# **Fairness with Dynamics**

# Min Wen <sup>1</sup> Osbert Bastani <sup>1</sup> Ufuk Topcu <sup>2</sup>

## **Abstract**

It has recently been shown that if feedback effects of decisions are ignored, then imposing fairness constraints such as demographic parity or equality of opportunity can actually exacerbate unfairness. We propose to address this challenge by modeling feedback effects as the dynamics of a Markov decision processes (MDPs). First, we define analogs of fairness properties that have been proposed for supervised learning. Second, we propose algorithms for learning fair decision-making policies for MDPs. We also explore extensions to reinforcement learning, where parts of the dynamical system are unknown and must be learned without violating fairness. Finally, we demonstrate the need to account for dynamical effects using simulations on a loan applicant MDP.

### 1. Introduction

Machine learning has the potential to substantially improve performance in legal and financial decision-making. However, it has been demonstrated that biases in the data can be reflected in a decision-making policy trained on that data (Dwork et al., 2012), which can result in decisions that unfairly discriminate against minorities. For example, consider the problem of giving loans to applicants (Hardt et al., 2016). If minorities are historically given loans less frequently, then there may be less data on how reliably they repay loans. Thus, a learned decision-making policy may unfairly label minorities as higher risk and deny them loans.

There have been several candidate definitions of fairness, including demographic parity (i.e., members of the majority and miniority subpopulations must on average have equal outcomes) (Calders et al., 2009), equality of opportunity (i.e., *qualified* members must on average have equal outcomes) (Hardt et al., 2016), and causal fairness (i.e., protected attributes should not influence outcomes) (Kusner et al., 2017; Kilbertus et al., 2017; Nabi & Shpitser, 2018).

The appropriate definition depends on the application.

So far, work on fairness has focused on supervised learning. However, it has recently been shown that naïvely imposing fairness constraints while ignoring even one-step feedback effects can actually harm minorities (Liu et al., 2018). Thus, it is critical that we extend existing definitions of fairness to account for the feedback effects of the decisions being made on members of the population. For example, denying loans to individuals may have consequences on their financial security that need be taken into account.

This paper proposes algorithms for learning fair decision-making policies that account for feedback effects of decisions. We model these effects as the dynamics of a Markov decision process (MDP), and extend existing fairness definitions to decision-making policies for a known MDP. Unlike supervised learning, we distinguish the quality of outcomes for the decision-maker (e.g., the bank) from the quality of the outcomes for individuals (e.g., a loan applicant). Then, fairness properties are constraints on the average quality of outcomes for individuals in different subpopulations (e.g., majorities and minorities are offered loans at the same frequency), whereas the reward measures the quality of outcomes for the decision-maker (e.g., the bank's profit).

The key challenge is that learning with a fairness constraint is much more challenging in the MDP setting due to the inherent non-convexity—indeed, constrained reinforcement learning is an active research area (Altman, 1999; Achiam et al., 2017; Wen & Topcu, 2018; Bastani et al., 2018). Building on this work, we propose both model-based (Altman, 1999) and model-free (Wen & Topcu, 2018) algorithms for learning policies that satisfy fairness constraints. We initially assume the dynamics are known, but propose extensions to settings where parts of the dynamics are unknown.

Prior work has studied one-step feedback effects (Liu et al., 2018), but do not propose fair learning algorithms. For learning unknown dynamics, prior work has focused on bandit settings with specific fairness constraints (Joseph et al., 2016; Hashimoto et al., 2018). There has been work on learning unknown dynamics in MDPs (Jabbari et al., 2017; Elzayn et al., 2019), but for a specific fairness contraint. Furthermore, for their constraint, the optimal policy is always fair. Thus, unlike our setting, solving for the optimal fair policy is trivial once the dynamics are known.

<sup>&</sup>lt;sup>1</sup>University of Pennsylvania <sup>2</sup>University of Texas. Correspondence to: Min Wen <wenm@seas.upenn.edu>, Osbert Bastani <obastani@seas.upenn.edu>.

We compare to two baselines that ignore dynamics: (i) an algorithm that optimistically pretends that actions do not affect the state distribution (i.e., supervised learning), and (ii) an algorithm that conservatively assumes the state distribution can change adversarially on each step. In a simulation study on a loan applicant MDP based on (Hardt et al., 2016), we show that compared to our algorithm, the optimistic algorithm learns unfair policies, and the conservative algorithm learns fair but poorly performing policies. Our results demonstrate the importance of accounting for dynamics.

### 2. Problem Formulation

Markov decision processes. A Markov decision process (MDP) is a tuple  $M=(S,A,D,P,R,\gamma)$ , where  $S=[n]=\{1,...,n\}$  are the states, A=[m] are the actions,  $D\in\mathbb{R}^{|S|}$  is the initial state distribution (i.e.,  $D_s$  is the probability of starting in state s),  $P\in\mathbb{R}^{|S|\times|A|\times|S|}$  are the transitions (i.e.,  $P_{s,a,s'}$  is the probability of transitioning from state s to state s' when taking action a),  $R\in\mathbb{R}^{|S|\times|A|}$  are the rewards (i.e.,  $R_{s,a}$  is the reward obtained taking action a in state s), and  $\gamma\in\mathbb{R}$  is the discount factor. Given a stochastic policy  $\pi\in\mathbb{R}^{|S|\times|A|}$  (i.e.,  $\pi_{s,a}$  is the probability of taking action a in state s), the induced state transtion probabilities are  $P^{(\pi)}\in\mathbb{R}^{|S|\times|S|}$ , where

$$P_{s,s'}^{(\pi)} = \sum_{a \in A} \pi_{s,a} P_{s,a,s'},$$

the induced distribution over states at time t is

$$D^{(\pi,t)} = \begin{cases} D & \text{if } t = 0\\ P^{(\pi)}D^{(\pi,t-1)} & \text{otherwise,} \end{cases}$$

and the time-discounted overall distribution over states is

$$D^{(\pi)} = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t D^{(\pi,t)} = (1 - \gamma) \sum_{t=0}^{\infty} (\gamma P^{(\pi)})^t D.$$

Then, the time-discounted overall distribution over state-action pairs is  $\Lambda \in \mathbb{R}^{|S| \times |A|}$ , where

$$\Lambda_{s,a}^{(\pi)} = D_s^{(\pi)} \pi_{s,a},$$

and the expected cumulative reward is

$$R^{(\pi)} = (1 - \gamma)^{-1} \langle R, \Lambda^{(\pi)} \rangle,$$

where  $\langle X, Y \rangle = \sum_{s \in S} \sum_{a \in A} X_{s,a} Y_{s,a}$ . Given a policy class  $\Pi$ , the *optimal policy* is  $\pi^* = \arg \max_{\pi \in \Pi} R^{(\pi)}$ .

**Example.** We describe an MDP  $M_{\rm loan}$  that models individuals applying for loans. We assume each individual has a true probability p of repaying their loan. On step t, the bank has a prior on p (e.g., a credit score); for simplicity, we assume this prior is a Beta distribution—i.e.,

 $p_t \sim \text{Beta}(\alpha_t, \beta_t)$ . Thus, the states of our MDP  $(\alpha_t, \beta_t)$ . <sup>1</sup> The actions are to offer (a=1) or deny (a=0) a loan. If the bank offers a loan, the transitions are

$$(\alpha_{t+1}, \beta_{t+1}) = \begin{cases} (\alpha_t + 1, \beta_t) & \text{with probability } p \\ (\alpha_t, \beta_{t+1}) & \text{with probability } 1 - p. \end{cases}$$

If the bank denies the loan, then mathematically, the parameters of the posterior are  $(\alpha_{t+1}, \beta_{t+1}) = (\alpha_t, \beta_t)$ . However, since we are interested in detrimental effects of the bank's decisions, we consider the possibility that this decision reduces the applicant's ability to pay for future loans:

$$(\alpha_{t+1}, \beta_{t+1}) = (\alpha_t, \beta_t + \tau),$$

for some  $\tau \in \mathbb{R}_+$ —e.g., if a loan is denied, the applicant may resort to more expensive loans, thus reducing their wealth. We assume the initial state distribution is  $z \sim \text{Bernoulli}(p_Z)$  and  $(\alpha,\beta) \sim p_0(\alpha,\beta \mid z)$  for some  $p_Z \in [0,1]$  and some distribution  $p_0$ . The bank's rewards are

$$\mathbb{E}_{\delta}[\delta I - (1 - \delta)P] - \lambda \sqrt{\text{Var}_{\delta}[\delta I - (1 - \delta)P]}, \quad (1)$$

where P is the principal (without loss of generality, we let P=1), I is interest,  $\delta$  indicates whether the loan is repaid, and  $\lambda \in \mathbb{R}_+$ . The first term is expected profit, and the second term is to avoid risk. We assume that the bank makes decisions with the goal of maximizing (1).

The fairness constraint refers to the rewards for individuals, which we call *agent rewards*. In our example, the agent rewards are  $\mathbb{I}[a=1]$ , where  $\mathbb{I}$  is the indicator function—i.e., a positive outcome is when the individual is offered a loan.

**Fairness.** Consider a population of individuals (e.g., loan applicants) interacting with a decision-maker (e.g., a bank) (Hardt et al., 2016). States S encode an individual's features (e.g., probability of repaying a loan), actions A are interventions (e.g., loan offer), and transitions P encode state changes (e.g., changes in ability to repay). We use rewards R to indicate quality of outcomes for the decision-maker (e.g., the bank's profit), and use agent rewards P0 ∈  $\mathbb{R}^{|S| \times |A|}$ 1 to indicate quality of outcomes for an individual (e.g., whether a loan is offered).

Our goal is to learn the optimal policy for the decision-maker under a fairness constraint. In particular, we want to ensure that  $\pi$  does not favor the *majority subpopulation* over the *minority subpopulation*. We assume the state space has the form  $S = \tilde{S} \times Z$ , where  $Z = \{\text{maj, min}\}$  encodes whether an individual  $(\tilde{s}, z) \in S$  is from the majority (z = maj) or minority (z = min) subpopulation, and  $\tilde{S}$  encodes non-sensitive individual characteristics (e.g., probability of repaying a loan). We base our constraints on those for supervised learning (Hardt et al., 2016).

 $<sup>^{1}</sup>$ Technically, our MDP is the belief MDP of the POMDP where the state p is unobserved.

**Definition 2.1.** Let M be an MDP with state space of the form  $S = Z \times \tilde{S}$ , where  $Z = \{\text{maj, min}\}$ , and let  $\rho \in \mathbb{R}^{|S| \times |A|}$  be the agent rewards. Then, a policy  $\pi$  satisfies demographic parity if

$$\begin{split} \mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{maj}}^{(\pi)}}[\rho_{s,a}] &= \mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{min}}^{(\pi)}}[\rho_{s,a}], \\ \text{where } \Lambda_z^{(\pi)} &= \Lambda^{(\pi)} \mid \exists \tilde{s} \in \tilde{S} \ . \ s_0 = (z,\tilde{s}) \ \text{for } z \in Z \text{—i.e.,} \\ (\Lambda_z^{(\pi)})_{s,a} &= (D_z^{(\pi)})_s \pi_{s,a} \qquad (\forall s \in S, a \in A) \\ (D_z^{(\pi)})_s &= (1-\gamma) \sum_{t=0}^{\infty} (\gamma P^{(\pi)})^t D_z \qquad (\forall s \in S) \\ (D_z)_{s_0} \propto D_{s_0} \cdot \mathbb{I}[\exists \tilde{s} \in \tilde{S} \ . \ s_0 = (z,\tilde{s})] \qquad (\forall s_0 \in S). \end{split}$$

That is,  $\Lambda_z^{(\pi)}$  is  $\Lambda^{(\pi)}$  conditioned on the initial state  $s_0$  being of the form  $s_0 = (z, \tilde{s}_0)$  for some  $\tilde{s}_0 \in \tilde{S}$ , and demographic parity says that the cumulative agent rewards are equal on average for the majority and minority subpopulations. For  $M_{\text{loan}}$ , demographic parity says that loans should be given to majority and minority members with equal frequency.

Our goal is to compute the optimal policy for the policy class  $\Pi_{DP}$  of policies that satisfy demographic parity:

$$\pi^* = \operatorname*{arg\,max}_{\pi \in \Pi_{\mathrm{DP}}} R^{(\pi)}. \tag{2}$$

**Remark 2.2.** We focus on demographic parity, but the techniques we develop are general. In particular, they apply to any fairness constraint saying that two subpopulations should have equal outcomes on average—i.e., for any set of subsets  $S_z \subseteq S$  for each  $z \in Z$ , letting

$$\tilde{\Lambda}_z^{(\pi)} = \tilde{\Lambda}^{(\pi)} \mid \mathbb{I}[s_0 \in S_z],$$

then we can handle the fairness constraint

$$\mathbb{E}_{(s,a)\sim\tilde{\Lambda}_{\mathrm{maj}}^{(\pi)}}[\rho_{s,a}] = \mathbb{E}_{(s,a)\sim\tilde{\Lambda}_{\mathrm{min}}^{(\pi)}}[\rho_{s,a}].$$

For example, our techniques also apply to equality of opportunity (Hardt et al., 2016) and path-specific causal fairness (Nabi & Shpitser, 2018); see Appendix A.

**Remark 2.3.** We sometimes consider a fairness with an  $\epsilon$  tolerance (for some  $\epsilon \in \mathbb{R}_+$ ):

$$\left| \mathbb{E}_{(s,a) \sim \Lambda_{-i}^{(\pi)}} [\rho_{s,a}] - \mathbb{E}_{(s,a) \sim \Lambda_{-i}^{(\pi)}} [\rho_{s,a}] \right| \le \epsilon. \tag{3}$$

We let  $\Pi_{DP,\epsilon}$  denote the policies satisfying (3).

**Remark 2.4.** We consider two subpopulations for simplicity; our techniques extend to multiple subpopulations.

**Existence and determinism.** We briefly discuss the existence of deterministic solutions to (2). Unconstrained MDPs always have a deterministic optimal policy (Sutton & Barto, 2018). With a fairness constraint, this result no longer holds:

**Theorem 2.5.** There exists an MDP such that  $\Pi_{DP} = \emptyset$ . There exists an MDP such that  $\pi^*$  in (2) is not deterministic.

We give a proof in Appendix B. For the following special case, we can prove existence of fair policies:

**Definition 2.6.** The agent rewards are *state-independent* if  $\rho_{s,a} = \tilde{\rho}_a$  for all  $s \in S$  and for some  $\tilde{\rho} \in \mathbb{R}^{|A|}$ .

Intuitively, this property captures settings where the decision-maker uses the state to choose actions (e.g., ability to repay), but the outcomes for the individuals only depend on whether the preferred action is taken (e.g., a loan offer). Our example  $M_{\rm loan}$  has state-independent agent rewards.

**Theorem 2.7.** *If the agent rewards are state-independent, then* (2) *has a solution.* 

*Proof.* Clearly, any policy  $\pi$  such that  $\pi_{s,a} = \tilde{\pi}_a$  for all  $s \in S$  and some  $\tilde{\pi} \in \mathbb{R}^{|A|}$ , satisfies  $\pi \in \Pi_{\mathrm{DP}}$ .

# 3. Model-Based Algorithm

We describe a model-based algorithm for solving (2), which has strong theoretical guarantees (i.e., it solves (2) exactly in polynomial time). On the other hand, it makes strong assumptions—i.e., that M has finite state and action spaces, and furthermore satisfies a separability property saying that the sensitive attribute  $z \in Z$  does not change over time:

**Definition 3.1.** An MDP with states  $S = Z \times \tilde{S}$  is *separable* if the transitions satisfy  $P_{(z,\tilde{s}),a,(z',\tilde{s}')} = \delta_{z,z'} \tilde{P}_{\tilde{s},a,\tilde{s}'}$ , where  $\delta_{z,z'} = \mathbb{I}[z=z']$  is the Kronecker delta and  $\tilde{P} \in \mathbb{R}^{|\tilde{S}| \times |A| \times |\tilde{S}|}$  is a transition matrix.

That is, the transitions do not affect z. This property is satisfied by many sensitive attributes (e.g., race and gender).

**Background.** When the policy class  $\Pi$  is unconstrained, then the optimal policy is deterministic, and can be expressed as the function

$$\pi^*(s) = \underset{a \in A}{\arg\max} \left( R_{s,a} + \gamma \sum_{s' \in S} P_{s,a,s'} V_{s'}^* \right)$$

where the value function  $V^* \in \mathbb{R}^{|S|}$  is the unique solution to the Bellman equation (Sutton & Barto, 2018):

$$V_s^* = \max_{a \in A} \left\{ R_{s,a} + \gamma \sum_{s' \in S} P_{s,a,s'} V_{s'}^* \right\}.$$
 (4)

Furthermore,  $V^*$  is the solution to the following linear program (LP) (Sutton & Barto, 2018):

$$\underset{V \in \mathbb{R}^{|S|}}{\arg\min} \langle D, V \rangle \tag{5}$$

subj. to 
$$V_s \ge R_{s,a} + \gamma \sum_{s' \in S} P_{s,a,s'} V_{s'} \ (\forall s \in S, a \in A).$$

### Algorithm 1 Model-based algorithm.

**Input:** Separable MDP M

Compute the solution  $\lambda^*$ ,  $c^*$  to the linear program

$$\underset{\lambda \in \mathbb{R}^{|S| \times |A|}, c \in \mathbb{R}}{\operatorname{arg \, max}} (1 - \gamma)^{-1} \sum_{s \in S} \sum_{a \in A} \lambda_{s,a} R_{s,a}$$

$$\operatorname{subj. to } \sum_{a \in A} \lambda_{s',a} = (1 - \gamma) D_{s'} + \gamma \sum_{s \in S} \sum_{a \in A} \lambda_{s,a} P_{s,a,s'}$$

$$(\forall s' \in S)$$

$$p_z^{-1} \sum_{\tilde{z} \in \tilde{S}} \sum_{a \in A} \lambda_{(z,\tilde{s}),a} \rho_{(z,\tilde{s}),a} = c \qquad (\forall z \in Z)$$

Output: Policy 
$$\pi_{s,a}^* = \frac{\lambda_{s,a}^*}{\sum_{a' \in A} \lambda_{s,a'}^*}$$

**Algorithm.** When we require that  $\pi \in \Pi_{DP}$ , then the Bellman equation (4) may no longer hold. Thus, if we add this constraint to (5), then (5) may become unsatisfiable. Instead, our approach is based on the dual of (5) (Altman, 1999). In particular, the objective and first set of constraints of the LP in Algorithm 1 form the dual. <sup>2</sup>

The last set of constraints in the LP in Algorithm 1 encodes demographic parity. These constraints exploit the separable structure of the underlying MDP. In particular, the component z of an initial state  $s=(z,\tilde{s})$  does not change over time, so the value of z for s equals the value of z for the initial state  $s_0 \sim D$ . Thus, randomly sampling a state  $s \sim D_z^{(\pi)}$  is equivalent to randomly sampling

$$s \sim D^{(\pi)} \mid \exists \tilde{s} \in \tilde{S} . s = (z, \tilde{s}).$$

Expanding the conditional probability, the probability of sampling  $s \sim D_z^{(\pi)}$  is

$$\frac{D_s^{(\pi)}\mathbb{I}[\exists \tilde{s} \in \tilde{S} . s = (z, \tilde{s})]}{p_z}, \quad \text{ where } p_z = \sum_{\tilde{s} \in \tilde{S}} D_{(z, \tilde{s})}.$$

It follows that

$$\mathbb{E}_{(s,a)\sim\Lambda_z^{(\pi)}}[\rho_{s,a}] = p_z^{-1} \sum_{\tilde{s}\in\tilde{S}} \sum_{a\in A} \lambda_{(z,\tilde{s}),a} \rho_{s,a}. \tag{6}$$

The last set of constraints in the LP in Algorithm 1 uses (6) to encode demographic parity.

**Theorem 3.2.** Given a separable MDP M, Algorithm 1 a solution  $\pi^*$  to (2) if and only if (2) is satisfiable.

We give a proof in Appendix C. Since we can solve an LP in polynomial time, Algorithm 1 runs in polynomial time.

# 4. Model-Free Algorithm

Our model-based algorithm makes strong assumptions about the given MDP—i.e., separability and finite state and action spaces. We propose a model-free algorithm for solving (2) that relaxes all these assumptions. Our algorithm learns policies that satisfy demographic parity with an  $\epsilon$  tolerance—i.e., the set of policies  $\Pi_{\mathrm{DP},\epsilon}$ . We need this tolerance since our model-free algorithm can only estimate the agent rewards  $\rho$ . Then, our algorithm is based on formulating (2) as the following optimization problem:

$$\underset{\pi \in \Pi, c \in \mathbb{R}}{\arg \max} R^{(\pi)} \quad \text{subj. to } |\rho_{\text{maj}}^{(\pi)} - \rho_{\text{min}}^{(\pi)}| \le \epsilon, \quad (7)$$

where  $\rho_z^{(\pi)}=\mathbb{E}_{(s,a)\sim\Lambda_z^{(\pi)}}[\rho]$ . Note that this optimization problem is non-convex. Thus, unlike our model-based algorithm, this algorithm may converge to a local optimum.

**Background.** Our algorithm relies on the cross-entropy (CE) method (Mannor et al., 2003; Hu et al., 2012), which is a general heuristic for solving optimization problems. Suppose our policies  $\pi_{\theta} \in \Pi$  are parameterized by  $\theta \in \Theta$ , and let a family F of probability distributions over  $\Theta$  parameterized by  $V \subseteq \mathbb{R}^d$ . In general, we use  $\theta$  and  $\pi_{\theta}$  interchangeably, e.g.,  $R^{(\theta)} = R^{(\pi_{\theta})}$ . In the unconstrained setting, CE aims to solve the following optimization problem:

$$v^* = \operatorname*{arg\,max}_{v \in V} \mathbb{E}_v[R^{(\theta)}],\tag{8}$$

where we have used the notation  $\mathbb{E}_v = \mathbb{E}_{\theta \sim f_v}$ . In other words, it aims to compute a distribution  $f_{v^*}$  that places high probability mass on  $\theta$  with high expected cumulative reward  $R^{(\theta)}$ . Then, it returns a sample  $\theta \sim f_{v^*}$ .

To solve (8), CE starts with initial parameters  $v_0 \in V$ . Then, on each iteration, it updates the current parameters  $v_k$  to move "closer" to  $v^*$ . More precisely, the update is

$$v_{k+1} = \arg\max_{v \in V} D_{\text{KL}}(g_{k+1} \parallel f_v)$$
 (9)

$$g_{k+1}(\theta') = \alpha \frac{R^{(\theta')} \mathbb{I}[R^{(\theta')} \ge \gamma_k] f_{v_k}(\theta')}{\mathbb{E}_{v_k}[R^{(\theta)} \mathbb{I}[R^{(\theta)} \ge \gamma_k]]} + (1 - \alpha) f_{v_k}(\theta')$$

where  $\gamma_k$  satisfies  $\Pr_{v_k}[R^{(\theta)} \geq \gamma_i] = \mu$ . Here,  $\alpha, \mu \in (0,1)$  are hyperparameters. Intuitively, the first term of  $g_i$  upweights  $\theta'$  with large values of  $R^{(\theta')}$  compared to  $f_{v_k}$ , both by directly weighting the probability of  $\theta'$  by  $R^{(\theta')}$ , and furthermore by placing zero probability mass on the bottom  $1-\mu$  fraction of the  $\theta'$ . The second term of  $g_k$  is a "smoothing" term that makes the update incremental.

To enable efficient optimization of (9), we assume that F is a (natural) exponential family.

**Definition 4.1.** A family  $\mathcal{F}$  of distributions over  $\Theta \subseteq \mathbb{R}^d$ 

<sup>&</sup>lt;sup>2</sup>Some variables are rescaled compared to the actual dual.

is an *exponential family* if, for a continuous  $\Gamma: \Theta \to \mathbb{R}^d$ ,

$$f_v(\theta) = \frac{1}{Z(\theta)} e^{v^\top \Gamma(\theta)}$$
 where  $Z(\theta) = \int e^{v^\top \Gamma(\theta)} d\theta$ .

In this case, it can be shown that (Hu et al., 2012)

$$v_{k+1} = m^{-1}(\eta_{k+1})$$

$$\eta_{k+1} = \alpha \frac{\mathbb{E}_{v_k}[R^{(\theta)}\mathbb{I}[R^{(\theta)} \ge \gamma_k]\Gamma(\theta)]}{\mathbb{E}_{v_k}[R^{(\theta)}\mathbb{I}[R^{(\theta)} \ge \gamma_k]]} + (1 - \alpha)\eta_k$$
(10)

where  $m(v) = \mathbb{E}_v[\Gamma(\theta)]$  is the moment map.

The CE algorithm approximates (10) by sampling rollouts  $\zeta = ((s_0, a_0), ..., (s_{T-1}, a_{T-1}))$  of length T according to policy  $\pi_{\theta}$ . Then, it computes estimate  $\hat{R}^{(\theta)} \approx R^{(\theta)}$ , where

$$\hat{R}^{(\theta)} = \frac{1}{m} \sum_{i=1}^{m} \hat{R}(\zeta^{(i)}) \quad \text{ where } \hat{R}(\zeta) = \sum_{t=0}^{T-1} \gamma^{t} R_{s_{t}, a_{t}},$$

and where  $\zeta^{(1)},...,\zeta^{(m)}$  are m sampled rollouts.

To estimate  $\eta_{k+1}$ , it takes n samples  $\theta^{(1)},...,\theta^{(n)} \sim f_v$  and computes  $\hat{R}^{(\theta^{(i)})}$  for each  $i \in [n]$ . It then ranks the  $\theta^{(i)}$  in decreasing order of  $\hat{R}^{(\theta^{(i)})}$ , and discards all but the top  $n' = \lceil n\mu \rceil$ . Now, it estimates the numerator in  $\eta_{k+1}$  as

$$\mathbb{E}_{v_k}[R^{(\theta)}\mathbb{I}[R^{(\theta)} \ge \gamma_k]\Gamma(\theta)] \approx \frac{1}{n} \sum_{i=1}^{n'} \hat{R}^{(\theta^{(i)})}\Gamma(\theta^{(i)}).$$

The denominator in  $\eta_{k+1}$  is estimated similarly.

Algorithm 2 computes this estimate of the update (10) assuming the condition on Line 11 is satisfied (as we discuss below, the check is needed to enforce the constraint in (7)). Line 6 of Algorithm 2 computes the estimates  $\hat{R}^{(\theta^{(i)})}$  for samples  $\theta^{(i)} \sim f_{v_k}$  for  $i \in [n]$ , and Line 14 estimates  $\eta_{k+1}$ . On Line 6 & 7, the notation  $\stackrel{\sim m}{\longleftarrow}$  means to estimate using m samples (in this case, rollouts  $\zeta^{(1)}, ..., \zeta^{(m)}$ ).

Finally, we use a constrained cross-entropy (CCE) method, which extends CE to handle constraints (Wen & Topcu, 2018). Intuitively, CCE prioritizes policies where the constraint in (7) is closer to holding, unless the constraint holds, in which case CCE prioritizes policies with higher expected cumulative reward. In particular, Algorithm 2 imposes this constraint by checking if a sufficient fraction of the  $\theta$  satisfies the constraint  $\hat{\epsilon}^{(\theta)} \leq \epsilon$  in Line 11, where  $\hat{\epsilon}^{(\theta)}$  is estimated from samples. As discussed below,  $\tilde{\epsilon}$  is used in place of  $\epsilon$  to enforce the constraint even though  $\hat{\epsilon}^{(\theta)}$  is inexact.

**Algorithm.** A key challenge to applying CCE is that it relies on estimates  $\hat{\epsilon}^{(\pi)}$  of  $\epsilon^{(\pi)}$ . These estimates are inexact for two reasons: (i) they are estimated from samples, and (ii) they are estimated based on a finite time horizon (whereas

Algorithm 2 Model-free algorithm.

```
1: Input: MDP M, Iters r, Parameter samples n, Top n',
         Rollout samples m, Smoothing \alpha, Tolerance \sigma
  2: \hat{n} \leftarrow \vec{0}
  3: for k \in [1, ..., r] do
               Sample \theta^{(1)}, ..., \theta^{(n)} \sim f_{m^{-1}(\hat{n})}
               for i \in [1, ..., n] do
                   \begin{array}{c} \hat{R}^{(\theta^{(i)})} \xleftarrow{\sim m} R^{(\theta^{(i)})} \\ \hat{\epsilon}^{(\theta^{(i)})} \xleftarrow{\sim m} |\rho_{\text{maj}}^{(\theta^{(i)})} - \rho_{\text{min}}^{(\theta^{(i)})}| \end{array}
  6:
  7:
  8:
              Sort \{\theta^{(i)}\}_{i=1}^n in increasing \hat{\epsilon}^{(\theta^{(i)})}
  9:
              i' \leftarrow \text{Largest } i \text{ such that } \hat{\epsilon}^{(\bar{\theta}^{(i)})} \leq (1 - \sigma)\epsilon
10:
              if n' \leq i' then
11:
                    Sort \{\theta^{(i)}\}_{i=1}^{i'} in decreasing \hat{R}^{(\theta^{(i)})}
12:
13:
              \hat{\eta} \leftarrow \alpha \cdot \frac{\frac{1}{n} \sum_{i=1}^{n'} \hat{R}^{(\theta^{(i)})} \Gamma(\theta^{(i)})}{\frac{1}{n} \sum_{i=1}^{n'} \hat{R}^{(\theta^{(i)})}} + (1 - \alpha) \cdot \hat{\eta}
16: Output: Policy \pi_{\hat{\theta}}, where \hat{\theta} \sim f_{m^{-1}(\hat{\eta})}
```

 $\epsilon^{(\pi)}$  is defined for an infinite time horizon). To account for this error, we use  $(1-\sigma)\epsilon$  (where  $\sigma\in(0,1)$ ) in place of  $\epsilon$  when checking the constraint on Line 11 of Algorithm 2.

We provide two guarantees for Algorithm 2. First, the errors in (i) estimating whether  $\epsilon^{(\pi)} \leq \epsilon$  and (ii) estimating  $R^{(\pi)}$  can be made arbitrarily small by making m and T large.

**Theorem 4.2.** Assume that  $R_{max}$  be an upper bound on R (i.e.,  $||R||_{\infty} = R_{max}$ ) and on  $\rho$ . Let  $\delta \in \mathbb{R}_+$  and  $\sigma \in (0, 1/2]$  be given, and suppose that

$$m \ge \frac{32R_{max}(1-\gamma)\log(6/\delta)}{\sigma^2 \epsilon^2}$$
$$T \ge \log \frac{4R_{max}}{\sigma^2 \epsilon(1-\gamma)}.$$

*Then, for any*  $\pi$ *, with probability at least*  $1 - \delta$ *, we have* 

$$(\hat{\epsilon}^{(\pi)} \leq \tilde{\epsilon}) \Rightarrow (\pi \in \Pi_{DP,\epsilon})$$

where  $\tilde{\epsilon} = (1 - \sigma)\epsilon$ , and

$$|R^{(\pi)} - \hat{R}^{(\pi)}| \le \frac{\sigma\epsilon}{2}.$$

Note that  $\hat{R}^{(\pi)}$  and  $\hat{\rho}_z^{(\pi)}$  (for each  $z \in Z$ ) must be estimated with independently sampled rollouts.

Second, we assume that Algorithm 2 solves the problem truncated to T steps exactly, and then consider how the solution compares to the untruncated problem (7).

**Theorem 4.3.** Assume that  $R_{max}$  is an upper bound on R (i.e.,  $||R||_{\infty} = R_{max}$ ) and on  $\rho$ . Let  $\sigma \in (0, 1/2]$  be given, and suppose that T is as in Theorem 4.2, and let

$$\pi^* = \argmax_{\pi \in \Pi_{D^p, (1-\sigma)^2\epsilon}} R^{(\pi)} \quad \text{ and } \quad \tilde{\pi}^* = \argmax_{\pi \in \tilde{\Pi}_{D^p, (1-\sigma)\epsilon}} \tilde{R}^{(\pi)},$$

where  $\tilde{R}^{(\pi)} = \sum_{t=0}^{T-1} \gamma^t \langle R, \Lambda^{(\pi,t)} \rangle$  is truncated to T steps, similarly  $\tilde{\rho}_z^{(\pi)} = \sum_{t=0}^{T-1} \gamma^t \langle R, \Lambda_z^{(\pi,t)} \rangle$ , and  $\tilde{\Pi}_{DP,\tilde{\epsilon}}$  is the set of  $\pi$  satisfying  $|\tilde{\rho}_{maj}^{(\pi)} - \tilde{\rho}_{min}^{(\pi)}| \leq \tilde{\epsilon}$ . Then,  $\tilde{\pi}^* \in \Pi_{DP,\epsilon}$ , and

$$|R^{(\pi^*)} - R^{(\tilde{\pi}^*)}| \le \frac{\sigma^2 \epsilon}{2}.$$

We give proofs for both theorems in Appendix D. Note that in Theorem 4.3, we have restricted (7) to guarantee fairness with tolerance  $(1 - \sigma)^2 \epsilon$  instead of  $\epsilon$ . This restriction is needed since Algorithm 2 must be conservative to ensure that the error due to the truncation error. This gap can be made arbitrarily small by taking  $\sigma$  sufficiently small.

# 5. Fairness Ignoring Dynamics

To demonstrate the importance of accounting for dynamics, we compare to two baseline algorithms that ignore dynamics when constraining fairness. More precisely, these algorithms solve an optimization problem of the form

$$\pi^* = \operatorname*{arg\,max}_{\pi \in \Pi_{\mathsf{DP}}^{\mathsf{OP}}} R^{(\pi)},$$

where  $\Pi^0_{\mathrm{DP}}$  does not take into account the MDP dynamics—i.e.,  $\Pi^0_{\mathrm{DP}}$  does not account for how actions affect the distribution of states  $D^{(\pi,t)}$  at future time steps t>0. We consider two algorithms, each using a different choice of  $\Pi^0_{\mathrm{DP}}$ .

The first algorithm optimistically pretends that actions do not affect the state distribution—i.e.,  $D^{(\pi,t)}$  does not change over time. This captures the supervised learning setting. Compared to our algorithm, this algorithm may learn a policy that is unfair but achieves higher reward.

The second algorithm conservatively assumes  $D^{(\pi,t)}$  can change arbitrarily on each step. Like our algorithm, this one learns a fair policy, but it may achieve much lower reward.

**Optimistic assumptions.** We can optimistically assume that the state distribution does not change over time, i.e.,

$$D^{(\pi,t)} = D \quad (\forall t > 0). \tag{11}$$

Given this assumption, the time-discounted distribution over states equals D regardless of the policy  $\pi$ —i.e.,  $D^{(\pi)} = D$  for any  $\pi$ . Then, we can let  $\pi^*$  be the solution to

$$\underset{\pi \in \Pi, c \in \mathbb{R}}{\arg \max} R^{(\pi)} \tag{12}$$

subj. to 
$$\mathbb{E}_{s \sim D_z} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right] = c \quad (\forall z \in Z),$$

where  $D_z = D \mid \exists \tilde{s} \in \tilde{S}$  .  $s_0 = (z, \tilde{s})$ . We can solve (12) using a straightforward modification of Algorithm 2.

**Theorem 5.1.** Assuming that (11) holds for the MDP M, then the solution of (12) is a solution to (2).

This theorem follows straightforwardly—in particular, the objective in (12) is the same as that in (2), and if (11) holds, then the constraint in (12) is equivalent to the demographic parity constraint. Of course, if (11) fails to hold, then we cannot provide any guarantees about  $\pi^*$ .

Conservative assumptions. We can conservatively assume that  $D^{(\pi,t)}$  for every t can be arbitrary, and to conservatively require that demographic parity holds for every possible sequence  $D^{(\pi,t)}$ . We focus on finite S. Then, we restrict to policies  $\pi$  that satisfy

$$\mathbb{E}_{s \sim D'_{\text{maj}}} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right] = \mathbb{E}_{s \sim D'_{\text{min}}} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right]$$
(13)
$$(\forall D' \in \Delta^{|S|}),$$

where  $D_z' = D' \mid \exists \tilde{s} \in \tilde{S}$ .  $s = (z, \tilde{s})$ , and  $\Delta^n$  is the standard n-simplex. Note that  $D_z'$  is conditioned on  $s = (z, \tilde{s})$  (i.e., the current state has sensitive attribute z) instead of  $s_0 = (z, \tilde{s})$  (i.e., the initial state has sensitive attribute z); if M is separable, these two conditions are equivalent. Finally, note that  $D_z'$  is undefined if the conditional has zero probability according to D'; we implicitly omit such D' from the universal quantification in (13).

The difficulty in enforcing (13) is handling the universal quantification over  $D' \in \Delta^{|S|}$ . In fact, we can equivalently enforce that the one-step rewards are independent of the state. Thus, we can solve the optimization problem

$$\underset{\pi \in \Pi, c \in \mathbb{R}}{\arg \max} R^{(\pi)}$$
 (14)  
subj. to 
$$\sum_{a \in A} \pi_{s,a} \rho_{s,a} = c \quad (\forall s \in S).$$

We use  $\pi_0$  to denote the solution to (14). When S is finite, we can solve (14) using a modification of Algorithm 2; however, to the best of our knowledge, the conservative approach is in general intractable when S is continuous. We have made conservative assumptions about  $D^{(\pi,t)}$ , so the solution to (14) satisfies demographic parity.

**Theorem 5.2.** The solution  $\pi^*$  of (14) satisfies  $\pi^* \in \Pi_{DP}$ .

We give a proof in Appendix E. Note that while the solution  $\pi^*$  of (14) is guaranteed to satisfy demographic parity, it may be suboptimal compared to taking into account the dynamics in the fairness constraint.

For separable M, (14) is equivalent to the LP

$$\underset{\lambda \in \mathbb{R}^{|S| \times |A|}, c \in \mathbb{R}^{|S|}}{\operatorname{arg\,max}} (1 - \gamma)^{-1} \sum_{s \in S} \sum_{a \in A} \lambda_{s,a} R_{s,a}$$

$$\text{subj. to } \sum_{a \in A} \lambda_{s',a} = (1 - \gamma) D_{s'} + \gamma \sum_{s \in S} \sum_{a \in A} \lambda_{s,a} P_{s,a,s'}$$

$$(\forall s' \in S)$$

$$\sum_{a \in A} \lambda_{(z,\tilde{s}),a} \rho_{(z,\tilde{s}),a} = c \qquad (\forall z \in Z, \ \forall \tilde{s} \in \tilde{S})$$

so we can return the policy  $\pi_{s,a} = \lambda_{s,a}/\sum_{a'\in A} \lambda_{s,a'}$ . The proof is analogous to the proof of Theorem 3.2.

## 6. Reinforcement Learning

We discuss extensions to the setting where parts of the MDP are initially unknown, and the goal is to ensure fairness while learning these quantities. The approaches we propose are naïve, leaving room for future work. We consider policies that satisfy demographic parity with an  $\epsilon$  tolerance—i.e., the set of policies  $\Pi_{DP,\epsilon}$ . Our proposed approaches rely on the policy  $\pi_0$  defined in (14), which is learned using conservative assumptions about the MDP dynamics.

Unknown dynamics. We propose an approach to fairness when the transitions P are unknown. Our goal is to ensure that with high probability, fairness holds for all time including during learning. We consider the episodic case where the system is reset after a fixed number of steps T, and take  $\gamma=1$ . That is, a finite sequence of interactions is performed repeatedly—e.g., each new loan applicant is a new episode. We assume there are a fixed total number of episodes N, and the goal is to perform well on average; the doubling trick can be used to generalize to unknown or unbounded N (see p. 99 of (Lattimore & Szepesvári)).

We use explore-then-commit (Lattimore & Szepesvári). First, we explore using the conservative policy  $\pi_0$  for  $N_0$  episodes. Then, we estimate P using the observed stateaction-state tuples (s, a, s') (i.e., transition to s' upon taking action a in state s):

$$\hat{P}_{s,a,s'} = \frac{\# \text{ observed tuples } (s,a,s')}{\# \text{ observed tuples } (s,a,s'') \text{ for some } s'' \in S}.$$

Finally, for the remaining  $N-N_0$ , it uses the optimal policy  $\hat{\pi}$  computed as if  $\hat{P}$  is the true transition matrix.

We assume that  $\pi_0$  explores all state-action pairs. In particular, let  $D^{(\pi)}=\frac{1}{T}\sum_{t=0}^{T-1}D^{(\pi,t)}$  and  $\Lambda_{s,a}^{(\pi)}=D_s^{(\pi)}\pi_{s,a}$ , where  $D^{(\pi,t)}$  is defined as before. Then, we assume that

$$\Lambda_{s,a}^{(\pi_0)} \ge \lambda_0 > 0 \quad (\forall s \in S, \ a \in A)$$

for some constant  $\lambda_0$ . We prove a bound on the *regret* 

$$\mathcal{R}(N) = \mathbb{E}\left[\sum_{n=1}^{N} R^{(\pi^*)} - R^{(\pi_n)}\right],$$

where the expectation is taken over the randomness of the observed tuples (s, a, s'),  $\pi^*$  is the optimal policy for known P that satisfies  $\pi^* \in \Pi_{\mathsf{DP}, \epsilon/4}$ , and

$$\pi_n = \begin{cases} \pi_0 & \text{if } n \le N_0 \\ \hat{\pi} & \text{otherwise} \end{cases}$$

is the policy our algorithm uses on episode n. Furthermore, we show that for a given  $\delta \in \mathbb{R}_+$ , we have  $\pi_n \in \Pi_{DP,\epsilon}$  for every  $n \in [N]$  with probability at least  $1 - \delta$ .

**Theorem 6.1.** Let  $\epsilon, \delta \in \mathbb{R}_+$  be given. Assume that  $R_{max}$  be an upper bound on R (i.e.,  $||R||_{\infty} = R_{max}$ ) and on  $\rho$ . Let

$$N_0 = \frac{128T^2 \cdot |S|^2 \cdot R_{max}^2 \cdot \log(2|S|^2|A|/\delta)}{\lambda_0^2 \epsilon^2}.$$

Let  $\hat{M}=(S,A,D,\hat{P},R,T)$ , and  $\hat{\pi}$  be the optimal policy for  $\hat{M}$  in  $\hat{\Pi}_{DP,\epsilon/2}$  (i.e., the set of policies satisfying demographic parity for  $\hat{M}$ ). Let M=(S,A,D,P,R,T), and  $\pi^*$  be optimal for M in  $\Pi_{DP,\epsilon/4}$ . Then,  $\hat{\pi}\in\Pi_{DP,\epsilon}$ , and  $R^{(\pi^*)}-R^{(\hat{\pi})}\leq \epsilon$ , where  $R^{(\pi)}$  is defined according to M.

We give a proof in Appendix G. Note that there is a gap between the fairness constraint of  $\pi^*$  (which is in  $\Pi_{DP,\epsilon/4}$ ) and that of  $\hat{\pi}$  (which is only in  $\Pi_{DP,\epsilon}$ )—i.e., we can only guarantee performance compared to a policy that satisfies a stricter level of fairness. Choosing  $\epsilon = N^{-2/3}$ , we have:

**Corollary 6.2.** For any  $\delta \in \mathbb{R}_+$ , we have regret  $\mathcal{R}(N) = O(N^{2/3}\log(1/\delta))$  with probability at least  $1 - \delta$ .

**Unknown initial distribution.** Suppose that the initial distribution D is unknown. We consider the non-episodic setting—i.e., we cannot restart. In this setting, initially use the conservative policy  $\pi_0$ , which is fair regardless of D. Then, we can improve performance once the Markov chain induced by  $\pi_0$  is *mixed*—i.e., close to its stationary distribution (Kearns & Singh, 2002; Even-Dar et al., 2005):

**Definition 6.3.** A policy  $\pi$  is ergodic if there exists  $d^{(\pi)} \in \Delta^{|S|}$  (the stationary distribution of  $\pi$ ) such that  $P^{(\pi)}d^{(\pi)} = d^{(\pi)}$ , and for all  $d \in \Delta^{|S|}$ ,  $\lim_{t \to \infty} (P^{(\pi)})^t d = d^{(\pi)}$ . Given  $\operatorname{ergodic} \pi$  and  $\epsilon \in \mathbb{R}_+$ , the  $\epsilon$  mixing time is  $T \in \mathbb{N}$  such that for all  $d \in \Delta^{|S|}$ ,  $\|(P^{(\pi)})^T d - d^{(\pi)}\|_{\infty} < \epsilon$ .

Let  $T_0$  be the  $\epsilon_0$  mixing time of  $\pi_0$ . While  $D^{(\pi_0,T_0)}$  is unknown, we know that it satisfies

$$||D^{(\pi_0,T_0)} - d^{(\pi_0)}||_{\infty} \le \epsilon_0.$$

Finally, we can run Algorithm 1 as if  $d^{(\pi_0)}$  is the initial distribution, to obtain a policy  $\tilde{\pi}$ . As long as  $\epsilon$  is sufficiently

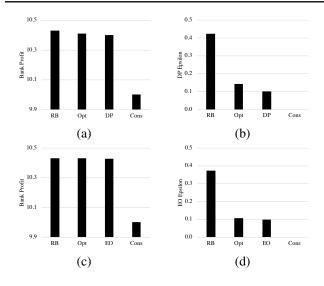


Figure 1. Demographic parity (a) objective value, (b) constraint value, and equal opportunity (c) objective value, (d) constraint value, for race-blind (RB), demographic parity (DP) or equal opportunity (EO), optimistic (Opt), and conservative (Cons).

small, we can show that  $\tilde{\pi}$  achieves reward close to the true optimal policy  $\pi^*$  (with initial distribution  $D^{(\pi_0, T_0)}$ ).

**Theorem 6.4.** Let  $\epsilon \in \mathbb{R}_+$  be given. Assume that  $R_{max}$  upper bounds R (i.e.,  $\|R\|_{\infty} = R_{max}$ ) and  $\rho$ , and that the  $\epsilon_0$  mixing time of  $\pi_0$  is  $T_0$ , where  $\epsilon_0 = \frac{(1-\gamma)\epsilon}{8|S| \cdot R_{max}}$ . Let  $\tilde{M} = (S,A,d^{(\pi_0)},P,R,\gamma)$ , and let  $\tilde{\pi}$  be the optimal policy for  $\tilde{M}$  in  $\tilde{\Pi}_{DP,\epsilon/2}$  (i.e., the set of policies satisfying demographic parity for  $\tilde{M}$ ). Similarly, let  $M = (S,A,D^{(\pi_0,T_0)},P,R,\gamma)$ , and let  $\tilde{\pi}^*$  be the optimal policy for M in  $\Pi_{DP,\epsilon/4}$ . Then,  $\tilde{\pi} \in \Pi_{DP,\epsilon}$ , and  $R^{(\tilde{\pi}^*)} - R^{(\tilde{\pi})} \leq \epsilon$ , where  $R^{(\pi)}$  is defined according to M.

We give a proof in Appendix F. In other words, we act close to optimally once the Markov chain has mixed. However, Theorem 6.4 only bounds the reward compared to the optimal policy  $\tilde{\pi}^*$  on  $D^{(\pi_0,T_0)}$ . Ideally, we would bound the reward compared to the optimal policy  $\pi^*$  trained with known initial dynamics. However, even though we are acting optimally after the first  $T_0$  steps,  $\pi_0$  may have moved the system into a suboptimal distribution  $D^{(\pi_0,T_0)}$  compared to  $D^{(\pi^*,T_0)}$ . Thus, we cannot prove such a bound.

Finally, as before, there is a gap between the fairness constraint of  $\pi^*$  (in  $\Pi_{DP,\epsilon/4}$ ) and that of  $\hat{\pi}$  (only in  $\Pi_{DP,\epsilon}$ ).

### 7. Experiments

**MDP parameters.** We run simulations using our loan example from Section 2. We estimated parameters based on FICO score data (Hardt et al., 2016). We consider the majority subpopulation to be Whites, and the minority sub-

population to be Blacks, Hispanics, and Asians. For the initial distribution  $p_0$ , we first fit parameters the parameters of the prior  $\text{Beta}(\alpha_z,\beta_z)$  based on the data. Then, we take a fixed number of steps  $T_z$  using action a=1 (i.e., offer loan) to force exploration. We choose  $T_{\text{maj}} > T_{\text{min}}$  to capture the idea that less data is available for minorities. We also estimate the probability  $p_Z$  of being a minority from the data. Similar to (Hardt et al., 2016), we choose I so the bank makes a profit on the average applicant. We manually choose the remaining parameters  $\lambda$ ,  $\tau$ ,  $T_{\text{maj}}$ , and  $T_{\text{min}}$  based on intuition. We give exact parameter values in Appendix I.

**Experimental setup.** We ran Algorithm 2 to learn fair policies for both the demographic parity and equal opportunity constraints, using fairness threshold  $\epsilon = 0.1$  (i.e., policy classes  $\Pi_{DP,\epsilon}$  and  $\Pi_{EO,\epsilon}$ ). For each constraint, we also use the optimistic and conservative algorithms described in Section 5. Note that the conservative algorithms we described do not apply since our state space is infinite. However, since our agent rewards are state-independent, the conservative assumption is in fact equivalent to optimizing over stateindependent policies—i.e., those of the form  $\pi_{s,a} = \tilde{\pi}_a$ , where  $\tilde{\pi} \in \mathbb{R}^{|A|}$ . Thus, we can apply a modified version of Algorithm 2 where we only learn state-independent policies. We also run a race-blind algorithm, which is unconstrained but where  $\pi$  ignores the sensitive attribute  $z \in Z$ . Note that the optimal policy is race-blind, since the portion  $\alpha, \beta$  of the state is a sufficient statistic, so it captures all information needed to determine whether to offer a loan.

**Results.** For demographic parity, Figure 1 (a) shows the reward achieved for the bank, and (b) shows the value of the fairness constraint—i.e., the smallest value of  $\epsilon$  for which  $\pi \in \Pi_{DP,\epsilon}$ . As expected, race-blind achieves the highest reward (10.43), followed by the optimistic algorithm (10.41), and then Algorithm 2 (10.40). Finally, the conservative algorithm performs substantially worse than the others (10.00). However, race-blind achieves a very poor constraint value (0.42), as does the optimistic algorithm (0.14), which performs performs 43% worse than Algorithm 2 (0.10). The conservative algorithm achieves constraint value 0.

For equal opportunity, Figure 1 (c) shows the bank reward, and (d) shows the value of the constraint. The bank's rewards are essentially the same for the race-blind algorithm, optimistic algorithm, and Algorithm 2 (10.43), but is substantially worse for the conservative algorithm (10.00). As with demographic parity, the constraint value for race-blind (0.37) is substantially worse than the others, but in this case optimistic (0.11) is fairly close to Algorithm 2 (0.10). The conservative algorithm achieves constraint value 0.

**Discussion.** Our results show that imposing demographic parity slightly reduces the bank's reward, but substantially

increases fairness compared to the race-blind and optimistic algorithms. Recall that the optimistic algorithm models supervised learning—thus, our results show the importance of accounting for dynamics when ensuring fairness. We find similar (but weaker) trends for equal opportunity. Like prior work (Hardt et al., 2016), we find that demographic parity reduces the bank's rewards more than equal opportunity.

Unlike the static case (Hardt et al., 2016), our model has dynamic parameters. Time series data would be needed to estimate these parameters; instead, we choose them manually. Also, (Hardt et al., 2016) uses the empirical CDF of the distribution over repayment probabilities  $p_0$  (whereas we assumed  $p_0$  is a Beta distribution). However, our goal is to understand the consequences of ignoring dynamics, not to study a real-world scenario.

### References

- Achiam, J., Held, D., Tamar, A., and Abbeel, P. Constrained policy optimization. In *ICML*, 2017.
- Altman, E. Constrained Markov decision processes, volume 7. CRC Press, 1999.
- Bastani, O., Pu, Y., and Solar-Lezama, A. Verifiable reinforcement learning via policy extraction. In *NIPS*, 2018.
- Calders, T., Kamiran, F., and Pechenizkiy, M. Building classifiers with independency constraints. In *Data min*ing workshops, 2009. ICDMW'09. IEEE international conference on, pp. 13–18. IEEE, 2009.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226. ACM, 2012.
- Elzayn, H., Jabbari, S., Jung, C., Kearns, M., Neel, S., Roth, A., and Schutzman, Z. Fair algorithms for learning in allocation problems. 2019.
- Even-Dar, E., Kakade, S. M., and Mansour, Y. Experts in a markov decision process. In *Advances in neural information processing systems*, pp. 401–408, 2005.
- Hardt, M., Price, E., Srebro, N., et al. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pp. 3315–3323, 2016.
- Hashimoto, T. B., Srivastava, M., Namkoong, H., and Liang,P. Fairness without demographics in repeated loss minimization. In *ICML*, 2018.
- Hu, J., Hu, P., and Chang, H. S. A stochastic approximation framework for a class of randomized optimization algorithms. *IEEE Transactions on Automatic Control*, 57(1): 165–178, 2012.

- Jabbari, S., Joseph, M., Kearns, M., Morgenstern, J., and Roth, A. Fairness in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1617–1626. JMLR. org, 2017.
- Joseph, M., Kearns, M., Morgenstern, J. H., and Roth, A. Fairness in learning: Classic and contextual bandits. In *Advances in Neural Information Processing Systems*, pp. 325–333, 2016.
- Kearns, M. and Singh, S. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232, 2002.
- Kilbertus, N., Carulla, M. R., Parascandolo, G., Hardt, M., Janzing, D., and Schölkopf, B. Avoiding discrimination through causal reasoning. In *Advances in Neural Information Processing Systems*, pp. 656–666, 2017.
- Kusner, M. J., Loftus, J., Russell, C., and Silva, R. Counterfactual fairness. In *Advances in Neural Information Processing Systems*, pp. 4066–4076, 2017.
- Lattimore, T. and Szepesvári, C. Bandit algorithms.
- Liu, L. T., Dean, S., Rolf, E., Simchowitz, M., and Hardt, M. Delayed impact of fair machine learning. In *ICML*, 2018.
- Mannor, S., Rubinstein, R. Y., and Gat, Y. The cross entropy method for fast policy search. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 512–519, 2003.
- Nabi, R. and Shpitser, I. Fair inference on outcomes. In *Proceedings of the... AAAI Conference on Artificial Intelligence. AAAI Conference on Artificial Intelligence*, volume 2018, pp. 1931. NIH Public Access, 2018.
- Pearl, J. Causality. Cambridge university press, 2009.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, 2019.
- Wen, M. and Topcu, U. Constrained cross-entropy method for safe reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 7461–7471, 2018.

# A. Additional Fairness Properties

**Equivalent of opportunity.** The following fairness property is based on the equality of opportunity fairness property from (Hardt et al., 2016) for supervised learning.

**Definition A.1.** Let M be an MDP with state space of the form  $S = Z \times Y \times \tilde{S}$ , where  $Z = \{\text{maj, min}\}$  and  $Y = \{\text{qual, unqual}\}$ , and let  $\rho \in \mathbb{R}^{|S| \times |A|}$  be the agent rewards. Then, a policy  $\pi$  satisfies *equality of opportunity* if

$$\mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{main}}^{(\pi)}}[\rho_{s,a}] = \mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{min}}^{(\pi)}}[\rho_{s,a}],$$

where 
$$\Lambda_z^{(\pi)} = \Lambda^{(\pi)} \mid \exists \tilde{s} \in \tilde{S} . s_0 = (z, \text{qual}, \tilde{s}).$$

This property is similar to demographic parity, but restricted to the *qualified* subpopulation—i.e., y = qual. For  $M_{\text{loan}}$ , opportunity says that loans should be given to qualified majority and minority members at equal rates (we assume an applicant is qualified if their true probability of repaying satisfies  $p \geq p_0$  for some  $p_0 \in [0,1]$ ). In particular, the policy can act arbitrarily for unqualified members y = unqual.

Path-specific causal fairness. We describe a causal notion of fairness from (Nabi & Shpitser, 2018), which we call path-specific causal fairness. We begin by giving background on causal graphs (Pearl, 2009). Our formulation differs from the one in (Nabi & Shpitser, 2018), but we believe it make the notion of mediated intervention more clear by isolating the source of randomness from the structural equations.

**Definition A.2.** A causal domain is a set of indices V = [k], where each index  $i \in V$  is associated with a random variable  $\epsilon_i \sim \mathcal{P}_i$ , which we call a noise term, with domain  $\mathcal{Z}_i$ . We use  $\epsilon_V = \begin{bmatrix} \epsilon_1 & \dots & \epsilon_n \end{bmatrix}^T \sim \mathcal{P}_V$  to denote the vector of all noise terms, and  $\mathcal{Z}_V$  to denote the domain of  $\epsilon_V$ .

Intuitively, a causal domain V indexes the set of variables  $X_1, ..., X_k$  of interest. In particular, these variables are functions of the noise terms indexed by V.

**Definition A.3.** Given a causal domain V, a *causal specification* is a set of equations  $H = \{h_1, ..., h_k\}$ , where  $h_i : \mathcal{Z}_V \to \mathcal{X}_i$  for each  $i \in V$ .

Intuitively, a causal specification describes how to construct the variables  $X_1,...,X_k$  given noise terms  $\epsilon_V$ , i.e.,  $X_i=h_i(\epsilon_V)$  for each  $i\in V$ . Typically, causal specifications are constructed from structural equations associated with a causal graph:

**Definition A.4.** Given a causal domain V, a causal graph is a directed acyclic graph G = (V, E, F), with vertices V, edges  $E \subseteq V \times V$ , and structural equations

 $F = \{f_1, ..., f_k\}$ , where

$$f_i: \left(\prod_{(j,i)\in E} \mathcal{X}_j\right) imes \mathcal{Z}_i o \mathcal{X}_i$$

for each  $i \in V$ . The corresponding causal specification is  $H^G$ , where  $h_i^G(\epsilon_V) = X_i$  is the solution to

$$X_i = f_i(\operatorname{pa}(X_i), \epsilon_i),$$

where  $pa(X_i) = \{X_j\}_{(j,i) \in E}$  are the parents of i; these equations have a unique solution since G is acyclic.

Intuitively, the vertices of a causal graph represent variables of interest, the edges represent dependence relationships between these variables, and the structural equations  $f_i$  capture the causal dependences of  $X_i$  on  $\operatorname{pa}(X_i)$ .

Given a causal graph G, we can modify the structural equations to specify changes in the causal structure of G. A key transformation is the following:

**Definition A.5.** Given causal graph G=(V,E,F), index  $i \in V$ , and value  $x \in \mathcal{X}_i$ , the *intervention specification* is the causal specification  $H^{G,\operatorname{do}(X_i=x)}$ , where  $h_i^{G,\operatorname{do}(X_i=x)}(\epsilon_V)=X_i$  is the solution to

$$X_j = \begin{cases} x & \text{if } j = i \\ f_j(\text{pa}(X_j), \epsilon_j) & \text{otherwise.} \end{cases}$$

We call *i* the *intervened variable*.

In other words,  $H^{G,do(X_i=x)}$  is computed replacing the previous value of  $X_i$  with a constant.

**Definition A.6.** Given causal graph G=(V,E,F), indices  $i,i'\in V$ , and values  $x,x'\in \mathcal{X}_i$ , the *mediated intervention specification* is the causal specification  $H^{G,\operatorname{do}(X_i=x),\operatorname{med}(X_{i'};X_i=x')}$ , where

$$h_j^{G,\operatorname{do}(X_i=x),\operatorname{med}(X_{i'};X_i=x')}(\epsilon_V) = X_j$$

is the solution to

$$X_{j} = \begin{cases} x & \text{if } j = i \\ h_{i'}^{G, \text{do}(X_{i} = x')}(\epsilon_{V}) & \text{if } j = i' \\ f_{j}(\text{pa}(X_{j}), \epsilon_{j}) & \text{otherwise.} \end{cases}$$

We call i' the *mediator variable*.

This causal specification is essentially the same as  $H^{G,\operatorname{do}(X_i=x)}$ , except  $X_{i'}$  is modified to equal its value according to the intervention specification  $H^{G,\operatorname{do}(X_i=x')}$ .

Now, we have the following fairness specification for supervised learning (Nabi & Shpitser, 2018):

**Definition A.7.** Given causal graph G, indices  $i, i' \in V$ , and values  $x, x' \in \mathcal{X}_i$ , a function  $\phi : \prod_{i=1}^n \mathcal{X}_i \to \mathcal{Y}$  satisfies *path-specific causal fairness* if

$$\mathbb{E}_{p(\epsilon_{V})}[\phi(\vec{h}^{G,\operatorname{do}(X_{i}=x),\operatorname{med}(X_{i'};X_{i}=x')}(\epsilon_{V}))]$$

$$= \mathbb{E}_{p(\epsilon_{V})}[\phi(\vec{h}^{G,\operatorname{do}(X_{i}=x')}(\epsilon_{V}))].$$
(15)

We call  $\mathcal{X}_i$  the *sensitive attribute* (where x represents the majority subpopulation and x' represents the minority subpopulation),  $\mathcal{X}_{i'}$  the *mediated attribute*, and  $\mathcal{Y}$  the *outcome*.

This specification captures the idea that the outcome should not change if we intervene on the sensitive attribute, except we ignore the effect of this intervention on the mediated attribute. For example, when a bank is deciding whether to give a loan, it should not discriminate against an individual based on their race (the sensitive attribute). However, they are allowed to base their decisions on attributes such as income (the mediated attribute), even if income is affected by race. Even if their income is lower because of past discrimination (e.g., they were unfairly denied a job in the past due to their race), this specification says it is not the responsibility of the bank to adjust for this discrepancy.

We have the following extension to the dynamical setting:

**Definition A.8.** Let M be an with states  $S = Z \times Y \times \tilde{S}$ , and a causal graph G with vertices  $V = \{Z, Y, \tilde{S}\}$ . Then, a policy  $\pi : S \to A$  satisfies path-specific causal fairness if

$$\mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{mai}}^{(\pi)}}[\rho_{s,a}] = \mathbb{E}_{\Lambda_{\mathrm{min}}^{(\pi)}}[\rho_{s,a}],$$

where  $\Lambda_{\rm mai}^{(\pi)}$  is  $\Lambda^{(\pi)}$  conditioned on starting from initial state

$$s_0 = \vec{h}^{G,\operatorname{do}(Z=\operatorname{maj}),\operatorname{med}(Y;Z=\operatorname{min})}(\epsilon_V)$$
  
 $\epsilon_V \sim p(\epsilon_V),$ 

and  $\Lambda_{\min}^{(\pi)}$  is  $\Lambda^{(\pi)}$  conditioned on starting from initial state

$$s_0 = \vec{h}^{G, \text{do}(Z=\text{min})}(\epsilon_V)$$
  
 $\epsilon_V \sim p(\epsilon_V).$ 

Unlike unlike demographic parity and equailty of opportunity,  $\Lambda_{\text{maj}}^{(\pi)}$  and  $\Lambda_{\text{min}}^{(\pi)}$  are asymmetric. We can impose symmetry between the two subpopulations by imposing a second constraint with maj and min swapped.

Finally, we can straightforwardly adapt both our modelbased and model-free reinforcement learning algorithms to work with path-specific causal fairness, since it only affects the initial state distribution. As with demographic parity and equal opportunity, the model-based algorithm requires separability to hold.

### B. Proof of Theorem 2.5

For the first claim, consider the MDP M. The states are  $s_0, s_1, s_2, s_3, s_4 \in \tilde{S} \times Z$ , where:

$$s_0 = (0, maj)$$
  
 $s_1 = (1, maj)$   
 $s_2 = (0, min)$   
 $s_3 = (1, min)$   
 $s_4 = (2, min)$ .

The actions are  $A = \{0, 1\}$ . The transitions are

$$\begin{aligned} P_{s_0,a,s_1} &= 1 \\ P_{s_1,a,s_1} &= 1 \\ P_{s_2,a,s_3} &= \mathbb{I}[a=0] \\ P_{s_2,a,s_4} &= \mathbb{I}[a=1] \\ P_{s_3,s_3} &= 1 \\ P_{s_4,s_4} &= 1 \end{aligned}$$

for all  $a \in A$ . The initial distribution is

$$D_{s_0} = D_{s_2} = \frac{1}{2}$$

$$D_{s_1} = D_{s_3} = D_{s_4} = 0.$$

The discount factor is  $\gamma = \frac{1}{2}$ . The agent rewards are

$$\rho_{s_0,a} = 0 
\rho_{s_1,a} = 1 
\rho_{s_2,a} = 0 
\rho_{s_3,a} = 0 
\rho_{s_4,a} = 2,$$

for all  $a \in A$ . Let  $\pi: S \to A$  be a deterministic policy. It is clear that the only value of  $\pi$  that matters is  $\pi(s_2)$ . Conditioned on z= maj, regardless of  $\pi$ , the expected cumulative agent reward is

$$\mathbb{E}_{(s,a) \sim \Lambda_{\text{maj}}^{(\pi)}}[\rho_{s,a}] = \left(1 - \frac{1}{2}\right) \sum_{t=1}^{\infty} \frac{1}{2^t}$$
$$= \frac{1}{2}.$$

Conditioned on  $z = \min$ , if  $\pi(s_2) = 0$ , then

$$\mathbb{E}_{(s,a) \sim \Lambda_{\min}^{(\pi)}}[\rho_{s,a}] = \begin{cases} 0 & \text{if } \pi(s_2) = 0\\ 1 & \text{if } \pi(s_2) = 1. \end{cases}$$

Therefore, it is impossible for the demographic parity constraint to be satisfied.

However, consider the stochastic policy

$$\pi_{s_2,0} = \pi_{s_2,1} = \frac{1}{2}.$$

Then,

$$\mathbb{E}_{(s,a)\sim\Lambda_{\min}^{(\pi)}}[\rho_{s,a}] = \frac{1}{2},$$

so this policy satisfies the demographic parity constraint.

For the second claim, consider the same MDP, except where

$$\rho_{s_4,a} = 0$$

for all  $a \in A$ . Then, it is clear that

$$\mathbb{E}_{(s,a)\sim\Lambda_{\min}^{(\pi)}}[\rho_{s,a}]=0$$

regardless of  $\pi$ . Thus, the demographic parity constraint cannot be satisfied—i.e.,  $\Pi_{DP} = 0$ .  $\square$ 

### C. Proof of Theorem 3.2

Our proof proceeds in three steps. First, we show that any feasible point of the LP in Algorithm 1 is the state-action distribution  $\Lambda^{(\pi)}$  for some policy  $\pi \in \Pi_{DP}$ . Second, we show that conversely, for any fair policy  $\pi \in \Pi_{DP}$ , the state-action distribution  $\Lambda^{(\pi)}$  is a feasible point of the LP. Finally, we combine these two results to prove the theorem.

**Step 1.** Let  $\pi \in \Pi_{DP}$  be any policy satisfying demographic parity. Then, we claim that the state-action distribution  $\Lambda^{(\pi)}$  is a feasible point of the LP in Algorithm 1.

First, we show that  $\Lambda^{(\pi)}$  satisfies the first constraint

$$\sum_{a \in A} \Lambda_{s',a}^{(\pi)} = (1 - \gamma) D_{s'} + \gamma \sum_{s \in S} \sum_{a \in A} \Lambda_{s,a}^{(\pi)} P_{s,a,s'}$$

for each  $s' \in S$ .

To this end, note that by induction,

$$D^{(\pi,t)} = (P^{(\pi)})^t D.$$

so

$$D^{(\pi)} = (1 - \gamma) \left[ \sum_{t=0}^{\infty} (\gamma P^{(\pi)})^t \right] D.$$
 (16)

Multiplying each side of (16) by  $I - \gamma P^{(\pi)}$  (where I is the  $|S| \times |S|$  identity matrix), we have

$$(I - \gamma P^{(\pi)})D^{(\pi)}$$

$$= (1 - \gamma) \left[ \sum_{t=0}^{\infty} (\gamma P^{(\pi)})^t - \sum_{t=1}^{\infty} (\gamma P^{(\pi)})^t \right] D$$

$$= (1 - \gamma) \cdot D.$$

Note that these algebraic manipulations are valid since the eigenvalues of  $\gamma P^{(\pi)}$  are bounded in norm by  $\gamma < 1$ , so all sums converge absolutely. Rearranging this equality gives

$$D^{(\pi)} = (1 - \gamma)D + \gamma P^{(\pi)}D^{(\pi)}.$$
 (17)

It follows that

$$\sum_{a \in A} \Lambda_{s',a}^{(\pi)} = (1 - \gamma) D_{s'} + \gamma \sum_{s \in S} \sum_{a \in A} \Lambda_{s,a}^{(\pi)} P_{s,a,s'}$$

for each  $s' \in S$ , where we have used the equalities

$$D_{s'}^{(\pi)} = \sum_{a \in A} \Lambda_{s',a}^{(\pi)}$$

and

$$\begin{split} (P^{(\pi)}D^{(\pi)})_{s'} &= \sum_{s \in S} P_{s,s'}^{(\pi)}D_s^{(\pi)} \\ &= \sum_{s \in S} \sum_{a \in A} P_{s,a,s'}\pi_{s,a}D_s^{(\pi)} \\ &= \sum_{s \in S} \sum_{a \in A} P_{s,a,s'}\Lambda_{s,a}^{(\pi)} \end{split}$$

that follow from the definition of  $\Lambda^{(\pi)}$ . Therefore,  $\Lambda^{(\pi)}$  satisfies the first constraint.

Next, we show that  $\Lambda^{(\pi)}$  satisfies the second constraint, which says that there exists  $c \in \mathbb{R}_+$  such that

$$p_z^{-1} \sum_{\tilde{s} \in \tilde{S}} \sum_{a \in A} \Lambda_{(z,\tilde{s}),a}^{(\pi)} \rho_{(z,\tilde{s}),a} = c$$

for all  $z \in \mathcal{Z}$ .

In particular, note that

$$D_z^{(\pi)} = D^{(\pi)} \mid \exists \tilde{s} \in \tilde{S} . s = (z\tilde{s}),$$

since the value of z for s equals the value of z for the initial state  $s_0 \sim D$ . Furthermore, the probability of sampling  $s \sim D^{(\pi)} \mid \exists \tilde{s} \in \tilde{S} \ . \ s = (z, \tilde{s})$  is

$$\frac{D_s^{(\pi)}\mathbb{I}[\exists \tilde{s} \in \tilde{S} . s = (z, \tilde{s})]}{p_z}.$$

Together with the definition of  $\Lambda_z^{(\pi)}$ , we have

$$\begin{split} (\Lambda_z^{(\pi)})_{s,a} &= (D_z^{(\pi)})_s \pi_{s,a} \\ &= \frac{D_s^{(\pi)} \mathbb{I}[\exists \tilde{s} \in \tilde{S} . s = (z, \tilde{s})]}{p_z} \cdot \pi_{s,a} \\ &= \frac{\Lambda_{s,a}^{(\pi)} \mathbb{I}[\exists \tilde{s} \in \tilde{S} . s = (z, \tilde{s})]}{p_z}. \end{split}$$

Therefore, we have

$$\mathbb{E}_{(s,a)\sim\Lambda_{z}^{(\pi)}}[\rho_{s,a}]$$

$$= \sum_{s\in S} \sum_{a\in A} \frac{\Lambda_{s,a}^{(\pi)}\mathbb{I}[\exists \tilde{s}\in \tilde{S} \cdot s = (z,\tilde{s})]}{p_{z}} \cdot \rho_{s,a}$$

$$= p_{z}^{-1} \sum_{\tilde{s}\in \tilde{S}} \sum_{a\in A} \Lambda_{s,a}^{(\pi)} \rho_{s,a}.$$
(18)

By assumption,  $\pi$  satisfies the demographic parity constraint, which says that (18) is constant for all  $z \in \mathcal{Z}$ . Equivalently, there exists  $c \in \mathbb{R}_+$  such that (18) equals c for all  $z \in \mathcal{Z}$ . Thus,  $\Lambda^{(\pi)}$  satisfies the second constraint.

Therefore,  $\Lambda^{(\pi)}$  is a feasible point of the LP, as claimed.

**Step 2.** Let  $\lambda \in \mathbb{R}^{|S| \times |A|}$  be a feasible point of the LP in Algorithm 1, and let

$$\pi_{s,a} = \frac{\lambda_{s,a}}{\sum_{a' \in A} \lambda_{s,a'}}$$

be the corresponding policy returned by Algorithm 1. Then, we claim that  $\lambda = \Lambda^{(\pi)}$ , that  $\pi \in \Pi_{DP}$ , and that the value of the objective for  $\lambda$  equals  $R^{(\pi)}$ .

To see the first claim, let  $d \in \mathbb{R}^{|S|}$  be defined by

$$d_s = \sum_{a \in A} \lambda_{s,a}.$$

We show that  $D^{(\pi)} = d$ . To this end, note that because  $\lambda$  satisfies the first constraint in the LP, we have

$$\sum_{a \in A} \lambda_{s',a} = (1 - \gamma)D_{s'} + \gamma \sum_{s \in S} \sum_{a \in A} \lambda_{s,a} P_{s,a,s'}.$$

Together with the equality

$$\pi_{s,a} = \frac{\lambda_{s,a}}{d_s},$$

we have

$$d_{s'} = (1 - \gamma)D_{s'} + \gamma \sum_{s \in S} \sum_{a \in A} d_s \pi_{s,a} P_{s,a,s'}$$
$$= (1 - \gamma)D_{s'} + \gamma (P^{(\pi)}d)_{s'}.$$

Thus,

$$d = (1 - \gamma)D + \gamma P^{(\pi)}d. \tag{19}$$

We note that  $I - \gamma P^{(\pi)}$  is invertible—in particular, the eigenvalues of  $\gamma P^{(\pi)}$  have norms bounded by  $\gamma$ , so the eigenvalues of  $I - \gamma P^{(\pi)}$  have norms bounded below by  $1 - \gamma$ ; therefore, the eigenvalues of  $I - \gamma P^{(\pi)}$  are nonzero, so it is invertible. As a consequence, we can solve for d in (19) to get

$$d = (1 - \gamma)(I - \gamma P^{(\pi)})^{-1}D.$$

Finally, from (17) in Step 1 of this proof, we established that  $D^{(\pi)}$  similarly satisfies

$$D^{(\pi)} = (1 - \gamma)D + \gamma P^{(\pi)}D^{(\pi)}.$$

As before, since  $I - \gamma P^{(\pi)}$  is invertible, we have

$$D^{(\pi)} = (1 - \gamma)(I - \gamma P^{(\pi)})^{-1}D = d.$$

Thus,

$$\lambda_{s,a} = d_s \pi_{s,a} = D_s \pi_{s,a} = \Lambda_{s,a}^{(\pi)},$$

so the first claim follows.

To see the second claim, note that since  $\lambda$  is feasible, it must satisfy the second constraint of the LP, which says that there exists  $c \in \mathbb{R}_+$  such that

$$p_z^{-1} \sum_{\tilde{s} \in \tilde{S}} \sum_{a \in A} \lambda_{(z,\tilde{s}),a} \rho_{(z,\tilde{s}),a} = c$$

for all  $z \in \mathcal{Z}$ . Since  $\lambda = \Lambda^{(\pi)}$ , the same holds true for  $\Lambda^{(\pi)}$ , i.e., there exists  $c \in \mathbb{R}_+$  such that

$$p_z^{-1} \sum_{\tilde{s} \in \tilde{S}} \sum_{a \in A} \lambda_{(z,\tilde{s}),a} \rho_{(z,\tilde{s}),a} = c$$
 (20)

for all  $z \in \mathcal{Z}$ . As shown in the first step of this proof, (20) is equivalent to the demographic parity constraint. Thus,  $\pi \in \Pi_{DP}$ , as claimed.

To see the third claim, note that

$$R^{(\pi)} = (1 - \gamma) \mathbb{E}_{(s,a) \sim \Lambda^{(\pi)}} [R_{s,a}]$$
$$(1 - \gamma) \sum_{s \in S} \sum_{a \in A} \Lambda_{s,a}^{(\pi)} R_{s,a}.$$

In other words, the value of the objective of the LP for the point  $\lambda$  is equal to  $R^{(\pi)}$ , as claimed.

**Step 3.** Finally, we use the results from the previous two steps to prove the theorem statement. First, let  $\pi^*$  be the solution to (2). Then, by the claim shown in the first step,  $\Lambda^{(\pi^*)}$  is a feasible point of the LP in Algorithm 1. Furthermore, by the claim shown in the second step, the value of the objective for  $\lambda = \Lambda^{(\pi^*)}$  is  $R^{(\pi^*)}$ .

Next, let  $\lambda^0$  be the solution to the LP in Algorithm 1. By the claim shown in the second step, (i)  $\lambda_0 = \Lambda^{(\pi_0)}$ , where  $\pi_0$  is the policy returned by Algorithm 1, (ii)  $\pi_0 \in \Pi_{\mathrm{DP}}$ , and (iii) the value of the objective for  $\lambda^0$  is  $R^{(\pi_0)}$ .

It follows that  $R^{(\pi^*)} \leq R^{(\pi_0)}$ , since  $\pi_0$  maximizes the objective of the LP over feasible points (and  $\Lambda^{(\pi^*)}$  is feasible). Since  $\pi_0 \in \Pi_{\mathrm{DP}}$ , it follows that  $\pi_0$  is also a solution to (2). Thus, we have proven the theorem statement.  $\square$ 

### D. Proof of Theorem1 4.2 & 4.3

Our proof proceeds in four steps. First, we bound the error  $|\tilde{R}^{(\pi)}-R^{(\pi)}|$  due to truncation. Second, we use Step 1 to prove Theorem 4.2. Third, we bound the estimation error  $|\hat{R}^{(\pi)}-\tilde{R}^{(\pi)}|$ . Fourth, we combine steps 1 and 3 to prove Theorem 4.3.

**Step 1.** Note that for any policy  $\pi$ , we have

$$\begin{split} |\tilde{R}^{(\pi)} - R^{(\pi)}| &= \left| \sum_{t=T}^{\infty} \gamma^t \langle R, \Lambda^{(\pi,t)} \rangle \right| \leq \sum_{t=T}^{\infty} \gamma^t R_{\max} \\ &\leq \frac{\gamma^T R_{\max}}{1 - \gamma} \\ &\leq \frac{\sigma^2 \epsilon}{4}. \end{split}$$

Similarly, we have

$$|\tilde{\rho}_z^{(\pi)} - \tilde{\rho}_z^{(\pi)}| \le \frac{\sigma^2 \epsilon}{4}$$

for all  $z \in Z$ .

**Step 2.** Now, we can prove Theorem 4.3. First, note that

$$\begin{split} &|\rho_{\mathrm{maj}}^{(\tilde{\pi}^*)} - \rho_{\mathrm{min}}^{(\tilde{\pi}^*)}| \\ &= |\rho_{\mathrm{maj}}^{(\tilde{\pi}^*)} - \tilde{\rho}_{\mathrm{maj}}^{(\tilde{\pi}^*)}| + |\tilde{\rho}_{\mathrm{maj}}^{(\tilde{\pi}^*)} - \tilde{\rho}_{\mathrm{min}}^{(\tilde{\pi}^*)}| + |\tilde{\rho}_{\mathrm{min}}^{(\tilde{\pi}^*)} - \rho_{\mathrm{min}}^{(\tilde{\pi}^*)}| \\ &\leq \frac{\sigma^2 \epsilon}{4} + (1 - \sigma)\epsilon + \frac{\sigma^2 \epsilon}{4} \\ &\leq \epsilon, \end{split}$$

so  $\tilde{\pi}^* \in \Pi_{DP,\epsilon}$ , so the first claim. Similarly,

$$\begin{split} &|\rho_{\text{maj}}^{(\pi^*)} - \rho_{\text{min}}^{(\pi^*)}| \\ &= |\rho_{\text{maj}}^{(\pi^*)} - \tilde{\rho}_{\text{maj}}^{(\pi^*)}| + |\tilde{\rho}_{\text{maj}}^{(\pi^*)} - \tilde{\rho}_{\text{min}}^{(\pi^*)}| + |\tilde{\rho}_{\text{min}}^{(\pi^*)} - \rho_{\text{min}}^{(\pi^*)}| \\ &\leq \frac{\sigma^2 \epsilon}{4} + (1 - \sigma)^2 \epsilon + \frac{\sigma^2 \epsilon}{4} \\ &< \epsilon. \end{split}$$

since  $\sigma \leq \frac{1}{2} \leq 1 - \sigma$ . So,  $\pi^* \in \tilde{\Pi}_{DP,(1-\sigma)\epsilon}$ . Thus, we have

$$\begin{split} &R^{(\pi^*)} - R^{(\tilde{\pi}^*)} \\ &\leq (R^{(\pi^*)} - \tilde{R}^{(\pi^*)}) + (\tilde{R}^{(\pi^*)} - \tilde{R}^{(\tilde{\pi}^*)}) + (\tilde{R}^{(\tilde{\pi}^*)} - R^{(\tilde{\pi}^*)}) \\ &\leq \frac{\sigma^2 \epsilon}{4} + 0 + \frac{\sigma^2 \epsilon}{4} \\ &= \frac{\sigma^2 \epsilon}{\epsilon}, \end{split}$$

where the second inequality follows because  $\tilde{\pi}^*$  maximizes  $\tilde{R}^{(\pi)}$  among policies  $\pi \in \tilde{\Pi}_{\mathrm{DP},(1-\sigma)\epsilon}$ . The second claim in the theorem follows. Thus, the theorem statement follows.

**Step 3.** Note that  $\hat{R}^{(\pi)}$  be an estimate of  $\tilde{R}^{(\pi)}$  using m sampled rollouts  $\zeta^{(1)},...,\zeta^{(m)}$ . First, note that

$$|\hat{R}^{(\pi)}| \le \frac{R_{\max}}{1 - \gamma}$$

is bounded, so we can apply Hoeffding's inequality (see Lemma H.1) to get

$$\begin{split} \Pr \left[ |\hat{R}^{(\pi)} - \tilde{R}^{(\pi)}| \geq \frac{\sigma \epsilon}{4} \right] \leq 2 \exp \left( -\frac{m \sigma^2 \epsilon^2}{32 R_{\max}/(1-\gamma)} \right) \\ \leq \frac{\delta}{3} \end{split}$$

Similarly, for all  $z \in Z$ , we have

$$\Pr\left[|\hat{\rho}_z^{(\pi)} - \tilde{\rho}_z^{(\pi)}| \ge \frac{\sigma\epsilon}{4}\right] \le \frac{\delta}{3}.$$

Since  $Z = \{\text{maj}, \text{min}\}$ , by a union bound,

$$\begin{split} |\hat{R}^{(\pi)} - \tilde{R}^{(\pi)}| &\leq \frac{\epsilon}{4} \\ |\hat{\rho}_z^{(\pi)} - \tilde{\rho}_z^{(\pi)}| &\leq \frac{\epsilon}{4} \quad \ (\forall z \in Z) \end{split}$$

with probability at least  $1 - \delta$ .

**Step 4.** Now, we can prove Theorem 4.2. First, note that with probability  $1 - \delta$ ,

$$\begin{split} |\hat{R}^{(\pi)} - R^{(\pi)}| &\leq |\hat{R}^{(\pi)} - \tilde{R}^{(\pi)}| + |\tilde{R}^{(\pi)} - R^{(\pi)}| \\ &\leq \frac{\sigma\epsilon}{4} + \frac{\sigma^2\epsilon}{4} \\ &\leq \frac{\sigma\epsilon}{2}, \end{split}$$

as well as

$$|\hat{\rho}_z^{(\pi)} - \rho_z^{(\pi)}| \le \frac{\sigma\epsilon}{2}$$

for all  $z \in Z$ . The first claim follows. Next, note that

$$\begin{split} &|\rho_{\text{maj}} - \rho_{\text{min}}| \\ &\leq |\rho_{\text{maj}} - \hat{\rho}_{\text{maj}}| + |\hat{\rho}_{\text{maj}} - \hat{\rho}_{\text{min}}| + |\hat{\rho}_{\text{min}} - \rho_{\text{min}}| \\ &\leq \frac{\sigma\epsilon}{2} + (1 - \sigma)\epsilon + \frac{\sigma\epsilon}{2} \\ &= \epsilon, \end{split}$$

which implies that  $\pi \in \Pi_{\mathrm{DP},\epsilon}$ , so the second claim follows. Thus, the theorem follows.  $\square$ 

#### E. Proof of Theorem 5.2

Our proof proceeds in two steps. First, we prove that the constraint in (14) is equivalent to (13). Second, we prove that if  $\pi$  satisfies (13), then  $\pi \in \Pi_{DP}$ . The theorem statement follows from these two claims, since they show that any solution to (14) must satisfy demographic parity.

**Step 1.** We claim that (13) and the constraint in (14) are equivalent.

First, we show that the constraint in (14) implies (13). To this end, let  $D' \in \Delta^{|S|}$ . By the constraint in (14), there exists  $c \in \mathbb{R}_+$  such that for any  $z \in \mathcal{Z}$ , we have

$$\mathbb{E}_{s \sim D_z'} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right] = \mathbb{E}_{s \sim D_z'}[c] = c.$$

Thus, (14) holds.

Second, we show that (13) implies (14). First, we show that for any two states  $s_{\text{maj}} = (\text{maj}, \tilde{s})$  and  $s_{\text{min}} = (\text{min}, \tilde{s}')$ , for any  $\tilde{s}, \tilde{s}' \in \tilde{S}$ , (13) implies that

$$\sum_{a \in A} \pi_{s_{\text{maj}},a} \rho_{s_{\text{maj}},a} = \sum_{a \in A} \pi_{s_{\text{min}},a} \rho_{s_{\text{min}},a}. \tag{21}$$

To this end, let  $D' \in \Delta^{|S|}$  be defined by

$$D_s' = \frac{1}{2} \delta_{s,s_{\text{maj}}} + \frac{1}{2} \delta_{s,s_{\text{min}}}.$$

Then, note that

$$(D'_z)_s = \frac{(D'_s)\mathbb{I}[\exists \tilde{s} \in \tilde{S} \cdot s = (z, \tilde{s})]}{\sum_{s' \in S} (D'_{s'}\mathbb{I}[\exists \tilde{s} \in \tilde{S} \cdot s' = (z, \tilde{s})])}$$
$$= \frac{(1/2)\delta_{s,s_z}}{1/2}$$
$$= \delta_{s,s_z}.$$

Thus,

$$\mathbb{E}_{s \sim D_z'} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right] = \sum_{a \in A} \pi_{s_z,a} \rho_{s_z,a},$$

so by (13), we have

$$\sum_{a \in A} \pi_{s_{\mathrm{maj}},a} \rho_{s_{\mathrm{maj}},a} = \sum_{a \in A} \pi_{s_{\mathrm{min}},a} \rho_{s_{\mathrm{min}},a},$$

so we have shown that (21) holds. Next, we show that for any two states  $s_0, s_1 \in S$ , (13) implies that

$$\sum_{a \in A} \pi_{s_0, a} \rho_{s_0, a} = \sum_{a \in A} \pi_{s_1, a} \rho_{s_1, a}. \tag{22}$$

Note that (21) implies (22) when the sensitive attribute of  $s_0$  and  $s_1$  are different. Thus, it remains to show that (22) holds when  $s_0$  and  $s_1$  have the same sensitive attribute. In this case, let  $s_2$  be any state with a different sensitive attribute than  $s_0$  and  $s_1$ . Then, by (21), we have

$$\sum_{a \in A} \pi_{s_0,a} \rho_{s_0,a} = \sum_{a \in A} \pi_{s_2,a} \rho_{s_2,a} = \sum_{a \in A} \pi_{s_1,a} \rho_{s_1,a},$$

so (22) holds. Finally, note that (21) implies that

$$\sum_{a \in A} \pi_{s,a} \rho_{s,a}$$

is constant for all  $s \in S$ . In other words, there exists  $c \in \mathbb{R}_+$  such that

$$\sum_{a \in A} \pi_{s,a} \rho_{s,a} = c$$

for all  $s \in S$ . This statement is equivalent to the constraint in (14), so (13) implies the constraint in (14).

Thus, (13) and the constraint in (14) are equivalent, so the claim follows.

**Step 2.** We claim that if (13) holds for a policy  $\pi$ , then  $\pi \in \Pi_{DP}$ . To this end, recall that the demographic parity constraint is

$$\mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{mai}}^{(\pi)}}[\rho_{s,a}] = \mathbb{E}_{(s,a)\sim\Lambda_{\mathrm{min}}^{(\pi)}}[\rho_{s,a}].$$

Note that for any  $z \in \mathcal{Z}$ ,

$$\mathbb{E}_{(s,a)\sim\Lambda_z^{(\pi)}}[\rho_{s,a}] = \sum_{s\in S} \sum_{a\in A} (\Lambda_z^{(\pi)})_{s,a} \rho_{s,a}$$
$$= \sum_{s\in S} \sum_{a\in A} (D_z^{(\pi)})_s \pi_{s,a} \rho_{s,a}$$
$$= \mathbb{E}_{s\sim D_z^{(\pi)}} \left[ \sum_{a\in A} \pi_{s,a} \rho_{s,a} \right].$$

Thus, the demographic parity constraint is equivalent to

$$\mathbb{E}_{s \sim D_{\text{maj}}^{(\pi)}} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right] = \mathbb{E}_{s \sim D_{\text{min}}^{(\pi)}} \left[ \sum_{a \in A} \pi_{s,a} \rho_{s,a} \right],$$

which is exactly (13) for  $D' = D^{(\pi)}$ . The claim follows.

### F. Proof of Theorem 6.4

Our proof proceeds in two steps. First, we prove that for any policy  $\pi$ , given two initial distributions D and  $\tilde{D}$  satisfying

$$||D - \tilde{D}||_{\infty} \le \epsilon_0,$$

then

$$|R^{(\pi)} - \tilde{R}^{(\pi)}| \le \frac{|S| \cdot R_{\max} \cdot \epsilon_0}{1 - \gamma},$$

where  $R^{(\pi)}$  (resp.,  $\tilde{R}^{(\pi)}$ ) is the expected cumulative distribution assuming the initial distribution is D (resp.,  $\tilde{D}$ ), and similarly for the agent rewards  $\rho$ . Second, we use this fact to prove the theorem statement.

### **Step 1.** We claim that assuming

$$||D - \tilde{D}||_{\infty} \le \epsilon_0,$$

then for any policy  $\pi$ , we have

$$|R^{(\pi)} - \tilde{R}^{(\pi)}| \le \frac{|S| \cdot R_{\max} \cdot \epsilon_0}{1 - \gamma},$$

where  $R^{(\pi)}$  is the expected cumulative reward for  $\pi$  in the MDP  $M=(S,A,D,P,R,\gamma)$  and  $\tilde{R}^{(\pi)}$  is the expected cumulative reward for  $\pi$  in the MDP  $\tilde{M}=(S,A,\tilde{D},P,R,\gamma)$ . In addition, for all  $z\in Z$ , we have

$$|\rho_z^{(\pi)} - \tilde{\rho}_z^{(\pi)}| \le \frac{\cdot |S| \cdot R_{\max} \cdot \epsilon_0}{1 - \gamma},$$

where

$$\rho_z^{(\pi)} = \mathbb{E}_{(s,a) \sim \Lambda_z^{(\pi)}} [\rho_{s,a}]$$

is the expected cumulative agent reward for the MDP M, and  $\tilde{\rho}_z^{(\pi)}$  is the expected cumulative agent reward for the MDP  $\tilde{M}$ . We only prove the claim for  $|R^{(\pi)} - \tilde{R}^{(\pi)}|$ ; the claim for  $|\rho_z^{(\pi)} - \tilde{\rho}_z^{(\pi)}|$  follows using the same argument.

Let  $W \in \mathbb{R}^{|S|}$  be

$$W_s = \langle \pi_{s,\cdot}, R_{s,\cdot} \rangle = \sum_{a \in A} \pi_{s,a} R_{s,a}.$$

Then, we have

$$\begin{split} R^{(\pi)} &= (1 - \gamma) \langle R, \Lambda^{(\pi)} \rangle \\ &= (1 - \gamma) \sum_{s \in S} \sum_{a \in A} D_s^{(\pi)} \pi_{s,a} R_{s,a} \\ &= (1 - \gamma) \langle D^{(\pi)}, W \rangle. \end{split}$$

Now, note that

$$D^{(\pi,t)} = (P^{(\pi)})D,$$

so we have

$$\begin{split} D^{(\pi)} &= \sum_{t=0}^{\infty} \frac{\gamma^t}{1 - \gamma} D^{(\pi, t)} \\ &= \frac{1}{1 - \gamma} \left[ \sum_{t=0}^{\infty} (\gamma P^{(\pi)})^t D \right]. \end{split}$$

Thus,

$$R^{(\pi)} = \sum_{t=0}^{\infty} \gamma^t \left\langle (P^{(\pi)})^t D, W \right\rangle.$$

Similarly,

$$\tilde{R}^{(\pi)} = \sum_{t=0}^{\infty} \gamma^t \left\langle (P^{(\pi)})^t \tilde{D}, W \right\rangle.$$

It follows that

$$R^{(\pi)} - \tilde{R}^{(\pi)} = \sum_{t=0}^{\infty} \gamma^t \left\langle (P^{(\pi)})^t D - (P^{(\pi)})^t \tilde{D}, W \right\rangle.$$

Thus,

$$|R^{(\pi)} - \tilde{R}^{(\pi)}|$$

$$\leq \sum_{t=0}^{\infty} \gamma^{t} \| (P^{(\pi)})^{t} D - (P^{(\pi)})^{t} \tilde{D} \|_{\infty} \cdot \|W\|_{1}$$

$$\leq \sum_{t=0}^{\infty} \gamma^{t} \| P^{(\pi)} \|_{\infty}^{t} \cdot \|D - \tilde{D}\|_{\infty} \cdot \|W\|_{1}, \qquad (23)$$

where the first line follows from Hölder's inequality and the second line follows from properties of the matrix norm. Note that

$$||P^{(\pi)}||_{\infty} = \max_{s \in S} \sum_{s' \in S} |P_{s,s'}^{(\pi)}| = 1,$$
 (24)

since  $P^{(\pi)}$  is a stochastic matrix. Furthermore,

$$\begin{aligned} |W_s| &= |\langle \pi_{s,\cdot}, R_{s,\cdot} \rangle| \leq \|\pi_{s,\cdot}\|_1 \cdot \|R_{s,\cdot}\|_{\infty} \\ &\leq \|R_{s,\cdot}\|_{\infty} \\ &\leq R_{\max}, \end{aligned}$$

where the first inequality follows from Hölder's inequality and the second inequality follows since  $\pi_{s,.}$  is a discrete probability distribution. Therefore,

$$||W||_1 = \sum_{s \in S} |W_s| \le |S| \cdot R_{\text{max}}.$$
 (25)

Plugging (24) and (25) into (23) gives

$$|R^{(\pi)} - \tilde{R}^{(\pi)}| \le |S| \cdot R_{\max} \cdot ||D - \tilde{D}||_{\infty} \cdot \sum_{t=0}^{\infty} \gamma^{t}$$

$$\le \frac{|S| \cdot R_{\max} \cdot \epsilon_{0}}{1 - \gamma},$$

as claimed.

**Step 2.** Now, we prove the theorem. Let  $\tilde{\pi}$  be the optimal policy for  $\tilde{M}$  (i.e., initial distribution  $\tilde{D}=d^{(\pi_0)}$ ) satisfying  $\tilde{\pi}\in \tilde{\Pi}_{\mathrm{DP},\epsilon/2}$ . Similarly, let  $\tilde{\pi}^*$  be the optimal policy for M (i.e., initial distribution  $D=D^{(\pi_0,T_0)}$ ) satisfying  $\tilde{\pi}^*\in \Pi_{\mathrm{DP},\epsilon/4}$ . Here,  $T_0$  is the  $\epsilon_0$  mixing time of  $\pi_0$  for

$$\epsilon_0 = \frac{(1 - \gamma)\epsilon}{8|S| \cdot R_{\text{max}}}.$$

Then, by definition of  $T_0$ , we have

$$\|\tilde{D} - D\|_{\infty} \le \epsilon_0.$$

Thus, by the first step, for all  $z, z' \in Z$ , we have

$$\begin{split} & \rho_z^{(\tilde{\pi})} - \rho_{z'}^{(\tilde{\pi})} \\ & \leq (\rho_z^{(\tilde{\pi})} - \tilde{\rho}_z^{(\tilde{\pi})}) + (\tilde{\rho}_z^{(\tilde{\pi})} - \tilde{\rho}_{z'}^{(\tilde{\pi})}) + (\rho_{z'}^{(\tilde{\pi})} - \tilde{\rho}_{z'}^{(\tilde{\pi})}) \\ & \leq \frac{|S| \cdot R_{\max} \cdot \epsilon_0}{1 - \gamma} + \frac{\epsilon}{2} + \frac{|S| \cdot R_{\max} \cdot \epsilon_0}{1 - \gamma} \\ & \leq \epsilon, \end{split}$$

where the inequality on the third line follows because  $\tilde{\pi} \in \tilde{\Pi}_{DP,\epsilon/2}$ . Thus, we guarantee that  $\tilde{\pi} \in \Pi_{DP,\epsilon}$ .

Next, note that similarly, for all  $z, z' \in Z$ , we have

$$\tilde{\rho}_z^{(\tilde{\pi}^*)} - \tilde{\rho}_{z'}^{(\tilde{\pi}^*)} \leq \frac{\epsilon}{2},$$

so  $\tilde{\pi}^* \in \tilde{\Pi}_{\mathrm{DP},\epsilon/2}.$  As a consequence, we have

$$\begin{split} &R^{(\tilde{\pi}^*)} - R^{(\tilde{\pi})} \\ &= (R^{(\tilde{\pi}^*)} - \tilde{R}^{(\tilde{\pi}^*)}) + (\tilde{R}^{(\tilde{\pi}^*)} - \tilde{R}^{(\tilde{\pi})}) + (\tilde{R}^{(\tilde{\pi})} - R^{(\tilde{\pi})}) \\ &\leq \frac{|S| \cdot R_{\text{max}} \cdot \epsilon_0}{1 - \gamma} + 0 + \frac{|S| \cdot R_{\text{max}} \cdot \epsilon_0}{1 - \gamma} \\ &\leq \epsilon, \end{split}$$

where the inequality on the third line follows because  $\tilde{\pi}$  maximizes  $\tilde{R}^{(\pi)}$  over  $\pi \in \tilde{\Pi}_{DP,\epsilon/2}$ , and  $\tilde{\pi}^* \in \tilde{\Pi}_{DP,\epsilon/2}$ .

Thus, the theorem statement follows.  $\Box$ 

### G. Proof of Theorem 6.1

Our proof proceeds in three steps. First, we prove that for any  $\epsilon_0$ ,  $\delta_0$ , we can choose  $N_0$  sufficiently large so that

$$||P - \hat{P}||_{\infty} \le \epsilon_0$$

with probability at least  $1 - \delta_0$ . Second, we prove that assuming  $||P - \hat{P}||_{\infty} \le \epsilon_0$ , then for any policy  $\pi$ , we have

$$|R^{(\pi)} - \hat{R}^{(\pi)}| \le T \cdot |S| \cdot R_{\max} \cdot \epsilon_0,$$

where  $R^{(\pi)}$  (resp.,  $\hat{R}^{(\pi)}$ ) is the expected cumulative distribution assuming the transitions are P (resp.,  $\hat{P}$ ), and similarly for the agent rewards  $\rho$ . Third, we use the first two steps to prove the theorem statement.

**Step 1.** Given  $\epsilon_0, \delta_0 \in \mathbb{R}_+$ , we claim that for

$$N_0 = \frac{2\log(2|S|^2|A|/\delta_0)}{\lambda_0^2 \epsilon_0^2},$$

then our estimate  $\hat{P}$  satisfies

$$\|\hat{P} - P\|_{\infty} \le \epsilon_0$$

with probability at least  $1 - \delta_0$ .

Let  $I_{s,a}$  be the random variable indicating whether our algorithm observes a tuple (s,a,s') (for some  $s' \in S$ ) on a single episode, and let  $I_{s,a,i}$  be samples of  $I_{s,a}$  for each of the  $N_0$  exploratory episodes taken by our algorithm. Let

$$\begin{split} & \mu_{s,a}^{(I)} = \mathbb{E}[I_{s,a}] \\ & \hat{\mu}_{s,a}^{(I)} = \frac{1}{N_0} \sum_{i=1}^{N_0} I_{s,a,i}. \end{split}$$

Then, by Hoeffding's inequality (see Lemma H.1), we have

$$\Pr\left[|\hat{\mu}_{s,a}^{(I)} - \mu_{s,a}^{(I)}| \ge \epsilon\right] \le 2e^{-2N_0\epsilon^2}.$$
 (26)

By assumption, we have

$$\mu_{s,a}^{(I)} = \Lambda_{s,a}^{(\pi_0)} \ge \lambda_0,$$

so using  $\epsilon = \lambda_0/2$  in (26), we have

$$\hat{\mu}_{s,a}^{(I)} \ge \frac{\mu_{s,a}^{(I)}}{2} \ge \frac{\lambda_0}{2} \tag{27}$$

with probability at least

$$1 - 2e^{-N_0(\mu_{s,a}^{(I)})^2/2} > 1 - 2e^{-N_0\lambda_0^2/2}$$

Taking a union bound over  $s \in S$  and  $a \in A$ , we have (27) holds for every  $s \in S$  and  $a \in A$  with probability at least

$$1 - 2|S| \cdot |A| \cdot e^{-N_0 \lambda_0^2/2}. (28)$$

In this event, we have at least  $\frac{N_0\lambda_0}{2}$  observations (s,a,s') (for some  $s'\in S$ ) for every  $s\in S$  and  $a\in A$ .

Now, for an observation (s,a,s''), let  $J_{s,a,s'}$  be the random variable indication whether s'=s''. Without loss of generality, we assume that we have exactly  $N_1=\frac{N_0\lambda_0}{2}$  samples  $J_{s,a,s',j}$  of  $J_{s,a,s'}$  for each  $s\in S$  and  $a\in A$ . Let

$$\mu_{s,a,s'}^{(J)} = \mathbb{E}[J_{s,a,s'}]$$

$$\hat{\mu}_{s,a,s'}^{(J)} = \frac{1}{N_1} \sum_{i=1}^{N_1} J_{s,a,s',j}.$$

Then, by Hoeffding's inequality (see Lemma H.1), we have

$$\Pr\left[|\hat{\mu}_{s,a,s'}^{(J)} - \mu_{s,a,s'}| \ge \epsilon\right] \le 2e^{-2N_1\epsilon^2}.$$
 (29)

Note that by definition,  $\mu_{s,a,s'}^{(J)} = P_{s,a,s'}$  and  $\hat{\mu}_{s,a,s'}^{(J)} = \hat{P}_{s,a,s'}$ . Thus, taking  $\epsilon = \epsilon_0$  in (29), we have

$$|P_{s,a,s'} - \hat{P}_{s,a,s'}| \le \epsilon_0$$
 (30)

with probability at least

$$1 - 2e^{-2N_1\epsilon_0^2}.$$

Taking a union bound over all  $s, s' \in S$  and  $a \in A$ , we have (30) for all  $s, s' \in S$  and  $a \in A$  with probability at least

$$1 - 2|S|^2|A| \cdot e^{-2N_1\epsilon_0^2}. (31)$$

In other words, in this event, we have  $||P - \hat{P}||_{\infty} \le \epsilon_0$ .

Taking a union bound over (28) and (31), we have

$$||P - \hat{P}||_{\infty} \le \epsilon_0$$

with probability at least

$$\begin{split} &1 - 2|S|^2|A| \cdot e^{-2N_1\epsilon_0^2} - 2|S| \cdot |A| \cdot e^{-N_0\lambda_0^2/2} \\ &= 1 - 2|S|^2|A| \cdot e^{-N_0\lambda_0\epsilon_0^2} - 2|S| \cdot |A| \cdot e^{-N_0\lambda_0^2/2} \\ &\geq 1 - 2|S|^2|A| \cdot e^{-N_0\lambda_0^2\epsilon_0^2/2} \\ &= \delta_0, \end{split}$$

as claimed.

## **Step 2.** We claim that assuming

$$||P - \hat{P}||_{\infty} \le \epsilon_0,$$

then for any policy  $\pi$ , we have

$$|R^{(\pi)} - \hat{R}^{(\pi)}| \le T \cdot |S| \cdot R_{\text{max}} \cdot \epsilon_0$$

where  $R^{(\pi)}$  is the expected cumulative reward for  $\pi$  in the MDP M=(S,A,D,P,R,T) and  $\hat{R}^{(\pi)}$  is the expected cumulative reward for  $\pi$  in the MDP  $\hat{M}=(S,A,D,\hat{P},R,T)$ . Note that we have replaced the discount factor  $\gamma$  with the time horizon T. In addition, for all  $z\in Z$ , we have

$$|\rho_z^{(\pi)} - \hat{\rho}_z^{(\pi)}| \le T \cdot |S| \cdot R_{\text{max}} \cdot \epsilon_0,$$

where

$$\rho_z^{(\pi)} = \mathbb{E}_{(s,a) \sim \Lambda_z^{(\pi)}} [\rho_{s,a}],$$

is the expected cumulative agent reward for the MDP M, and  $\hat{\rho}_z^{(\pi)}$  is the expected cumulative agent reward for the MDP  $\hat{M}$ . We only prove the claim for  $|R^{(\pi)} - \hat{R}^{(\pi)}|$ ; the claim for  $|\rho_z^{(\pi)} - \hat{\rho}_z^{(\pi)}|$  follows using the same argument.

Let  $W \in \mathbb{R}^{|S|}$  be

$$W_s = \langle \pi_{s,\cdot}, R_{s,\cdot} \rangle = \sum_{a \in A} \pi_{s,a} R_{s,a}.$$

Then, we have

$$\begin{split} R^{(\pi)} &= \langle R, \Lambda^{(\pi)} \rangle \\ &= \sum_{s \in S} \sum_{a \in A} D_s^{(\pi)} \pi_{s,a} R_{s,a} \\ &= \langle D^{(\pi)}, W \rangle. \end{split}$$

Now, note that

$$D^{(\pi,t)} = (P^{(\pi)})D,$$

so we have

$$D^{(\pi)} = \frac{1}{T} \sum_{t=0}^{T-1} D^{(\pi,t)}$$
$$= \frac{1}{T} \left[ \sum_{t=0}^{T-1} (P^{(\pi)})^t D \right].$$

Thus,

$$R^{(\pi)} = \sum_{t=0}^{T-1} \left\langle (P^{(\pi)})^t D, W \right\rangle.$$

Similarly,

$$\hat{R}^{(\pi)} = \sum_{t=0}^{T-1} \left\langle (\hat{P}^{(\pi)})^t D, W \right\rangle.$$

It follows that

$$R^{(\pi)} - \hat{R}^{(\pi)} = \sum_{t=0}^{T-1} \left\langle (P^{(\pi)})^t D - (\hat{P}^{(\pi)})^t D, W \right\rangle.$$

Thus,

$$|R^{(\pi)} - \hat{R}^{(\pi)}|$$

$$\leq \sum_{t=0}^{T-1} \|(P^{(\pi)})^t D - (\hat{P}^{(\pi)})^t D\|_{\infty} \cdot \|W\|_{1}$$

$$\leq \sum_{t=1}^{T-1} \|P^{(\pi)} - \hat{P}^{(\pi)}\|_{\infty}^t \cdot \|D\|_{\infty} \cdot \|W\|_{1}, \quad (32)$$

where the first line follows from Hölder's inequality and the second line follows from properties of the matrix norm. Also, in the second line, we have used the fact that the summand is zero for t=0 since  $(P^{(\pi)})^0=(\hat{P}^{(\pi)})^0=I$ . Note that

$$||D||_{\infty} \le 1 \tag{33}$$

Furthermore,

$$\begin{aligned} |W_s| &= |\langle \pi_{s,\cdot}, R_{s,\cdot} \rangle| \leq \|\pi_{s,\cdot}\|_1 \cdot \|R_{s,\cdot}\|_{\infty} \\ &\leq \|R_{s,\cdot}\|_{\infty} \\ &\leq R_{\max}, \end{aligned}$$

where the first inequality follows from Hölder's inequality and the second inequality follows since  $\pi_{s,\cdot}$  is a discrete probability distribution. Therefore,

$$||W||_1 = \sum_{s \in S} |W_s| \le |S| \cdot R_{\text{max}}.$$
 (34)

Plugging (33) and (34) into (32) gives

$$|R^{(\pi)} - \hat{R}^{(\pi)}| \le |S| \cdot R_{\text{max}} \cdot \sum_{t=1}^{T-1} \epsilon_0^t$$
$$\le T \cdot |S| \cdot R_{\text{max}} \cdot \epsilon_0.$$

**Step 3.** Now, we prove the theorem. Let  $\hat{\pi}$  be the optimal policy for  $\hat{M}$  (i.e., transitions  $\hat{P}$ ) satisfying  $\hat{\pi} \in \hat{\Pi}_{DP,\epsilon/2}$ . Similarly, let  $\pi^*$  be the optimal policy for M (i.e., transitions P) satisfying  $\pi^* \in \Pi_{DP,\epsilon/4}$ . We apply the second step with

$$\epsilon_0 = \frac{\epsilon}{8T \cdot |S| \cdot R_{\text{max}}}$$
$$\delta_0 = \delta.$$

Then, by the first step, for all  $z, z' \in Z$ , we have

$$\begin{split} & \rho_{z}^{(\hat{\pi})} - \rho_{z'}^{(\hat{\pi})} \\ & \leq (\rho_{z}^{(\hat{\pi})} - \hat{\rho}_{z}^{(\hat{\pi})}) + (\hat{\rho}_{z}^{(\hat{\pi})} - \hat{\rho}_{z'}^{(\hat{\pi})}) + + (\rho_{z'}^{(\hat{\pi})} - \hat{\rho}_{z'}^{(\hat{\pi})}) \\ & \leq T \cdot |S| \cdot R_{\max} \cdot \epsilon_{0} + \frac{\epsilon}{2} + T \cdot |S| \cdot R_{\max} \cdot \epsilon_{0} \\ & \leq \epsilon, \end{split}$$

where the inequality on the third line follows because  $\hat{\pi} \in \hat{\Pi}_{DP,\epsilon/2}$ . Thus, we guarantee that  $\hat{\pi} \in \Pi_{DP,\epsilon}$ .

Next, note that similarly, for all  $z, z' \in Z$ , we have

$$\hat{\rho}_z^{(\pi^*)} - \hat{\rho}_{z'}^{(\pi^*)} \le \frac{\epsilon}{2},$$

so  $\pi^* \in \hat{\Pi}_{\mathrm{DP},\epsilon/2}$ . As a consequence, we have

$$\begin{split} &R^{(\pi^*)} - R^{(\hat{\pi})} \\ &= (R^{(\pi^*)} - \hat{R}^{(\pi^*)}) + (\hat{R}^{(\pi^*)} - \hat{R}^{(\hat{\pi})}) + (\hat{R}^{(\hat{\pi})} - R^{(\hat{\pi})}) \\ &\leq T \cdot |S| \cdot R_{\text{max}} \cdot \epsilon_0 + 0 + T \cdot |S| \cdot R_{\text{max}} \cdot \epsilon_0 \\ &< \epsilon, \end{split}$$

where the inequality on the third line follows because  $\hat{\pi}$  maximizes  $\hat{R}^{(\pi)}$  over  $\pi \in \hat{\Pi}_{DP,\epsilon/2}$ , and  $\pi^* \in \hat{\Pi}_{DP,\epsilon/2}$ .

Thus, the theorem statement follows.  $\Box$ 

### H. Technical Lemmas

**Lemma H.1.** (Hoeffding's inequality) Let  $X \sim p_X$  be a random variable with domain  $[a,b] \subseteq \mathbb{R}$  and mean  $\mu_X$ , and let  $\hat{\mu}_X = n^{-1} \sum_{i=1}^n X_i$  be an estimate of  $\mu_X$  a using n i.i.d. samples  $X_i \sim p_X$ . Then, we have

$$Pr[|\hat{\mu}_X - \mu_X| \ge \epsilon] \le 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right),$$
 (35)

where the probability is taken over the randomness in the i.i.d. samples  $X_1, ..., X_n \sim p_X$ .

*Proof.* See (Wainwright, 2019) for a proof.  $\Box$ 

# I. Experiment Parameters

We use the following parameters for our loan MDP:

$$\begin{split} I &= 0.17318629 \\ p_Z &= 0.29294318 \\ \alpha_{\text{maj}} &= 0.65338681 \\ \beta_{\text{maj}} &= 0.20783559 \\ \alpha_{\text{min}} &= 0.48824268 \\ \beta_{\text{min}} &= 0.48346869 \\ \lambda &= 0.01 \\ \tau &= 0.1 \\ \epsilon &= 0.1 \\ T &= 50 \\ T_{\text{maj}} &= 10 \end{split}$$

 $T_{\min} = 7$ .