decision making, including pre-processing the labels or data and reweighting costs based on groups [16], adversarial de-biasing [28], introducing regularizers based on mutual information [17], addition of constraints that satisfy fairness criteria [26], learning representations that obfuscate group information [27] and many more.

Most of the above studies focus on a static framework where the long-term effects of decisions on the population are not explored. However, in many practical applications, decisions may affect the feature distributions across groups and influence the future rewards, that will eventually affect the dynamics of the decision-making loop [19].

The long-term dynamic study of such systems can be modeled through a reinforcement learning framework based on Markov Decision Process (MDP) as considered in our work. Our framework is motivated by real-world examples such as the following. Consider a company receiving applications every month, which wants to hire good candidates in an unbiased manner (e.g., by ultimately selecting equal numbers of male and female applicants). With the total monthly intake fixed based on a budget, the company selects a certain proportion of candidates from each group. The hiring decisions affect the subsequent pool of applicants: admitting more candidates from a particular group might encourage more such candidates to apply, or successful candidates from a group might inspire other such candidates, providing positive feedback into the decision-making loop. Such a strategy could not only enhance diversity and equity, but also enable the company to learn more about a minority group so as to eventually have a richer pool of well-qualified applicants. Another motivating example is college admissions, where the goal may be to admit students with the best academic records, while accounting for socio-economic background and reducing bias based on sensitive attributes such as race or gender. Could one, for example, reverse the trend in the decrease in the proportion of women in science, technology, engineering and mathematics (STEM) as documented in [5]? It reported that 18% of bachelor's degrees in computer science were awarded to women in 2010, down from 37% in 1985. Studies also point out that fewer

women choose to apply to such fields as result of societal influences. We suggest here a structured framework for fair selection aimed at combating such systemic imbalances by encouraging a larger number of people from minority groups to participate in the selection process.

Contributions. Based on a simple model for evolution of the composition of the applicant pool, we develop a framework for fair selection by formulating the problem as a Markov Decision Process (MDP) with two objectives - maximizing the utility by admitting candidates with the highest scores, and minimizing the disparity between the proportions of selected candidates from each group. We present two policies for fair selection: an optimal policy based on value-iteration that maximizes the utility accumulated over multiple rounds, where the utility comprises of a greedy term that maximizes sum of scores of selected applicants and a fair-only term that minimizes disparity; and second, a computationally simple and effective policy, which we term the Fair-Greedy (FG) policy, that optimizes for instantaneous utility. We characterize the structure of the FG policy and show convergence and also prove that the applicant pool proportion approaches the target proportion that is desired by the system under identical score distributions across the two groups. We provide experimental results on interesting scenarios with synthetic data, as well as with dynamic data created from the static law school [23] and German credit [9] datasets.

## 2 RELATED WORK

Recent work on fairness in sequential decision making includes settings such as online classification [2], Bayesian decision making [7] and predictive policing [10]. Several works address the notion of imposing fairness in multi-armed bandit and online learning problems [6, 12, 13, 15, 21]. This body of work focuses on the design of policies and the effects of fairness constraints on them. However, in these frameworks, decisions do not affect future samples.

The importance of introducing dynamics into notions of fairness is highlighted by studies indicating that static fairness criteria may lead to undesired long-term effects on minority groups [18], [29]. While we focus in this paper on the participation rates of different groups in the selection process, prior work on fairness in sequential decision making has focused, either explicitly or implicitly, on the impact of decisions on the qualifications or score distributions of the different groups.

In particular, [18] models the effect of fairness-aware decisions via a one-step feedback model: for example, they might model the mean change in credit score in a disadvantaged segment of the population as a function of the rate at which bank loans are granted. It is shown in [18] that, depending on the specific model for the change, "fair" policies (e.g., equalizing selection rates or true positive rates across disadvantaged and advantaged groups) may sometimes lead to negative outcomes. The work [29] studies how the imposition of hard fairness constraints leads to changes in the underlying feature distributions and the group representation. In particular, they show that imposing typical notions of fairness such as statistical parity or equality of opportunity could lead to exacerbation of the disparity between the group proportions of samples, and the disadvantaged group may even exit the system.

Modeling the long-term impact in the sense of the updated population distributions feeding into the subsequent examples seen by the system and studying such feedback effects have been traditionally investigated using reinforcement learning frameworks via Markov Decision Processes (MDPs), and introducing fairness constraints in the reward functions [11, 14, 22]. Departing from conventional statistical notions of fairness based on independence or separation, [14] adopts a 'weakly meritocratic' notion where they devise policies such that, their algorithm never (probabilistically) prefers an action over another, if the latter has larger long-term utility, which for example in a hiring process, can be viewed as the selction process cannot target one group over another if selection from either groups leads to similar long-term utility or benefit to the institution.

Recent works such as [20, 24, 30] examine the long-term impact of decisions on the features of the population. Building on the work of [18], the authors of [24] propose a dynamic model with the motivation of loan lending decisions. They model the group-wise distributions of the likelihood of loan repayment (analogous to score distributions in our framework), termed the payback probabilities, and consider dynamics governed by the hypothesis: granting loans produces upward mobility for a population when they are repaid. Along with examining the impact of fair decisions on the likelihood of loan repayment, they also highlight the detrimental effects of unequal misestimation of the payback probabilities across groups under their model, even under fair decisions. A fundamental notion of fairness is that of 'affirmative action', which is viewed in [20] as balancing the long-term qualification across groups. The authors in [20] study the evolution of qualification rates while attempting to maintain the social equity of selecting an equal number of applicants from both groups. They assume that the selection decisions could act as either an incentive or impediment, causing a change in the proficiency of a group: for example, systemic rejection of a particular group may cause the group's population to lose the interest to participate altogether. The long-term dynamics of group wise qualification rates are also investigated in [30]. Under a partially observable MDP setting, they introduce a myopic policy, characterize the equilibrium of dynamics and study their effects on population under two regimes: one where accepted individuals feel less motivated to remain qualified, and another where accepted individuals get access to better resources and hence remain or become more qualified.

We adopt an outlook complementary to the preceding body of work, seeking to influence the participation of under-represented groups in the selection process. We do not assume that the score distributions change as a consequence of our decisions, but our model can be extended to accommodate such changes, as long as we can estimate them. Rather than studying the impact of fair policies as in [20, 24, 30], we provide a generic framework for achieving long-term fairness dynamically. While we also consider a score-based selection problem as in [18], our notion of fairness is that the proportion of applicants and also that of admissions is equitable across groups or approaches a target set by the policy-maker. We adopt the MDP framework as well, but instead of imposing fairness as a hard static constraint at every round in the sequential decision-making process, we define our reward as a composition of

two-fold objectives of maximization of scores of accepted individuals and minimizing disparity between the proportion of accepted individuals from a target set by the decision-maker. We model the proportion of applicants as states of the MDP, thus the state space is different from that considered in other works.

#### 3 PROBLEM SETTING

Given that there are two groups u and v within the population, based on a binary valued sensitive attribute, we denote the total number of applicants in round t by  $N_t$ , out of which  $N_t^u$  belong to group u and  $N_t^v = N_t - N_t^u$  belong to group v. We wish to admit a fixed proportion  $\bar{a}$  of the total applicants, leading to  $A_t = \bar{a}N_t$  number of total applicants accepted in round t. We denote by  $A_t^u$  and  $A_t^v = A_t - A_t^u$  the number of applicants selected in round t from groups u and v respectively.

**Score distributions.** The qualification of an applicant is measured by the *score*, assumed to be an increasing function of the proficiency of a candidate. Let  $\mathcal{P}_u$  and  $\mathcal{P}_v$  denote the score distributions of the two groups. Thus the scores for groups u and v are  $\{X_i^u\}_{i=1}^{N_t^u}$  and  $\{X_j^v\}_{j=1}^{N_t^v}$ , generated from  $\mathcal{P}_u$  and  $\mathcal{P}_v$  respectively. We denote the ordered scores by  $\{X_{(i)}^u\}_{i=1}^{N_t^u}$  and  $\{X_{(j)}^v\}_{j=1}^{N_t^v}$ , where  $X_{(i)}^u$  and  $X_{(j)}^v$  denote the  $i^{th}$  and  $j^{th}$  largest scores out of  $N_t^u$  and  $N_t^v$  respectively.

**Fairness-aware utility**. The goal is to optimize the *utility*, which comprises of two parts: a greedy term (to be maximized) which is the expected sum of scores of selected candidates, and a fair term (to be minimized) measuring disparity between groups based on a *target* proportion.

**MDP formulation.** We define the MDP state  $s_t \in [0, 1]$  as the proportion of applicants from group u out of the total, and the action  $a_t \in [0, 1]$  as the proportion of selected candidates from group u out of the total selected candidates:

$$s_t = \frac{N_t^u}{N_t}, a_t = \frac{A_t^u}{A_t}.$$

We denote by  $\bar{s} \in (0,1)$  the long-term target of the proportion of group u among the selected applicants. For example, if group u is under-represented in the applicant pool, we may set  $\bar{s}$  as the proportion of group u in society at large. Instead, if our long-term goal is to admit equal number from both groups, we set  $\bar{s} = 0.5$ . Note that formulating the states and actions as proportions of group u is sufficient since the proportion of applicants and admitted candidates from group v is naturally  $1 - s_t$  and  $1 - a_t$  respectively. The overall utility or reward is:

$$R(s_t, a_t) = R_G(s_t, a_t) - \lambda L_{\mathcal{F}}(a_t), \tag{1}$$

where the *greedy* reward term is the expected sum of scores of admitted candidates, given by:

$$R_{\mathcal{G}}(s_{t}, a_{t}) = \frac{1}{A_{t}} \mathbb{E} \left[ \sum_{i=1}^{A_{t}^{u}} X_{(i)}^{u} + \sum_{i=1}^{A_{t}^{v}} X_{(i)}^{v} \right]$$
$$= \frac{1}{A_{t}} \mathbb{E} \left[ \sum_{i=1}^{a_{t}A_{t}} X_{(i)}^{u} + \sum_{i=1}^{(1-a_{t})A_{t}} X_{(i)}^{v} \right],$$

and the fairness loss term is

$$L_{\mathcal{F}}(a_t) = (a_t - \bar{s})^2. \tag{2}$$

Since the accepted candidates are the ones with the largest scores, the ordered statistics of the score distributions come into play. In (1),  $\lambda \geq 0$  is a parameter used to control the weight given to the fairness objective relative to the greedy objective. The greedy objective promotes the admission of good candidates, while the fairness objective promotes fairness in selection proportion. The fairness objective is balanced: it pushes the selection proportion towards  $\bar{s}$  regardless of whether group u is under-represented or over-represented among the selected applicants. Note that the dependence of the greedy reward on state  $s_t$  is through  $N_t^u$  and  $N_t^v$ , where the ordered scores for groups u and v are specifically out of  $N_t^u$  and  $N_t^v$  applicants for groups u and v respectively.

**Applicant pool evolution.** We illustrate our ideas with a simple linear model for positive reinforcement. The effect of decisions on subsequent applicant pools would, in reality, be far more complex; we hope that our work stimulates the major effort in experimentation and data collection required to build such models. We model the positive reinforcement provided by decision-making as a set of transition probabilities  $\mathcal{P}(s_{t+1}|s_t,a_t)$ . The total number of applicants  $N_t$  to the system at round t can be any sequence of numbers and the number of applicants from group u to the system is sampled from a Poisson distribution based on the mean parameter and overall number of applicants (which is variable) as

$$N_t^u \sim Pois(\theta_t N_t),$$
 (3)

where  $Pois(\cdot)$  is the Poisson distribution with mean  $\theta_t N_t$ . Thus,  $\theta_t$  is the mean proportion of group u in the applicant pool in round t. We consider the following model for positive reinforcement:

$$\theta_{t+1} = [\theta_t + \eta(a_t - s_t)]_C, \tag{4}$$

where  $\eta$  is a step-size parameter and  $[]_C$  is the projection on the convex set C = [0,1]. Thus the update is such that when the admission rate  $a_t$  of the group u is higher than the application rate  $s_t$ , more applicants from the group are encouraged in future rounds, and vice versa. The state then evolves as

$$s_{t+1} = \frac{N_{t+1}^u}{N_{t+1}}.$$

The model for positive reinforcement is relevant to many real-world selection systems and is inspired by the social behavior that the successful admission of candidates from a particular group encourages more such candidates to apply to the institution. For instance, a large number of female college graduates in society serve as role-models, encouraging the future generations of women to go to college. However, if a particular program is known for admitting women at a rate smaller than the application rate, lesser women might consider the institution as worth applying to.

**Optimal Policy.** The maximum long-term reward accumulated by the system through the horizon H is given by

$$\max_{\pi} \mathbb{E}\left[\sum_{t=0}^{H} R(s_t, a_t) | \pi\right] \tag{5}$$

where  $\pi$  is the policy or mapping from the set of states to the set of actions. The optimal policy  $\pi^*(s)$  can be found by exact methods

such as value iteration [4], where the optimal value function is defined as:

$$V^*(s) = \max_{\pi} \mathbb{E} \left[ \sum_{t=0}^{H} \gamma^t R(s_t, a_t) | \pi, s_0 = s \right], \tag{6}$$

which is the cumulative reward earned by playing policy  $\pi$ , and starting from initial state s, with  $0 < \gamma < 1$  being the discount factor. The optimal policy  $\pi^*(s)$  is found by iteratively solving the Bellman equation:

$$V_k^*(s) = \max_{a} \sum_{s'} \mathcal{P}(s'|s, a) [R(s, a) + \gamma V_{k-1}^*(s')], \forall s$$
 (7)

and the optimal policy is computed iteratively as below:

$$\pi_k^*(s) = \arg\max_{a} \sum_{s'} \mathcal{P}(s'|s, a) [R(s, a) + \gamma V_{k-1}^*(s')], \qquad (8)$$

until the optimal policy converges to  $\pi^*(s)$ . It is also known that the value iteration algorithm converges as long as the reward is bounded in magnitude [4]. However, analyzing the equilibrium state of the MDP under this optimal policy is intractable.

We observe through simulations that the structure of the optimal policy  $\pi^*(s)$  is similar to that of the simpler *Fair-Greedy* policy proposed next, and that the applicant pool evolution converges to an equilibrium point.

#### 4 FAIR-GREEDY POLICY

Finding an optimal policy is computationally expensive as the state space grows larger. We therefore propose a simple, yet effective, *Fair-Greedy* policy that optimizes the instantaneous overall utility in (1):

$$\pi_{FG}^*(s_t) = \arg\max_{a_t} R(s_t, a_t). \tag{9}$$

We provide insight into this policy by considering its performance for a large applicant pool ( $N_t$  large) with identical score distributions across the two groups. In this regime, we first prove that the greedy reward term is optimized when the admission proportion is the same as the applicant proportion. We then derive some key properties of the FG policy, and provide theoretical guarantees for the convergence of the applicant pool to the target proportion.

Theorem 4.1. If the score distributions  $\mathcal{P}_u$  and  $\mathcal{P}_v$  of the two groups are identical, the greedy reward  $R_{\mathcal{G}}(s_t, a_t)$  is optimized by the action:

$$a_{\mathcal{G}}^* = \arg\max_{a_t} R_{\mathcal{G}}(s_t, a_t) = s_t. \tag{10}$$

PROOF. Recall that the greedy reward is given by:

$$R_{\mathcal{G}}(s_t, a_t) = \frac{1}{A_t} \mathbb{E}\left[\sum_{i=1}^{A_t^u} X_{(i)}^u + \sum_{i=1}^{A_v^v} X_{(i)}^v\right]$$
(11)

Since we assume the space of actions as  $a_t \in [0, 1]$ , the number of admitted candidates from each group, more formally, are  $A_t^u = \lfloor a_t A_t \rfloor$  and  $A_t^v = \lfloor (1-a_t)A_t \rfloor$ . For simplicity of presentation, we omit the 'floor' without loss of generality of our results since we are interested in the regime that  $N_t$  is large. Therefore, we write:

$$R_{\mathcal{G}}(s_t, a_t) = a_t \mathbb{E}\left[\frac{\sum_{i=1}^{a_t A_t} X_{(i)}^u}{a_t A_t}\right] + (1 - a_t) \mathbb{E}\left[\frac{\sum_{i=1}^{(1 - a_t) A_t} X_{(i)}^v}{(1 - a_t) A_t}\right]$$

By the law of large numbers, the collection of score variables  $\{X_i^u\}_{i=1}^{N_t^u}$  and  $\{X_i^v\}_{i=1}^{N_t^v}$  converge to their respective distributions  $\mathcal{P}_u$  and  $\mathcal{P}_v$  as  $N_t$  increases. Choosing the top  $A_t^u = a_t A_t$  candidates out of  $N_t^u$  (similarly top  $A_t^v$  out of  $N_t^v$ ) is equivalent to setting a threshold  $t_u$  (similarly,  $t_v$ ) and admitting all candidates with scores above the threshold. This holds for generic score distributions and they need not necessarily be identical across the groups. Thus for large  $N_t$ , the average score of the admitted candidates from each group approaches its expected value as:

$$\lim_{N_t \to \infty} \frac{\sum_{i=1}^{a_t A_t} X_{(i)}^u}{a_t A_t} = \mathbb{E}[X^u | X^u \ge t_u]$$
 (12)

$$\lim_{N_t \to \infty} \frac{\sum_{i=1}^{(1-a_t)A_t} X_{(i)}^v}{(1-a_t)A_t} = \mathbb{E}[X^v | X^v \ge t_v]$$
 (13)

Rewriting the greedy reward in terms of the above conditional expectations leads to the following equation:

$$R_{\mathcal{G}}(s_t, a_t) = a_t \frac{\int_{t_u}^{\infty} u \mathcal{P}_u(u) du}{\int_{t_u}^{\infty} \mathcal{P}_u(u) du} + (1 - a_t) \frac{\int_{t_v}^{\infty} v \mathcal{P}_v(v) dv}{\int_{t_v}^{\infty} \mathcal{P}_v(v) dv}$$
(14)

with the additional constraint being that the thresholds  $t_u$  and  $t_v$  are such that the total number of admitted candidates is equal to  $A_t = \bar{a}N_t$ . Note that  $t_u$  and  $t_v$  depend on the current state  $s_t$  and action  $a_t$ .

Since the acceptance is decided by a group-wise threshold, the fraction of applicants from a group who are admitted is precisely determined by the area under its score distribution beyond the threshold. Formalizing the above, for large  $N_t$ , we have:

$$\int_{t_u}^{\infty} \mathcal{P}_u(u) du = 1 - F_u(t_u) = \frac{a_t A_t}{s_t N_t}$$
$$\int_{t_u}^{\infty} \mathcal{P}_v(v) dv = 1 - F_v(t_v) = \frac{(1 - a_t) A_t}{(1 - s_t) N_t}$$

and the constraint on the total number of candidates admitted can now be expressed through the following equivalent statements:

$$\begin{array}{rcl} a_t A_t + (1-a_t) A_t & = & \bar{a} N_t \\ s_t N_t (1-F_u(t_u)) + (1-s_t) N_t (1-F_v(t_v)) & = & \bar{a} N_t, \end{array}$$

and finally, we have:

$$s_t N_t \int_{t_u}^{\infty} \mathcal{P}_u(u) du + (1 - s_t) N_t \int_{t_v}^{\infty} \mathcal{P}_v(v) dv = \bar{a} N_t.$$
 (15)

Let us now consider the maximization of the greedy reward. Given state  $s_t$ , and generic distributions  $\mathcal{P}_u$  and  $\mathcal{P}_v$ , we need to set the thresholds  $t_u$  and  $t_v$  for the respective groups such that the sum of scores of all admitted candidates is maximized. We show by contradiction that to maximize the greedy reward, we require  $t_u = t_v$ .

Assume a pair of thresholds  $(t_u,t_v)$  that result in the maximization of the greedy reward, and  $t_u < t_v$ . Let us denote the expected sum of scores of the admitted candidates by  $S(t_u,t_v)$ , which is the optimum. One can construct thresholds  $t_u' = t_u + \epsilon_1$  and  $t_v' = t_v - \epsilon_2$  (where  $\epsilon_1,\epsilon_2>0$ , infinitesimally small for large  $N_t$ ), such that we admit one more candidate from group v (as a result of the decreased threshold) and one less from group u (as a result of the increased threshold) as compared to the case with thresholds  $(t_u,t_v)$ . As long as  $t_v'>t_u'$ , we have  $S(t_u',t_v')>S(t_u,t_v)$ , which contradicts the

assumption that  $(t_u, t_v)$  maximize the greedy reward. Similarly, if we begin with a pair of optimal  $(t_u, t_v)$  such that  $t_u > t_v$ , we can construct thresholds  $t'_u = t_u - \epsilon_3$  and  $t'_v = t_v + \epsilon_4$ , so that we admit one more candidate from group u and one less from group v. As long as  $t'_u > t'_v$ , we arrive at the contradiction  $S(t'_u, t'_v) > S(t_u, t_v)$ . Thus the greedy reward is optimized when thresholds across the groups are equal, irrespective of the nature of  $\mathcal{P}_u$  and  $\mathcal{P}_v$ .

Thus, for arbitrary score distributions, the action that maximizes the greedy reward is such that:

$$t_{u} = t_{v}$$

$$\implies F_{u}^{-1} \left( 1 - \frac{a_{t} A_{t}}{s_{t} N_{t}} \right) = F_{v}^{-1} \left( 1 - \frac{(1 - a_{t}) A_{t}}{(1 - s_{t}) N_{t}} \right)$$
 (16)

If  $\mathcal{P}_u$  and  $\mathcal{P}_v$  are identical, the arguments of the inverse CDFs in (16) need to be equal. Thus the optimal action should be such that:

$$1 - \frac{a_t A_t}{s_t N_t} = 1 - \frac{(1 - a_t) A_t}{(1 - s_t) N_t}$$

$$\implies a_t = s_t.$$

Thus, the greedy reward is maximized by choosing the admission proportion of group u to be same as the applicant proportion of group u:

$$a_{\mathcal{G}}^* = s_t.$$

Employing theorem 4.1, we arrive at the the following theorem which informs us about the convergence of the applicant pool and characterizes the FG policy.

THEOREM 4.2. For identical score distributions across the groups, the Fair-Greedy policy satisfies the following properties:

$$\begin{array}{lcl} s_t & < & \pi^*_{FG}(s_t) < \bar{s}, \ if \, s_t < \bar{s} \\ \bar{s} & < & \pi^*_{FG}(s_t) < s_t, \ if \, s_t > \bar{s} \\ & & \pi^*_{FG}(s_t) = \bar{s}, \ if \, s_t = \bar{s} \end{array}$$

Furthermore, if the step-size  $\eta_t$  decays with time and satisfies the conditions (i)  $\sum_t \eta_t = \infty$  and (ii)  $\sum_t \eta_t^2 < \infty$ , the applicant pool proportion converges to the target proportion  $\bar{s}$ . This implies that the admission or action at equilibrium also approaches the societal or target proportion, in the asymptotic regime that the total applicants in every round are large.

Proof. Under the FG policy,  $a_t = \pi_{FG}^*(s_t)$ . The applicant pool update for the mean parameter is:

$$\theta_{t+1} = [\theta_t + \eta(\pi_{FG}^*(s_t) - s_t)]_C. \tag{17}$$

The fairness loss in (2) is minimized when the admission proportion is same as the target, formalized as:

$$a_{\mathcal{F}}^* = \arg\min_{a_t} L_{\mathcal{F}}(a_t) = \bar{s}$$

The overall reward  $R(s_t, a_t)$  is a sum of the greedy reward and fairness loss (scaled by  $\lambda$ ). The fairness loss is convex (hence  $-L_{\mathcal{F}}(a_t)$  is concave) in  $a_t$ . It can be seen that the greedy reward monotonically decreases in either directions around  $a_t = s_t$ , and in addition it possesses continuity in  $a_t$ . When at state  $s_t$ , suppose the optimal action  $a^*$  of the FG policy is such that  $a^* < s_t$ , when  $s_t < \bar{s}$ . Then by continuity and since the greedy reward is maximized at  $s_t$ ,  $\exists$  some  $a' > s_t$ , such that  $R_{\mathcal{G}}(s_t, a') \geq R_{\mathcal{G}}(s_t, a^*)$ , and moreover has

a smaller fairness loss, i.e.,  $L_{\mathcal{F}}(a') < L_{\mathcal{F}}(a^*)$ , which violates the optimality of  $a^*$ . Thus the optimal action for the FG policy must be  $a^* > s_t$ , if  $s_t < \bar{s}$ . Similar arguments hold if  $s_t > \bar{s}$ , and here we can show that the optimal action must be such that  $a^* < s_t$ . Hence, it follows that the optimal action for overall utility lies between the optimal actions for greedy and fairness terms:

$$s_t < \pi_{FG}^*(s_t) < \bar{s}, \text{ if } s_t < \bar{s}$$
 (18)

$$\bar{s} < \pi_{FG}^*(s_t) < s_t, \text{ if } s_t > \bar{s}$$
 (19)

$$\pi_{FG}^*(s_t) = \bar{s}, \text{ if } s_t = \bar{s} \tag{20}$$

Now we show the convergence of the applicant pool to its equilibrium. Let us consider a step-size that decays with time such that  $\sum_t \eta_t = \infty$  and  $\sum_t \eta_t^2 < \infty$ . Consider the case when  $s_t < \bar{s}$ , where we have:  $s_t < \pi_{FG}^*(s_t) < \bar{s}$ . From (17), we can see that the mean proportion parameter  $\theta_{t+1}$  increases. Similarly, when  $s_t > \bar{s}$ , it follows that  $\bar{s} < \pi_{FG}^*(s_t) < s_t$ , and the mean proportion parameter decreases. Note that the target proportion is a fixed point of the FG policy, i.e.,  $\pi_{FG}^*(\bar{s}) = \bar{s}$ . Due to the above characterization of  $\pi_{FG}^*(s_t)$  and the model for the update of the applicant pool, the mean parameter  $\theta_t$  grows or reduces in the direction of  $\bar{s}$ . Hence, as the step-size is decaying, one can show that the mean parameter  $\theta_t$  converges to  $\bar{s}$  (see appendix A for details). Moreover, the variance of the number of group u applicants is  $var(N_t^u) = \theta_t N_t$ due to the Poisson distribution. Thus, the state  $s_t = N_t^u/N_t$  has variance  $O(1/N_t)$ . Consequently, in the asymptotic regime that  $N_t$  is large, using Chebyshev's inequality one can show that  $s_t$ also converges to  $\theta_t$  in probability. This implies that the applicant proportion approaches  $\bar{s}$ , which completes the proof.

## 5 EXPERIMENTAL EVALUATION

# 5.1 Evaluation on synthetic data

We begin by employing synthetic data to demonstrate the fairness framework we develop in this paper, and study interesting scenarios.

Optimal policy based on value iteration. Let us first consider the MDP setting from Section 3, where the policy learnt is the optimal policy (8) maximizing the accumulated utilities. Consider the case where the two groups have identical score distributions. This may often be the case in real-world scenarios when there is no inherent reason for the sensitive attribute to influence the scores or proficiency of a candidate. Let the score distributions be Gaussian with means  $\mu_u = \mu_v = 5$  and variances  $\sigma_u^2 = \sigma_v^2 = 1$ . The societal/target proportion can be set by employing guidance from the societal state or based upon the long-term target that the selector has in mind. For example, suppose our application is to hire software engineers, then representing women as group u, one can set the target proportion to be the proportion of women in computer science, or in the society in general. Or, if we target to have a certain proportion of women in the company in the longterm, we could set  $\bar{s}$  accordingly. In this experiment, we set  $\bar{s} = 0.4$ and the admission rate is fixed to  $\bar{a} = 0.3$ , or in other words, the selector aims to admit only 30% of the total applied candidates. The other parameter values used for this experiment are  $\gamma = 0.99$ ,  $\lambda = 1.5$ , a fixed step-size of  $\eta = 0.05$ . Figure 1 shows how the proportion of applicants, admitted candidates and mean parameter

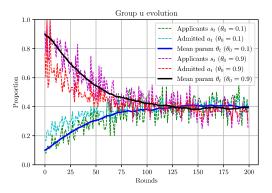


Figure 1: Optimal policy under identical score distributions across the groups.

 $\theta_t$  vary for group u. We see in the figure that beginning the process from different initial states  $\theta_0 = 0.1, 0.9$ , we observe convergence of the applicant pool proportion for group u. The optimal policy under the evolution model considered has resulted in close to 40% of the applicants belonging to group u, and also approximately the same proportion of the admitted candidates are from group u.

FG policy under identical score distributions. Let us now consider the same score distributions, target proportion and overall admission rate  $\bar{a}$  as above, but with the FG policy in (9), described in Section 4. The step-size for changes in applicant pool mean parameter is fixed to  $\eta = 0.05$ , though a decaying step-size would in fact aid in smoother convergence behavior. Figure 2 shows the convergence of the applicant pool to the target proportion of 40%, and the proportion of admitted candidates belonging to group u is also around 0.4, as guaranteed by our analysis of the FG policy. We also observe that the FG policy follows the structure stated in Theorem 4.2. The framework is capable of handling an inversion in the majority and minority proportions as supported by the evolutions shown from two distinct initial applicant mean proportion parameters  $\theta_0 = 0.1$ and  $\theta_0 = 0.9$ . We report on the dynamics for the proportion of applicants and admitted candidates for individual sample paths in which the number of applicants is randomly drawn as in (3). We do not smooth over multiple sample paths in such figures because our objective is to highlight the convergence of the mean parameter  $\theta_t$ over each sample path. Note that tuning of the hyperparameter  $\lambda$ is not required when score distributions are identical (here we set  $\lambda = 2$ ). As long as  $\lambda > 0$ , the applicant pool converges to the target proportion, with only the rate of convergence increasing with  $\lambda$ , as we depict in Figure 3.

FG policy: under selective applications. Next, we focus on a setting where the underprivileged class u has larger variance, but slightly smaller mean ( $\sigma_u^2=1.5$ ,  $\mu_u=4.9$ ). We set  $\bar{s}=0.4$ , and consider a more selective process, with  $\bar{a}=0.1$ . Typically, such cases might occur when the data about unprivileged group is unreliable or there is imbalance in the amount of samples available, leading to a larger variance. From Figure 4, we note that the applicant mean and also the group admission converges to a proportion larger than  $\bar{s}$ . This is due to the fact that as the admission rate gets selective,

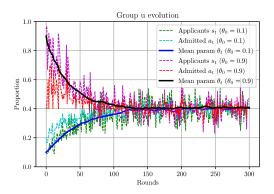


Figure 2: FG policy under identical score distribution across groups, showing convergence from distinct initial mean parameters  $\theta_0 = 0.1, 0.9$ .

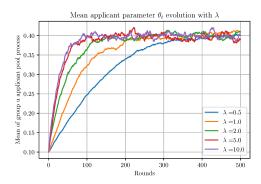


Figure 3: Applicant pool converges to the target proportion for identical score distributions under the FG policy.

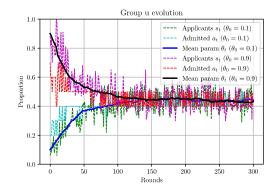


Figure 4: FG policy under selective system, lower mean and larger variance for group u. Shows convergence from  $\theta_0 = 0.1, 0.9$ .

the greedy part of the reward is optimized by an action that admits more from the group with longer tail (larger variance), which is the unprivileged group u in this case. Hence the greedy reward

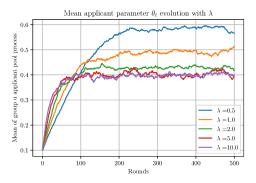


Figure 5: Applicant pool convergence for the selective system under FG policy.

promotes more admission from group u. This is also evident in Figure 5, where we observe that for smaller values of  $\lambda$ , i.e., when more weight is assigned to the greedy reward, the mean parameter  $\theta_t$ , which measures the expected proportion of group u applicants, converges to larger values. However with enough weight being given to fairness, the applicant pool still converges to the desired ratio.

# 5.2 Dynamically adapted real-world datasets

We simulate the dynamics by considering the following: (i) the law school bar study dataset, applying our framework from the viewpoint of a recruiter selecting candidates who are likely to be successful in the bar exam, based on features such as LSAT scores, undergraduate GPA, law school GPA and others, while maintaining equity based on race as the sensitive attribute. The aim of positive reinforcement is to drive the system towards a richer pool of applicants. (ii) German credit dataset with gender as the sensitive attribute, where the motivation is to encourage higher levels of participation of women in the financial lending system.

The law school bar study dataset [23] consists of data collected by a Law School Admission Council survey across law schools in the United States. The predictions indicate whether or not a candidate would pass the bar exam based on features such as LSAT scores, undergraduate GPA, law school GPA, race, sex, family income, age and so on. We consider race as the sensitive attribute, and though originally there are 8 distinct races in the dataset, we group the samples by combining samples corresponding to all others except 'white', giving rise to binary groups 'white' and 'non-white'. We observe that the data is imbalanced – about only 25% of the samples belong to group 'non-white', which we will label as group u. This proportion will serve as a starting point for the applicant pool composition. We use a version of the bar study dataset found at [1] with around 1800 instances. A longer and more popular version of the same in fairness literature is the law school GPA admissions dataset which comprises of about 21,790 samples and the labels indicate if an applicant will have a high first year average GPA. We fit score distributions on this dataset as well, but choose the bar study dataset to study the dynamics of positive reinforcement and observe how the decisions of admitting candidates who are more

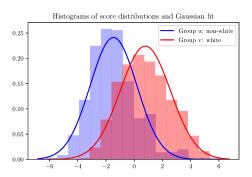


Figure 6: Histograms and Gaussian fit for score distributions of law school bar study dataset

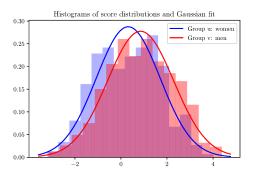


Figure 7: Histograms and Gaussian fit for score distributions of German credit dataset

likely to eventually succeed in the bar exam affects the composition of the applicant pool.

The German credit dataset [9] consists of 1000 instances, with 20 features (both numeric and qualitative), such as credit history, account history, employment status, age, gender and so on. This is typically used to assess the risk of lending loans to people, i.e., to determine if granting credit is risky or not. We consider gender as the binary valued sensitive attribute, labeling women as group u and men as group v. The dataset is imbalanced – about 31% of the instances belong to group u.

After pre-processing the datasets to suit our usage, our first step is to learn score distributions that measure the proficiency of candidates. To achieve this, we fit a predictor based on logistic regression that uses the features and labels to fit scores, which are the derived as the product of the model coefficients and the features. We observe that the histograms of the scores of the two groups reveal that they are indeed Gaussian in nature. We fit a Gaussian for each of the histograms, to obtain the mean and variance parameters of the score distributions  $\mathcal{P}_u$  and  $\mathcal{P}_v$ .

The histograms and the Gaussian fit for the score distributions for the law school bar study and German credit dataset are depicted in Figures 6 and 7 respectively. For the law school bar study dataset the parameters of scores are  $\mu_u = -1.46$ ,  $\sigma_u^2 = 2.73$ ,  $\mu_v = 0.79$ ,

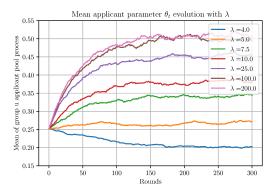


Figure 8: Law school bar study dataset: applicant pool convergence with initial mean proportion parameter  $\theta_0=0.25$ , as  $\lambda$  is varied.

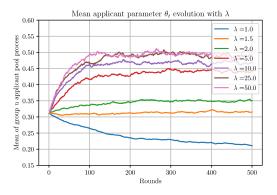


Figure 9: German credit dataset: applicant pool convergence with initial mean proportion parameter  $\theta_0=0.31$ , as  $\lambda$  is varied.

 $\sigma_v^2 = 3.16$ . For the German credit dataset the score distributions are closer with parameters  $\mu_u = 0.32$ ,  $\sigma_u^2 = 1.93$ ,  $\mu_v = 0.85$ ,  $\sigma_v^2 = 2.06$ .

We now simulate the dynamics of the application process, under the FG policy, by sampling from these distributions with initial state of the applicant process  $\theta_0$  determined by the number of instances of respective groups, which is 0.25 for the law school bar study and 0.31 for the German credit datasets respectively. The variation of the applicant pool for different values of hyperparameter  $\lambda$  are shown for the datasets in Figures 8 and 9 respectively. The evolution step size used in these simulations is  $\eta = 0.025$ , admission rate is set to  $\bar{a} = 0.3$  and the target proportion is set to  $\bar{s} = 0.5$ , which is equivalent to demographic parity, i.e., admitting same number proportion of candidates from both groups. In both the figures, we observe that when the greedy reward is favored (lower values of  $\lambda$ ), the applicant pool in fact converges to a point lesser than the target, while it approaches the target as  $\lambda$  increases. This means that for maximizing the utility, more samples need to be admitted from group v, due to the nature of their score distributions, when less importance is allotted to fairness objective. The tuning of the

Table 1: Gaussian score distribution parameters for different datasets

Dataset	Sensitive attribute	$\mu_{u}$	$\mu_v$	$\sigma_u^2$	$\sigma_v^2$
LS bar study	included	-1.46	0.79	2.73	3.16
LS bar study	excluded	-1.33	0.76	2.85	3.23
German credit	included	0.32	0.85	1.93	2.06
German credit	excluded	0.62	0.84	2.03	2.14
LS GPA admissions	included	1.45	3.14	2.44	1.89
LS GPA admissions	excluded	1.51	3.13	2.50	1.93

hyperparameter  $\lambda$  to achieve desired level of applicant pool proportion depends on the order statistics of  $\mathcal{P}_u$  and  $\mathcal{P}_v$ . The step-size parameter  $\eta$  can be set appropriately based on how quickly we wish to achieve convergence.

These experiments with real-world datasets indicate that scores which are fit after learning predictors based on logistic regression are distributed like Gaussians. Once we have the parameters of the scores, the application of the FG policy and the applicant pool evolution follows.

It is interesting to examine how the score distributions change when we approach fairness through unawareness, that is, by omitting the sensitive attributes while learning the logistic regression based predictor. Note that we learn a single predictor based on all samples and then distinguish the scores based on the sensitive attribute. Table 1 lists the score parameters when the predictor is learnt with or without the inclusion of the sensitive attribute for the law school bar study, the law school GPA admissions and the German credit datasets. Similar to bar study dataset, we employ race as the sensitive attribute to the GPA admissions dataset as well. For both the law school datasets, we observe that the score distributions are not very different, although the difference between the means of minority and majority groups has decreased slightly when the sensitive attribute is dropped during the learning. For the German credit dataset, the distributions are significantly closer when the sensitive attribute is omitted, and there is a clear drop in the difference between the group means.

## 6 CONCLUSION

As AI-based decision-making becomes increasingly impactful on human society, the study of the influence of fairness-aware policies on the population becomes important. In this paper, we proposed a framework for fair selection of applicants to a system, and studied the long-term effects of decisions on the composition of the applicant pool. We proposed an optimal policy based on dynamic programming, and also a simple Fair-Greedy policy that optimizes for instantaneous utility. We characterized the FG policy for general score distributions and proved that the applicant pool approaches the target proportion when score distributions are identical across groups. Experimental evaluation reveals that by appropriately choosing the hyperparameter  $\lambda$ , a desired equilibrium point in the applicant pool composition can be achieved for generic score distributions.

Our results indicate the potential of achieving long-term fairness objectives through positive reinforcement via decision making. We hope that this work stimulates the collaboration between machine learning researchers and social scientists required for these ideas to make real-world impact. A key future direction is to devise and conduct experiments for measuring, understanding and shaping the evolution dynamics posited in our framework.

#### ACKNOWLEDGMENTS

This work was supported by the Army Research Office under grant W911NF-19-1-0053, and by the National Science Foundation under grant CCF 1909320.

## REFERENCES

- 2018. Law School Bar Study Dataset. Retrieved February 2022 from https://github.com/algowatchpenn/GerryFair/blob/master/dataset/lawschool.csv
- [2] Yahav Bechavod, Katrina Ligett, Aaron Roth, Bo Waggoner, and Steven Z Wu. 2019. Equal opportunity in online classification with partial feedback. Advances in Neural Information Processing Systems 32 (2019).
- [3] James A Berkovec, Glenn B Canner, Stuart A Gabriel, and Timothy H Hannan. 2018. Mortgage discrimination and FHA loan performance. In Mortgage Lending, Racial Discrimination, and Federal Policy. Routledge, 289–305.
- [4] Dimitri Bertsekas. 2007. Dynamic programming and optimal control: Volume II.
   Vol. 2. Athena scientific.
- [5] Steven Broad and Meredith McGee. 2014. Recruiting Women into Computer Science and Information Systems. Association Supporting Computer Users in Education (2014).
- [6] Yifang Chen, Alex Cuellar, Haipeng Luo, Jignesh Modi, Heramb Nemlekar, and Stefanos Nikolaidis. 2020. Fair contextual multi-armed bandits: Theory and experiments. In Conference on Uncertainty in Artificial Intelligence. PMLR, 181– 190.
- [7] Christos Dimitrakakis, Yang Liu, David C Parkes, and Goran Radanovic. 2019.
   Bayesian fairness. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 509–516.
- [8] Julia Dressel and Hany Farid. 2018. The accuracy, fairness, and limits of predicting recidivism. Science advances 4, 1 (2018), eaao5580.
- [9] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. http://archive.ics.uci.edu/ml
- [10] Danielle Ensign, Sorelle A Friedler, Scott Neville, Carlos Scheidegger, and Suresh Venkatasubramanian. 2018. Runaway feedback loops in predictive policing. In Conference on Fairness, Accountability and Transparency. PMLR, 160–171.
- [11] Ganesh Ghalme, Vineet Nair, Vishakha Patil, and Yilun Zhou. 2021. State-Visitation Fairness in Average-Reward MDPs. arXiv preprint arXiv:2102.07120 (2021).
- [12] Stephen Gillen, Christopher Jung, Michael Kearns, and Aaron Roth. 2018. Online learning with an unknown fairness metric. Advances in neural information processing systems 31 (2018).
- [13] Hoda Heidari and Andreas Krause. 2018. Preventing Disparate Treatment in Sequential Decision Making.. In IJCAI. 2248–2254.
- [14] Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, and Aaron Roth. 2017. Fairness in reinforcement learning. In *International conference on machine learning*. PMLR, 1617–1626.
- [15] Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. 2018. Meritocratic fairness for infinite and contextual bandits. In Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society. 158–163.
- [16] Faisal Kamiran and Toon Calders. 2012. Data preprocessing techniques for classification without discrimination. Knowledge and information systems 33, 1 (2012), 1–33.
- [17] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2012. Fairness-aware classifier with prejudice remover regularizer. In Joint European conference on machine learning and knowledge discovery in databases. Springer, 35–50
- [18] Lydia T Liu, Sarah Dean, Esther Rolf, Max Simchowitz, and Moritz Hardt. 2018. Delayed impact of fair machine learning. In *International Conference on Machine Learning*. PMLR, 3150–3158.
- [19] Lydia T Liu, Ashia Wilson, Nika Haghtalab, Adam Tauman Kalai, Christian Borgs, and Jennifer Chayes. 2020. The disparate equilibria of algorithmic decision making when individuals invest rationally. In Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency. 381–391.
- [20] Hussein Mouzannar, Mesrob I Ohannessian, and Nathan Srebro. 2019. From fair decision making to social equality. In Proceedings of the Conference on Fairness, Accountability, and Transparency. 359–368.

- [21] Vishakha Patil, Ganesh Ghalme, Vineet Nair, and Yadati Narahari. 2020. Achieving Fairness in the Stochastic Multi-Armed Bandit Problem.. In AAAI. 5379–5386.
- [22] Min Wen, Osbert Bastani, and Ufuk Topcu. 2021. Algorithms for fairness in sequential decision making. In *International Conference on Artificial Intelligence* and Statistics. PMLR, 1144–1152.
- [23] Linda F Wightman. 1998. LSAC National Longitudinal Bar Passage Study. LSAC Research Report Series. (1998).
- [24] Joshua Williams and J Zico Kolter. 2019. Dynamic modeling and equilibria in fair decision making. arXiv preprint arXiv:1911.06837 (2019).
- [25] Sirui Yao and Bert Huang. 2017. Beyond parity: Fairness objectives for collaborative filtering. Advances in neural information processing systems 30 (2017).
- [26] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. 2017. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In Proceedings of the 26th international conference on world wide web. 1171–1180.
- [27] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *International conference on machine learning*. PMLR. 325–333.
- [28] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. 2018. Mitigating unwanted biases with adversarial learning. In Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society. 335–340.
- [29] Xueru Zhang, Mohammadmahdi Khaliligarekani, Cem Tekin, and Mingyan Liu. 2019. Group retention when using machine learning in sequential decision making: the interplay between user dynamics and fairness. Advances in Neural Information Processing Systems 32 (2019).
- [30] Xueru Zhang, Ruibo Tu, Yang Liu, Mingyan Liu, Hedvig Kjellstrom, Kun Zhang, and Cheng Zhang. 2020. How do fair decisions fare in long-term qualification? Advances in Neural Information Processing Systems 33 (2020), 18457–18469.

## A APPENDIX: PROOF DETAILS

Lemma A.1. If the step-size  $\eta_t$  decays with time and satisfies the conditions (i)  $\sum_t \eta_t = \infty$  and (ii)  $\sum_t \eta_t^2 < \infty$ , the mean of the applicant pool proportion for group u converges to the target proportion  $\bar{s}$  under the FG policy, when the score distributions across the groups are identical.

PROOF. We wish to show that  $\theta_t \to \bar{s}$  as  $t \to \infty$ . Let  $d_t = \frac{1}{2}(\theta_t - \bar{s})^2$ . Fix an  $\epsilon > 0$ . We need to show that there exists some  $t_0(\epsilon)$  such that when  $t \ge t_0(\epsilon)$ ,

$$d_{t+1} \leq d_t - \gamma_t$$
, if  $d_t \geq \epsilon$  (21)

$$d_{t+1} < c\epsilon$$
, if  $d_t < \epsilon$  (22)

where c is a positive constant. Moreover  $\gamma_t > 0$  and  $\sum_t \gamma_t = \infty$ . If the above hold, then eventually for some  $t = t_1(\epsilon) \ge t_0(\epsilon)$ , one has  $d_t < \epsilon$ . But due to (21) and (22)  $d_t < c\epsilon$  for all  $t > t_1(\epsilon)$ . Since  $\epsilon$  is arbitrary,  $\theta_t \to \bar{s}$  as  $t \to \infty$ .

We first show that (22) holds.

$$\begin{split} d_{t+1} &= \frac{1}{2}(\theta_{t+1} - \bar{s})^2 \\ &= \frac{1}{2}([\theta_t - \eta_t(s_t - a_t)]_C - \bar{s})^2 \\ &\leq \frac{1}{2}(\theta_t - \eta_t(s_t - a_t) - \bar{s})^2 \\ &= d_t + \eta_t(\bar{s} - \theta_t)(s_t - a_t) + \frac{1}{2}\eta_t^2(s_t - a_t)^2 \\ &\leq d_t + \eta_t(\bar{s} - \theta_t)(s_t - a_t) + \frac{1}{2}\eta_t^2 \\ &\leq d_t + \frac{\eta_t}{2}((\bar{s} - \theta_t)^2 + 1) + \frac{1}{2}\eta_t^2 \end{split}$$

Since  $\eta_t$  is arbitrarily small, if  $d_t < \epsilon$ , we have:

$$d_{t+1} < c\epsilon. (23)$$

When  $d_t \geq \epsilon$ , we want to first show that

$$(\bar{s} - \theta_t)(\theta_t - a_t) \le -\delta(\epsilon)$$
 (24)

where  $\delta(\epsilon) > 0$ . If this holds, we have,

$$d_{t+1} \le d_t - \eta_t \delta(\epsilon) + \frac{1}{2} \eta_t^2. \tag{25}$$

Let us denote  $\gamma_t = \eta_t \delta(\epsilon) - \frac{1}{2} \eta_t^2$ . Since  $\eta_t \to 0$ , there exists some  $t_2(\epsilon)$  such that  $\gamma_t > 0$  for  $t > t_2(\epsilon)$ . Moreover, due to conditions on step size, we have  $\sum_t \gamma_t = \infty$ .

Next, we will account for the stochasticity of  $s_t$ . We have  $s_t - a_t = \theta_t + (s_t - \theta_t) - a_t$ . Denoting  $z_t = s_t - \theta_t$ , we have

$$d_{t+1} \le d_t + \eta_t (\bar{s} - \theta_t) (\theta_t + z_t - a_t) + \frac{1}{2} \eta_t^2$$
 (26)

 $z_t$  is a zero-mean random variable. Also  $E[z_t^2] = var(s_t) = \theta_t/N_t$ , which is bounded. Therefore  $v_t := \sum_{m=0}^t \eta_m z_m$  is a martingale, and

 $E[v_t^2]$  is also bounded. This implies, by the martingale convergence theorem, that  $v_t$  converges to a finite random variable. Therefore, we have  $\sum_{m=t}^{\infty} \eta_m z_m \to 0$ . Since  $|\theta_t - \bar{s}|$  is bounded, the effect of noise  $z_t$  is asymptotically negligible.

What remains to be shown is (24). In the regime of large number of applicants  $N_t$ , we can see that the state  $s_t$  is equal to its mean  $\theta_t$  with probability approaching one, through the Chebyshev inequality.

When  $d_t \ge \epsilon$ , since  $s_t$  is equal to  $\theta_t$ , we need to consider only the cases (i)  $s_t > \bar{s}$  and (ii)  $s_t < \bar{s}$ . Under both these cases, we have  $(\bar{s} - \theta_t)(\theta_t - a_t) < 0$  due to the structure of the FG policy in (18) and (19), when the score distributions across the groups are identical