

On Welfare-Centric Fair Reinforcement Learning

Cyrus Cousins*
cbcousins@umass.edu

Kavosh Asadi
Amazon

Elita Lobo*
elobo@umass.edu

Michael L. Littman†

Abstract

We propose a novel welfare-centric fair reinforcement-learning setting, in which an agent enjoys *vector-valued* reward from a set of beneficiaries. Given a *welfare function* $W(\cdot)$, the task is to select a policy π that is favorable to all beneficiaries, in the sense that it approximately optimizes the welfare of their value functions from state s_0 , i.e., $\arg\max_{\pi} W(\mathbf{V}_1^{\pi}(s_0), \mathbf{V}_2^{\pi}(s_0), \dots, \mathbf{V}_g^{\pi}(s_0))$. We find that welfare-optimal policies are stochastic and start-state dependent. Whether an individual action is a mistake thus depends on the agent’s policy, therefore mistake bounds, regret analysis, and the PAC-MDP framework do not readily generalize to our setting. We thus develop the *adversarial-fair KWIK* (KWIK-AF) learning model, wherein at each timestep, an agent either takes an *exploration action* or outputs an *exploitation policy*. We require that each exploitation policy be ε -welfare optimal, and the number of exploration actions be bounded. We reduce PAC-MDP to KWIK-AF, introduce the *Equitable Explicit Explore Exploit* (E^4) learner, and show that it is a KWIK-AF learner, thus demonstrating that fair RL is theoretically and practically tractable.

Keywords: Fair RL · Vector-Valued MDP · PAC-MDP · KWIK Learning

1 Introduction

As the negative societal consequences of machine learning (ML) run amok become increasingly apparent, fair ML methods have seen increased attention for tasks like facial recognition (Buolamwini and Gebru, 2018; Cook et al., 2019; Cavazos et al., 2020) and hiring (Kleinberg et al., 2018; Raghavan et al., 2020). Despite this positive trend, most attention on the theory side has been focused on fair supervised (Agarwal et al., 2018; Thomas et al., 2019; Cousins, 2021) and unsupervised (Chierichetti et al., 2017; Chhabra et al., 2021) learning, whereas the second-order societal-welfare impact of ML models, such as the runaway positive feedback loops in settings like predictive policing (Ensign et al., 2018; Alikhademi et al., 2021), are more naturally posed as reinforcement learning (RL) problems.

We apply ideas from welfare-centric supervised learning to the RL setting; in particular, we assume an agent receives a vector-valued reward signal from a set of *beneficiaries*, each representing, e.g., different racial, gender, or religious groups, and the task is to learn a single policy that treats beneficiaries fairly. We argue it is not the role of the algorithm designer to dictate what fairness means in the sense of *how to compromise between beneficiaries*, but rather to optimize for a given fairness notion (ideally one agreed upon by society, government, political philosophers, and other interested parties), as encapsulated by a metric of *societal welfare*. In supervised learning, doing so is relatively straightforward, as we generally maximize the welfare of *expected per-beneficiary value* (Cousins, 2023b), and in our setting, we take utility to be the standard *geometrically discounted reward* (value) for each beneficiary. In general, optimizing welfare is referred to as the *social planner’s problem*, so in a sense our work addresses this problem in the context of RL.

While optimizing the welfare of beneficiary value functions is a well-specified goal for *planning* and *asymptotic learning*, we also ask *how quickly* we can learn to act fairly in an unknown MDP.

*University of Massachusetts Amherst, College of Information and Computer Sciences

†Brown University, Department of Computer Science

Quantifying whether an action is *fair* is substantially more difficult than quantifying whether an action is *optimal* to a single agent, because fairness depends on the *context* of the agent’s policy (i.e., tradeoffs between beneficiaries should be balanced). To address this issue, we combine ideas from the PAC-MDP framework and KWIK (Know What It Knows) learning (Li et al., 2011) to create the *adversarial fair KWIK MDP* learning framework, termed KWIK-AF, which represents a substantially more difficult learning task. We require a KWIK-AF agent to explicitly output at each step either a *fair policy* or an *exploration action*, and our concept of learnability requires that the agent with high probability always outputs ε -optimal fair policies, while taking only a bounded number of exploration actions over its infinite lifetime. For the sake of generality, we allow an adversary to move the agent arbitrarily after it outputs a policy. At any step, the adversary is allowed to select a new welfare function, representing changing societal ideals of how fairness should work, and the agent is expected to output either an exploration action, or a policy optimizing said welfare function. Finally, we introduce an algorithm inspired by the classic E^3 algorithm of Kearns and Singh (2002), which we call *Equitable Explicit Explore Exploit* (E^4), and show that it is a KWIK-AF learner.

We summarize our contributions below.

1. We frame the traditionally egocentric challenge of reinforcement learning as a social problem, where the actions taken by an agent impact a *set of beneficiaries*, each with their own reward function.
2. Using ideas from vector-valued RL, econometrics, and social welfare theory, we establish the goal of learning policies to optimize the *welfare* of per-beneficiary expected discounted rewards.
3. Section 3 introduces the *adversarial fair KWIK MDP* (KWIK-AF) learning framework, in which an agent learns only from *exploration actions*, and an adversary moves the agent when the agent outputs an *exploitation policy*. W.h.p., a learner must output only ε -optimal exploitation policies and take polynomially many exploration actions. We assess *policies* rather than *actions*, since welfare-optimal policies may be *stochastic* and *start-state dependent*, thus actions can not be assessed without context.
4. In section 4, we present the E^4 algorithm, and prove that it is a KWIK-AF learner.

1.1 Related Work

With the rapid increase in the adoption of ML algorithms, authors such as Thomas et al. (2019) note that it is imperative to ensure that such algorithms are well-behaved, and that they do not perpetuate harmful biases. There is thus a need to think about fairness when developing ML algorithms, rather than treating the fairness problem as an afterthought (Kleinberg et al., 2018). Many works study the fairness problem in supervised and unsupervised learning with different — and sometimes inconsistent — fairness definitions. The welfare-centric approach has recently seen success as a generic solution to fair compromise between the wants and needs of various groups, but has thus far been studied primarily in supervised learning (Cousins, 2021; 2022; 2023b; Cousins et al., 2024). Defining fairness in the RL setting is particularly challenging due to the sequential nature of RL decision-makers (Thomas et al., 2019; Jabbari et al., 2017), as we must also decide how fair decisions should be distributed over time.

There is a rich body of literature on multi-objective sequential decision making, which arises naturally in bandit settings (Metevier et al., 2019; Chen et al., 2020), and more generally in planning and RL (Roijers et al., 2013). One approach to fairness is to optimize some objective subject to *fairness constraints*, usually requiring approximate parity between groups. In *contextual bandit settings*, Metevier et al. (2019) learn and plan while (probabilistically) satisfying various fairness constraints. Similarly, Wen et al. (2021) show guarantees for learning and planning in MDPs under parity constraints on per-group value functions, and Satija et al. (2022) generalize their setting by allowing not just rewards, but also the transition function, to differ between groups.

In welfare-centric RL, the final objective is a (nonlinear) function of per-group objectives (value functions). Assuming monotonicity of welfare, the notion of optimal policy is understood to lie on a convex hull or a Pareto frontier (Van Moffaert and Nowé, 2014). Lizotte et al. (2012) show how to identify globally dominated actions in the multi-objective case and in the presence of linear function approximation, and Weng (2019) connects vector-valued reward to welfare objective optimization.

More recently, Siddique et al. (2020); Yu et al. (2023) study the problem of learning fair policies for multi-objective deep RL. Satija et al. (2021) proposed a framework for finding policies that maximize returns while also satisfying certain group fairness constraints. Cousins et al. (2022) consider a tabular similar setting, and Fan et al. (2023) consider a setting similar to ours, where they successfully plan for Nash social welfare, leveraging its differentiability and linearizability. This work theoretically treats fair RL, in particular analyzing *sample complexity*, which has been a core tool for studying exploration in RL (Kearns and Singh, 2002; Kakade, 2003; Li et al., 2011).

1.2 Background

We now present relevant background material on RL, fairness, and welfare-centric ML.

Reinforcement Learning Reinforcement learning (RL) is the study of an environment and an agent that learns to maximize reward through environmental interaction. The Markov decision process (Puterman, 1994), or MDP, is the standard mathematical formalism of RL. In the single-beneficiary case, an MDP is specified by the tuple $\langle \mathcal{S}, \mathcal{A}, R, \mathbf{P}, \gamma \rangle$, where \mathcal{S} is the set of states and \mathcal{A} is the set of actions. The functions $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and $\mathbf{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$, where $\mathcal{P}(X)$ denotes probability distributions over some set X . Finally, γ is called the discount rate, which *geometrically downweights* future rewards, thus incentivizing the agent to seek more near-term rewards.

The standard goal in the RL problem is to learn a policy $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ that can achieve high sums of future discounted rewards. An important concept in RL is the value function, defined as

$$V^\pi(s) \doteq \mathbb{E}_{\substack{a_t \sim \pi(s_t) \\ s_{t+1} \sim \mathbf{P}(s_t, a_t, \cdot)}} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s \right] = \mathbb{E}_{\substack{a_0 \sim \pi(s_0) \\ s_1 \sim \mathbf{P}(s_0, a_0, \cdot)}} \left[R(s_0, \pi(s_0)) + \gamma V^\pi(s_1) \mid s_0 = s \right] .$$

The value function $V^\pi(s)$ describes the *expected utility* of the following policy π at state s . RL often adopts a *egocentric view*, in which the scalar-valued reward function $R(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is intrinsic to the agent, who selfishly wishes to optimize their wellbeing (as measured by the value function).

On Welfare In this paper, we are interested in vector-valued or multi-beneficiary MDPs, denoted $\langle \mathcal{S}, \mathcal{A}, \mathbf{R}, \mathbf{P}, \gamma \rangle$. The state set \mathcal{S} , action set \mathcal{A} , transition function \mathbf{P} , and discount factor γ are exactly as in the standard RL setting. We consciously use the term *beneficiary* to explicitly extricate the *passive nature* of preferences of those impacted by the system from the *active role* played by the *agent*. In this setting, there exist g beneficiaries, each with a corresponding reward function \mathbf{R}_i , value function \mathbf{V}_i^π , and Q function \mathbf{Q}_i^π . In this work, we define the utility of a beneficiary to be the standard RL target of their geometrically discounted accumulated reward (value).

A welfare function $W(\cdot) : \mathbb{R}_{\geq 0}^g \rightarrow \mathbb{R}_{\geq 0}$ then summarizes the *utility* of each of g beneficiaries as a single number, thus establishing a *preference* or *ranking* over possible policies, and the goal is generally to select a policy to *maximize welfare*. For example, the *utilitarian welfare* and *egalitarian welfare* of value vector \mathbf{v} are defined as

$$W_{\text{Util}}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_g) = \frac{1}{g} \sum_{i=1}^g \mathbf{v}_i \quad \text{and} \quad W_{\text{Egal}}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_g) = \min_{i \in \{1, \dots, g\}} \mathbf{v}_i .$$

Utilitarian welfare draws on classical ideas of utilitarian philosophy (Bentham, 1789; Mill, 1863), wherein all members of society should be treated equally, and the goal is to maximize overall utility. On the other hand, *egalitarian welfare* draws from Rawlsian theory (Rawls, 1971; 2001), where the idea is that society should seek to uplift its most disadvantaged (or impoverished) members. Both can be interpreted through a mechanism design or game theoretic lens (Cousins, 2023a), wherein a Dæmon creates a society populated by the beneficiaries, and an Angel then banishes the Dæmon to join the society. If the Angel uniformly randomly selects who the Dæmon becomes, the Dæmon should maximize utilitarian welfare to maximize their expected utility. However, if the Angel adversarially selects the worst-off beneficiary, the Dæmon should instead maximize egalitarian welfare.

Utilitarian welfare is “fair” in the sense that it treats everyone *ostensibly equally*, however it has no preference for *equity*, and can thus incentivize high utility for some beneficiaries at the cost of low utility for others. This often benefits large groups and those who can achieve high utility, while harming disadvantaged or minority groups. On the other hand, egalitarian welfare is “fair” in the sense that it allocates resources optimally to help those most in need first, however it can also be viewed as unfair, as beneficiaries that are difficult or impossible to satisfy may be catered to exclusively, at the expense of all others. A plethora of alternative welfare functions exist in the literature. A spectrum of *prioritarian* welfare concepts (Parfit, 1997; Arneson, 2000) seek a middle ground, incentivizing equity by prioritizing the needs of disadvantaged people, but not to the extreme degree of egalitarianism. We now describe two families, both of which contain utilitarian and egalitarian welfare as extreme cases, as well as a continuum of intermediate cases.

Definition 1.1 (Power-Mean Welfare). *We define the power-mean family $W_p(\mathbf{v})$, for power $p \leq 1$, for any utility vector $\mathbf{v} \in \mathbb{R}_{\geq 0}^g$, as*

$$W_p(\mathbf{v}) \doteq \sqrt[p]{\frac{1}{g} \sum_{i=1}^g v_i^p} \ , \quad W_{-\infty}(\mathbf{v}) \doteq \min_{i \in \{1, \dots, g\}} v_i \ , \quad \text{or} \quad W_0(\mathbf{v}) \doteq \sqrt[g]{\prod_{i=1}^g v_i} \ .$$

Definition 1.2 (Gini Social Welfare). *Given a decreasing stochastic weight vector $\mathbf{w} \in \Delta^g$ (i.e., $\mathbf{w} \in [0, 1]^g$ s.t. $\|\mathbf{w}\|_1 = 1$), the Gini social welfare on \mathcal{M} at policy π from state s_0 is defined as*

$$W_{\mathbf{w}}(\mathcal{M}, \pi) \doteq \sum_{i=1}^g \mathbf{w}_i v_i^\uparrow \ ,$$

where v_i^\uparrow denotes the entries in \mathbf{v}_i in ascending sorted order.

From a prioritarian perspective, these classes make intuitive sense, as the marginal gain of providing additional utility to a beneficiary is larger for low-utility groups than for high-utility groups. This preference for equity is captured by the Pigou (1912)-Dalton (1920) transfer principle, which requires that equitable redistribution of utility has nonnegative impact on welfare. Cardinal welfare theory provides axiomatizations for both the power-mean (Debreu, 1959; Gorman, 1968; Cousins, 2021; 2023b) and Gini (Weymark, 1981; Gajdos and Weymark, 2005) classes.

While these axiomatizations overlap significantly, they are necessarily in conflict. In other words, real humans disagree not only on how to *measure* welfare, but even on how to *axiomatize* welfare. Thus rather than making normative claims as to the “correct” welfare function, we seek to show broad sufficient conditions on welfare under which group-fair RL is possible. For technical reasons, we assume the welfare function must be $\lambda \|\cdot\|_\infty$ Lipschitz continuous, and *concavity* is often convenient for planning, but our methods treat any welfare function that meets these conditions. Cousins (2023b) shows that the power-mean family is Lipschitz continuous except when $p \in [0, 1)$, and the entire Gini family is known to be $1 \|\cdot\|_\infty$ Lipschitz continuous.

In the context of fair RL, our goal is, roughly speaking, to learn a policy π to maximize the welfare of MDP \mathcal{M} (start state s_0 , transition function $\mathbf{P}(\dots)$, reward functions $\mathbf{R}_{1:g}(\dots)$). In other words, we want to find $\hat{\pi}$ to approximate π^* , where

$$W\left(\mathbf{V}_1^{\hat{\pi}}(s_0), \dots, \mathbf{V}_g^{\hat{\pi}}(s_0)\right) \geq W\left(\mathbf{V}_1^{\pi^*}(s_0), \dots, \mathbf{V}_g^{\pi^*}(s_0)\right) - \varepsilon \ .$$

In section 3, we make precise what it means to learn such policies efficiently; in particular, we define how agents interact with their environments, and discuss how much experience an agent may need to make such policy estimates.

2 Illuminating Examples

Here we present a few simple examples (visualized in figure 1) to illustrate that intuition from the standard scalar-reward RL setting can be misleading. We consider the simple *egalitarian welfare* objective (maximize the minimum utility) for two beneficiaries, on essentially stateless

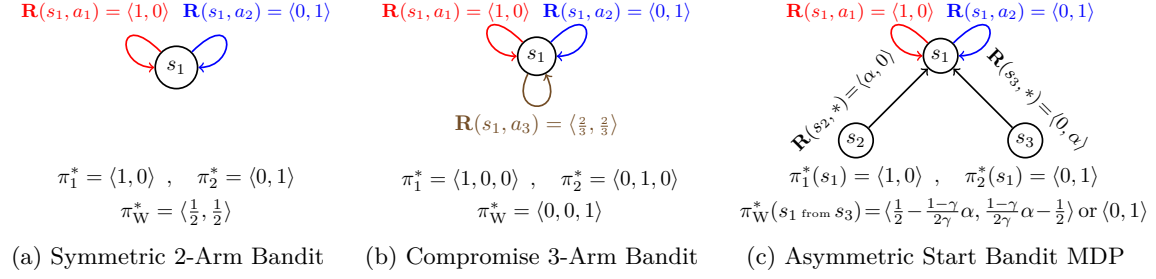


Figure 1: Small MDPs that exhibit surprising behavior under multi-beneficiary objectives.

(single recurrent state) MDPs with deterministic rewards. Even in this elementary setting, we draw surprising conclusions as to the nature of welfare-optimal policies π_W^* (as compared to per-beneficiary optimal policies π_1^* and π_2^*) and the behavior of RL algorithms (both in planning and in exploration). Section 2.1 presents these simple MDPs, and section 2.2 then discusses the challenges of evaluating fair policies (and the learners that produce them).

2.1 Simple Multi-Agent MDPs

We first consider a basic 2-armed bandit, in which the beneficiaries prefer different arms. We then extend our analysis to allow for a third “compromise” arm. Finally, we also allow for additional transient states that immediately reward one of the beneficiaries to represent an “unfair start,” wherein one beneficiary or the other is “privileged” and fair agents must learn to compensate.

One might expect, or at least hope, that convenient properties from standard RL would be preserved in the fair-RL setting. In particular, one might expect the following.

1. We need only consider deterministic stationary policies, i.e., we can assume there always exists an optimal deterministic stationary policy.
2. We can explore by letting individually beneficiaries take turns controlling the agent (thus mitigating potentially challenging learning problems with well-studied techniques).
3. A single policy is optimal from all starting states.

Unfortunately, *none of these properties hold* in the welfare setting. The examples of this section are presented to disabuse the reader of such notions.

Example 2.1 (Symmetric 2-Arm Bandit; Figure 1a). *Suppose a 2-arm bandit, with reward $\mathbf{R}(a_1) = \langle 1, 0 \rangle$, and $\mathbf{R}(a_2) = \langle 0, 1 \rangle$. The unique welfare-optimal stationary policy is $\pi_W^* = \langle \frac{1}{2}, \frac{1}{2} \rangle$.*

There are several surprises here:

1. The (unique) optimal policy is *stationary* (see lemma 3.1), but not *deterministic* (i.e., stochastic).
2. Policy iteration iteratively selects the greedy welfare-optimal policy, i.e., selects the policy

$$\pi^{(t+1)} \leftarrow \operatorname{argmax}_{\pi \in \Pi_{\mathcal{M}}} W \left(\mathbb{E}_{\pi, s_1} [\mathbf{R}_1(s_0, \pi(s_0)) + \gamma \mathbf{V}_1^{\pi^{(t)}}(s_1)], \dots, \mathbb{E}_{\pi, s_1} [\mathbf{R}_g(s_0, \pi(s_0)) + \gamma \mathbf{V}_g^{\pi^{(t)}}(s_1)] \right), \quad (1)$$

where $\pi^{(t)}$ is the policy selected at iterate t . This would optimize the policy in one step if updating the policy did not impact the value function, and this strategy is *convergent* for linear (value) MDP objectives. However, policy iteration for the egalitarian welfare objective, initiated at either deterministic policy, *oscillates* between $\pi(s_1) = \langle 1, 0 \rangle$ and $\pi(s_1) = \langle 0, 1 \rangle$ for any $\gamma > \frac{1}{2}$. This occurs since, assuming $\pi^{(t)}(s_1) = \langle 1, 0 \rangle$, the (stale) value function is $\mathbf{V}^{\pi^{(t)}}(s_1) = \langle \frac{1}{1-\gamma}, 0 \rangle$, thus taking $\pi^{(t+1)}(s_1) = \langle 0, 1 \rangle$ maximizes egalitarian welfare at $\min(\frac{\gamma}{1-\gamma}, 1) = 1$ in (1). In other words, each iteration *overcorrects* for initial policy unfairness, yielding oscillatory behavior. Notably, for $\gamma < \frac{1}{2}$, the oscillation is damped and (1) actually converges to the optimal stochastic policy, but this is case-specific, and policy iteration is not in general a valid planning strategy for welfare objectives.

We now consider an extension to this MDP that includes a third arm (action), which is not preferable to either beneficiary, but is more effective as a compromise than any mixture of the first two options.

Example 2.2 (Compromise 3-Arm Bandit; Figure 1b). *Suppose reward $\mathbf{R}(a_1) = \langle 1, 0 \rangle$, $\mathbf{R}(a_2) = \langle 0, 1 \rangle$, and $\mathbf{R}(a_3) = \langle \frac{2}{3}, \frac{2}{3} \rangle$. Then, the unique optimal stationary per-beneficiary and egalitarian welfare-optimal policies are $\pi_1^* = \langle 1, 0, 0 \rangle$, $\pi_2^* = \langle 0, 1, 0 \rangle$, and $\pi_W^* = \langle 0, 0, 1 \rangle$, respectively.*

This example starkly illustrates how different the welfare-optimal policy π_W may be from the per-beneficiary optimal policies π_1^* and π_2^* . In particular, despite there being unique optimal stationary policies π_1^* , π_2^* , and π_W^* , clearly, π_W^* is not a linear combination of π_1^* and π_2^* . Furthermore, these policies are in some sense totally disjoint, as no two optimal policies will ever take the same action.

This divergence in optimal policies also has implications for how the MDP should be explored. For instance, if beneficiaries 1 and 2 are independently allowed to run a UCB-style algorithm (Auer et al., 2008), in all likelihood, neither will even bother to adequately a_3 , thus even together they do not collect the appropriate information for welfare-optimal planning. We can conclude that, not only the planning, but also the exploration aspect of RL is “more than the sum of its parts,” as under welfare objectives, there is an obligation to explore the MDP more thoroughly.

We now extend the 2-armed bandit example further by adding two additional states, which represent disparate starting conditions that favor one beneficiary or the other.

Example 2.3 (Symmetric 2-Arm Bandit, with Asymmetric Starting Conditions; Figure 1c). *Suppose an MDP with recurrent state s_1 and transient states s_2 and s_3 , thus the environment is a 2-armed bandit upon reaching s_1 . Any action from s_2 yields reward α to beneficiary 1, and any action from s_3 yields reward α to beneficiary 2. Upon reaching state s_1 , the MDP is identical to example 2.1.*

From s_1 , neither beneficiary is privileged, and the recurrent MDP matches example 2.1, but from s_2 and s_3 , some beneficiary begins with an advantage of α utility. To uniquely optimize any non-utilitarian power-mean or Gini welfare, the stochastic stationary policy compensates by selecting $\pi(s_1)$ to benefit the disadvantaged group. Starting at s_3 , we choose a_1 with probability x such that $\frac{\gamma}{1-\gamma}x + \alpha = \frac{\gamma}{1-\gamma}(1-x)$, thus $x = \frac{1}{2} - \frac{1-\gamma}{2\gamma}\alpha$ and perfect equity is achieved, or $x = 0$ if the initial disparity can not be overcome. Starting at s_2 , we instead choose a_2 with probability x by symmetry.

2.2 On Evaluating the Optimality of Fair Learners

When we consider example 2.3, two extremely subtle points arise as to how we are to *evaluate the performance* of a learner and the actions it makes. First, for any $\alpha < \frac{\gamma}{\gamma-1}$, welfare-optimal stationary policies are stochastic at s_1 , (i.e., actions a_1 and a_2 are both taken with nonzero probability). It is thus impossible to determine whether *individual actions* taken by a learner are fair *in isolation*, and a simple mistake-bound style of analysis thus seems inapplicable. This criticism is not unique to our fair RL setting, as it arises whenever it holds that *no deterministic policy is optimal*, for instance in game-theoretic multi-agent RL settings, and also with various constrained learning or risk-averse (e.g., variance-discounted) objectives. In particular, whether a policy is good or not depends on the probability that a given action is taken, not just the action that is taken at some timestep during the execution of the policy.

In settings where stochasticity is optimal, in some cases *statistics* of individual decisions could still be aggregated to assess model performance. An issue more specific to our fair RL setting is that the optimal policy π_W^* depends on the *start state*, so it is meaningless to decompose the learning process into a sequence of individual per-timestep decisions, and evaluate each independently, as this erases the *context* (i.e., as welfare-optimal actions from which starting state) in which decisions are made.

The next section explores these issues further, and derives an appropriate learning model that evaluates agents — not just on individual actions, but on their ability to output nearly welfare-optimal policies. Evaluating fair RL agents is deceptively tricky, particularly due to the contextual nature of start-state dependent policies. Due to the complexity of introducing the context of a starting state, we adopt an *adversarial setting*, in which many design decisions — in particular those

regarding episodic vs. continuous learning, choice of start states or distributions, and the welfare function — are made adversarially. Section 4 then shows that even in this general adversarial setting, fair learning remains possible and algorithmically practical.

3 A Model of Fair Adversarial Reinforcement Learning

In this section, we review the PAC-MDP framework, explain why a straightforward generalization to fairness-sensitive settings is troublesome, and define the KWIK-AF framework. Our guarantees are similar to the classical E^3 policy-centric guarantees of [Kearns and Singh \(2002\)](#), but are adapted to *adversarial* state and welfare-function selection. Both are important to the welfare-centric RL setting, as the adversary can be used to model how policies generated by the agent are actually used (and thus how they impact society), as well as shifts in human fairness concepts over time. Furthermore, we show constructively via reduction that our setting is at least as hard as the PAC-MDP framework. Before further describing the learning setting, we lead with a key lemma that allows us to restrict our attention to start-state dependent stationary MDPs.

Lemma 3.1 (Optimality of Stationary Policies: Lemma 3.1 of [Siddique et al. \(2020\)](#)). *For any start state $s_0 \in \mathcal{S}$, there exists some $W(\cdot)$ -optimal policy*

$$\pi_{s_0}^* \doteq \operatorname{argmax}_{\pi \in \Pi} W(\mathbf{V}_1^\pi(s_0), \dots, \mathbf{V}_g^\pi(s_0))$$

that is a stationary stochastic policy, i.e., given the current state s_t , $\pi_{s_0}^(s_t)$ may prescribe distributions over actions, but they may not depend on the history other than s_0 (i.e., on s_1, s_2, \dots, s_{t-1}).*

3.1 Motivation

Our framework introduces two major ideas. First, we explicitly model the explore-exploit tradeoff by requiring our learners to either take *exploration actions* when uncertain about how to behave, or to output *exploitation policies* when they can near-optimally plan from the current state. In particular, we require that, with high probability, the agent takes a *bounded* (usually polynomial) number of exploration actions, and every exploitation policy is *approximately welfare-optimal*. Second, many decisions in our learning model are made adversarially, and thus our model encompasses a plethora of related settings, including episodic, continual, teacher-assisted, fair, and single-agent RL settings. Consequently, our algorithms and analysis can be directly applied to more specific settings. The central motivation for our policy-centric framework is that simple per-action regret or mistake bounds don’t translate to the fair RL setting. This is because, as discussed in section 2.2 it is not possible to evaluate the optimality of *individual actions* of a fair learner, as they may be stochastic, and they may also depend on the *context* of the start-state s_0 .

Ideally, we could still ensure the agent behaved near-optimally during learning, however, because fair policies are inherently contextual, it also does not make sense to have the learner follow its own policies *at each timestep*, as these policies may disagree, so from where would we even measure suboptimality? While resetting the start state at each step ignores historical context, indefinitely using the agent’s start state puts *too much emphasis* on the past, as from a geometric discounting perspective, we are only planning for optimal behavior in a geometric-length episode, and as time progresses, the start state should become irrelevant in any recurrent MDP. In either case, the agent behaves poorly in some sense during learning; one may consider example 2.3 starting from state s_2 , where keeping the start state indefinitely favors beneficiary 2, whereas resetting it each step favors beneficiary 1 (as their initial privilege is never addressed).

There are many reasonable ways to resolve this issue, but we wish not to limit our framework by committing to one of them. For example, running the agent’s policy for a geometric-length episode before returning control to the agent (to choose either an exploration action, or to output another exploitation policy) would ensure that *behavior* during policy execution is fair. However, even here, reasonable design decisions abound: After a policy execution episode, should we continue from the current state, or start afresh from a new state? If we restart, should the start state be drawn i.i.d., or

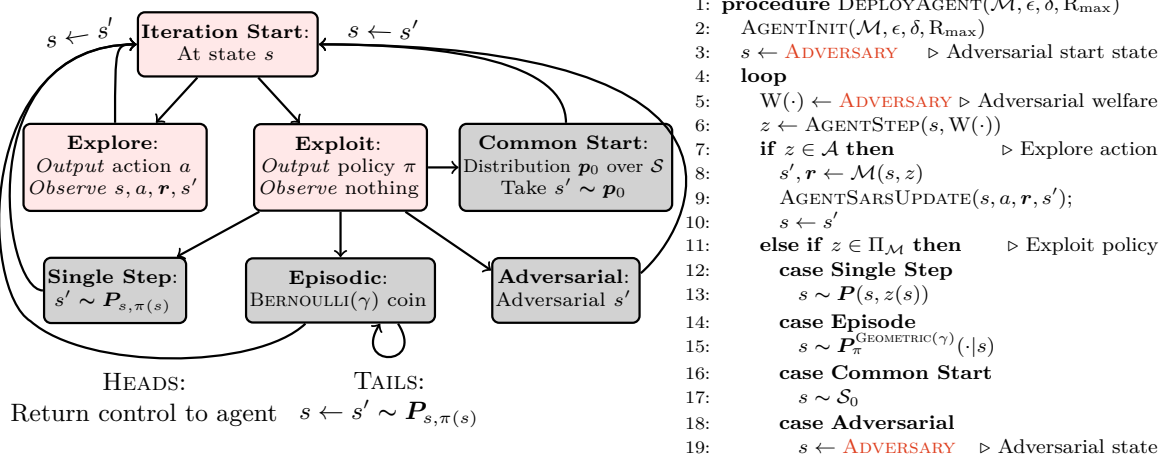


Figure 2: Illustration of MDP Policy Agent control flow. The flowchart (left) describes the control flow represented by the pseudocode (right). An MDP policy agent must implement the `AGENTINIT(...)`, `AGENTSTEP(...)`, and `AGENTSARSUPDATE(...)` subroutines to interact with the environment.

might its distribution change over time? Should the welfare function be fixed, or could it too change over time to reflect evolving societal values or shifting demographics? Rather than adopt some fixed control flow, we require agents to behave near-optimally against a largely-adversarial system.

Essentially, the adversary provides modular flexibility to fairness-sensitive decisions and parameters, and robustness against a learning agent exploiting any possible structure in the learning procedure. This preempts fairness issues arising from a limited model, by requiring that the agent itself must operate under *general* (adversarial) conditions, which the model designer may select to fit a domain-specific ideal of fairness. Furthermore, while *exploitation policies* are guaranteed to be *nearly welfare-optimal*, how they are *actually used* is equally important to fairness. In this context, *adversarial state selection* should be interpreted as taking arbitrary real-world actions informed by the agent’s policy, which should be approximately optimal, before returning control to the agent.

3.2 MDP Policy Agents and the Fair Adversarial KWIK Framework

We now define the MDP policy agent, which codifies how a learner interacts with its environment. This interface is more complicated than standard PAC-MDP learners, because it explicitly models both exploration and exploitation, but this complexity seems necessary to disambiguate good from bad decisions in rich environments where individual actions do not suffice.

Definition 3.2 (MDP Policy Agent). *An MDP policy agent interacts with an MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{R}, \gamma, \mathbf{p}_0 \rangle$ by starting from a start state $s_0 \sim \mathbf{p}_0(\cdot)$. At any timestep t , at state s_t , the agent then produces either an exploration action or an exploitation policy $z \in \mathcal{Z}$ from the space*

$$\mathcal{Z} \doteq \underbrace{\mathcal{A}}_{\text{EXPLORATION ACTION}} \cup \underbrace{\Pi_{\mathcal{M}}}_{\text{EXPLOITATION POLICY}}.$$

If the agent outputs an exploration action a_t , it is executed in state s_t of \mathcal{M} to produce reward $r_{t+1} \sim \mathbf{R}(s_t, a_t)$ and subsequent state $s_{t+1} \sim \mathbf{P}(\cdot | s_t, a_t)$, and the agent observes $\langle s_t, a_t, r_{t+1}, s_{t+1} \rangle$.

Alternatively, if the agent outputs an exploitation policy π_t , a new state s_{t+1} is produced, the agent observes the new state, but the agent does not observe any reward or action. There are many reasonable models for selecting s_{t+1} , and we propose any of:

1. **Single step model:** $s_{t+1} \sim \mathbf{P}(\cdot | s_t, \pi_t(s_t))$;
2. **Episode model:** $s_{t+1} \sim \mathbf{P}_{\pi_t}^k(\cdot | s_t)$ for $\text{GEOMETRIC}(\gamma)$ episode length k ;
3. **Common start model:** Given some start-state distribution \mathbf{p}_0 , we take $s_{t+1} \sim \mathbf{p}_0$; or
4. **Adversarial model:** s_{t+1} is selected adversarially.

Figure 2 illustrates control flow of this system. Note that definition 3.2 resembles the interface of the standard E^3 algorithm, in which an agent takes individual exploration actions until it is ready to output an ε -optimal policy from its current state, at which point it terminates. We require our agents to continue operating after producing a policy, and based on the discussion of section 3.1, we present several reasonable modes of operation following an agent producing an exploitation policy, but all are encompassed by *adversarial choice* of subsequent state. To describe *successful* or *efficient* MDP-policy agents, we define the *policy-KWIK* class, which resembles the KWIK framework for supervised learning (Li et al., 2011), in the sense that the agent is issued “queries” (what to do at the current state), and the agent may either say “I don’t know” to receive information (i.e., issue an exploration action to receive reward and transition samples), or answer the query (give a policy).

Definition 3.3 (Policy-KWIK Learner). *An MDP policy agent is a policy-KWIK learner with sample complexity $m(\varepsilon, \delta, |\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max})$ if, for any error tolerance $\varepsilon > 0$ and failure probability $\delta \in (0, 1)$, the following pair of conditions hold with probability at least $1 - \delta$.*

1. **Exploration condition:** *The number of exploration actions is bounded, i.e.,*

$$\sum_{t=1}^{\infty} \mathbb{1}_{\mathcal{A}}(z_t) \leq m(\varepsilon, \delta, |\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max}) .$$

2. **Exploitation condition:** *All exploitation policies are ε -optimal, i.e., if $z_t \in \Pi_{\mathcal{M}}$, then*

$$\forall t : z_t \in \Pi_{\mathcal{M}} \implies V^{z_t}(s_t) \geq V^*(s_t) - \varepsilon .$$

Definitions 3.2 and 3.3 explicitly delineate between exploration and exploitation; in particular, the agent output space \mathcal{Z} is explicitly factored into *exploration actions*, which are used to take a single step and learn from the environment, and *exploitation policies*, through which the agent demonstrates that it knows how to act near-optimally from its current state. This is codified in conditions 1 and 2 of definition 3.3, as condition 1 requires that an agent may not take too many exploration actions, and condition 2 requires that each policy an agent dares to output must be ε -optimal.

Theorem 3.4 (Policy-KWIK and PAC-MDP Learners). *Every policy-KWIK learner that outputs deterministic policies is a PAC-MDP learner, in the sense that executing an exploitation policy or exploration action at each timestep produces no more than $m(\varepsilon, \delta, |\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max})$ total mistakes (i.e., ε -suboptimal actions) with probability at least $1 - \delta$.*

Proof Sketch. Essentially, this result follows by noting that a policy-KWIK learner can be converted to a PAC-MDP learner by executing each *exploration action* a , or $\pi(s)$ for each *exploitation policy* π at state s , and in doing so, with probability at least $1 - \delta$, no exploitation action is ε -suboptimal. See appendix A for full proof of this result. \square

Group-Fair Models of Reinforcement Learning Definitions 3.2 and 3.3 describe standard (scalar-valued) learning settings, so we now generalize them to definitions 3.5 and 3.6 to model efficient fair learning with welfare objectives for multiple beneficiaries.

Definition 3.5 (Fair Adversarial MDP Policy Agent). *At each timestep t , at state s_t , the adversary presents a welfare function $W_t(\cdot)$ from some class \mathcal{W} . The agent then produces either an exploration action or an exploitation policy $z \in \mathcal{Z}$ from the space*

$$\mathcal{Z} \doteq \underbrace{\mathcal{A}}_{\text{EXPLORATION ACTION}} \cup \underbrace{\Pi_{\mathcal{M}}}_{\text{EXPLOITATION POLICY}} .$$

At this point, if the agent selected an exploration action a_t , the action is executed in state s_t of the MDP to produce reward $\mathbf{r}_{t+1} \sim \mathbf{R}(s_t, a_t)$ and subsequent state $s_{t+1} \sim \mathbf{P}(\cdot \mid s_t, a_t)$, and the agent observes the triplet $\langle s_t, \mathbf{r}_{t+1}, s_{t+1} \rangle$. Alternatively, if the agent selected exploitation policy π_t , the adversary then selects the next state s_{t+1} , and the agent does not observe any reward or action.

Definition 3.6 (KWIK-AF Learner). *An agent is KWIK-AF over welfare class \mathcal{W} with sample complexity $m(\varepsilon, \delta, |\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max}, g)$ if, for any error tolerance $\varepsilon > 0$ and failure probability $\delta \in (0, 1)$, the following pair of conditions hold with probability at least $1 - \delta$:*

1. **Exploration condition:** The number of exploration actions is bounded, i.e.,

$$\sum_{t=1}^{\infty} \mathbb{1}_{\mathcal{A}}(z_t) \leq m(\varepsilon, \delta, |\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max}, g) .$$

2. **Exploitation condition:** All exploitation policies are ε -optimal (with respect to the welfare function $W_t \in \mathcal{W}$ provided by the adversary at each timestep t), i.e., if $z_t \in \Pi_{\mathcal{M}}$, then

$$W_t(\mathbf{V}_1^{z_t}(s_t), \mathbf{V}_2^{z_t}(s_t), \dots, \mathbf{V}_g^{z_t}(s_t)) \geq \sup_{\pi^* \in \Pi_{\mathcal{M}}} W_t(\mathbf{V}_1^{\pi^*}(s_t), \mathbf{V}_2^{\pi^*}(s_t), \dots, \mathbf{V}_g^{\pi^*}(s_t)) - \varepsilon .$$

In other words, the key differences are that reward is now vector-valued, and optimal policies may now be stochastic and must now be welfare-optimal.

4 Algorithms for Fair Planning and Learning

We now present algorithms for fair planning and learning in our multi-beneficiary MDP setting. We first demonstrate how to plan in an MDP to maximize a given concave welfare objective in section 4.1. We then introduce the *Equitable Explicit Explore Exploit* (E^4) fair adversarial MDP policy agent (algorithm 1) in section 4.2. Finally, we bound the sample complexity of E^4 and show that it is a KWIK-AF learner in section 4.3.

4.1 On Welfare-Optimal Planning

For a given start-state distribution $\mathbf{p}_0 \in \Delta^{\mathcal{S}}$, let $\mathbf{d}^{\pi} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ denote the geometrically-discounted state-action occupancy measure of the policy π , defined as

$$\mathbf{d}_{s,a}^{\pi} \doteq \mathbb{E}_{\substack{a_t \sim \pi(s_t) \\ s_0 \sim \mathbf{p}_0 \\ s_{t+1} \sim P(s_t, a_t, \cdot)}} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{1}_s(s_t) \mathbb{1}_a(a_t) \right] = \pi(s, a) \left(\mathbf{p}_0(s) + \gamma \sum_{\substack{s' \in \mathcal{S} \\ a' \in \mathcal{A}}} P(s', a', s) \mathbf{d}_{s',a'}^{\pi} \right) . \quad (2)$$

We now apply the state-action occupancy measure to the welfare-optimal planning problem.

Proposition 4.1 (Welfare-Optimal Planning). *For a concave welfare function $W(\cdot)$, the welfare-optimal policy $\pi^* = \operatorname{argmax}_{\pi \in \Pi_{\mathcal{M}}} W(\mathbf{V}^{\pi}(s))$ can be computed by first solving*

$$\begin{aligned} \mathbf{d}^* = \operatorname{argmax}_{\substack{\mathbf{d} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \\ \mathbf{d} \geq 0}} W \left(\sum_{s \in \mathcal{S}, a \in \mathcal{A}} \mathbf{d}_{s,a} \mathbf{R}_1(s, a), \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \mathbf{d}(s, a) \mathbf{R}_2(s, a), \dots, \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \mathbf{d}_{s,a} \mathbf{R}_g(s, a) \right) \quad (3) \\ \text{such that } \forall s \in \mathcal{S} : \sum_{a \in \mathcal{A}} \mathbf{d}_{s,a} = \mathbf{p}_0(s) + \gamma \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}} P(s', a', s) \mathbf{d}_{s',a'} , \end{aligned}$$

and then setting $\pi^*(s, a) = \frac{\mathbf{d}_{s,a}^*}{\sum_{a' \in \mathcal{A}} \mathbf{d}_{s,a'}^*}$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$.

Note that an ε -optimal policy π_{ε} for (3) can be identified using standard convex optimization techniques, e.g., subgradient ascent, in $\operatorname{Poly}(|\mathcal{S}|, |\mathcal{A}|, \varepsilon)$ time. See appendix B for proof of this claim, as well as the details and interesting special cases of the resulting optimization problems.

4.2 The E^4 Algorithm

We first briefly describe the E^4 algorithm, for which we give pseudocode in algorithm 1. The key to understanding E^4 is that the state space is divided into three sets: The unknown set \mathcal{S}_{unk} , the outer-known set \mathcal{S}_{out} , and the inner-known set \mathcal{S}_{inn} . Initially, all states are unexplored, and thus in \mathcal{S}_{unk} (line 8). After visiting a state $s \in \mathcal{S}$ and taking all actions sufficiently many times (line 24) using balanced wandering (line 13), s becomes *known*, entering either \mathcal{S}_{inn} or \mathcal{S}_{out} . We then construct an *empirical MDP* $\hat{\mathcal{M}}$ using the empirical transition frequencies and average reward from each known state and each action (lines 25 and 26), and self-loop probability 1 and reward $\mathbf{0}$ for all unknown

Algorithm 1 Equitable Explicit Explore Exploit (E^4)

```

1: procedure AGENTINIT( $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathbf{R}, \mathbf{P}, \gamma \rangle, \epsilon, \delta, R_{\max}$ )
2:    $T \leftarrow \left\lceil \log_{\frac{1}{\gamma}} \left( \frac{6\lambda R_{\max}}{\epsilon(1-\gamma)} \right) \right\rceil$  ▷ Set escape time
3:    $\alpha \leftarrow \frac{2\epsilon(1-\gamma)^2}{3\lambda\sqrt{R_{\max}}(2+\gamma\sqrt{R_{\max}}+6T(1-\gamma)\sqrt{R_{\max}})}$  ▷ Set error tolerance  $\alpha$  for transitions
4:    $\beta \leftarrow \alpha\sqrt{R_{\max}}$  ▷ Set error tolerance  $\beta$  for rewards
5:    $E \leftarrow 2\alpha T$  ▷ Initialize escape threshold  $E$ 
6:    $t \leftarrow 0$  ▷ Escape timer  $t$ 
7:    $M \leftarrow \left\lceil \ln \left( \frac{2|\mathcal{S}||\mathcal{A}|(2^{|\mathcal{S}|} - 2 + 2g)}{\delta} \right) \max \left( \frac{1}{2\beta^2}, \frac{R_{\max}^2}{2\alpha^2} \right) \right\rceil$  ▷ Compute sufficient per-state-action pair sample size
8:    $\forall s \in \mathcal{S}, a \in \mathcal{A} : m_{s,a} \leftarrow 0$  ▷ Initialize per- $(s, a)$  visitation counters
9:    $\mathcal{S}_{\text{unk}} \leftarrow \mathcal{S}; \mathcal{S}_{\text{out}} \leftarrow \emptyset; \mathcal{S}_{\text{inn}} \leftarrow \emptyset$  ▷ Initialize all states to unknown
10:   $\hat{\mathcal{M}} = \langle \mathcal{S}, \mathcal{A}, \hat{\mathbf{R}}, \hat{\mathbf{P}}, \gamma \rangle \leftarrow \langle \mathcal{S}, \mathcal{A}, (s, a) \mapsto \mathbf{0}, s \mapsto \mathbb{1}_s, \gamma \rangle$  ▷ Initialize empirical MDP  $\hat{\mathcal{M}}$  to  $\mathbf{0}$ -reward recurrent states
11:  procedure AGENTSTEP( $s, W(\cdot)$ )
12:    case  $s \in \mathcal{S}_{\text{unk}}$  ▷ Successful escape attempt has reached  $\mathcal{S}_{\text{unk}}$ 
13:       $t \leftarrow 0$ ; return  $a_{\text{xpr}} \leftarrow \underset{a \in \mathcal{A}}{\operatorname{argmin}} m_{s,a}$  ▷ Balanced walk step
14:    case  $t > 0$  ▷ Ongoing attempt to escape to  $\mathcal{S}_{\text{unk}}$ 
15:       $t \leftarrow t - 1$ ; return  $a_{\text{xpr}} \leftarrow \pi_{\text{esc}}(s, t)$  ▷ Explore  $a$  from escape  $\pi$ 
16:    case  $s \in \mathcal{S}_{\text{inn}}$  ▷ Return exploit policy
17:      return  $\pi_{\text{xpt}} \leftarrow \underset{\pi \in \Pi_{\hat{\mathcal{M}}}}{\operatorname{argmax}} W(\hat{\mathbf{V}}^{\pi}(s))$ 
18:    case  $s \in \mathcal{S}_{\text{out}}$  ▷ Begin escape attempt
19:       $t \leftarrow T$ ; return  $a_{\text{xpr}} \leftarrow \pi_{\text{esc}}(s, t)$ 
20:  procedure AGENTSARSUPDATE( $s, a, r, s'$ )
21:    if  $s \in \mathcal{S}_{\text{unk}}$  then
22:       $m_{s,a} \leftarrow m_{s,a} + 1$  ▷ Increment visitation count
23:       $(\mathbf{E}_{s,a,m_{s,a},s}, \mathbf{E}_{s,a,m_{s,a},r}) \leftarrow (s', r)$  ▷ Add to experience buffer
24:      if  $\min_{a \in \mathcal{A}} m_{s,a} = M$  then ▷ State  $s$  is learned
25:         $\forall a \in \mathcal{A}, s' \in \mathcal{S} : \hat{\mathbf{P}}_{s,a,s'} \leftarrow \frac{1}{M} \sum_{i=1}^M \mathbb{1}_{s'}(\mathbf{E}_{s,a,i,s})$  ▷ Empirical transition model  $\hat{\mathbf{P}}$ 
26:         $\forall a : \hat{\mathbf{R}}_{s,a} \leftarrow \frac{1}{M} \sum_{i=1}^M (\mathbf{E}_{s,a,i,r})$  ▷ Empirical reward function  $\hat{\mathbf{R}}$ 
27:         $\pi_{\text{esc}} \leftarrow \underset{\pi \in \Pi_T}{\operatorname{argmax}} \sum_{s \in \mathcal{S}} \mathbb{P} \left( \bigvee_{i=1}^T s_i \in \mathcal{S}_{\text{unk}} \mid \pi, s_1 = s, s_t = \hat{\mathbf{P}}s_{t-1} \right)$  ▷  $T$ -step deterministic escape policy using  $\hat{\mathbf{P}}$ 
28:         $\mathcal{S}_{\text{unk}} \leftarrow \mathcal{S}_{\text{unk}} \setminus \{s\}$  ▷ Remove  $s$  from the unknown set
29:         $\mathcal{S}_{\text{out}} \leftarrow \left\{ s \in (\mathcal{S} \setminus \mathcal{S}_{\text{unk}}) \mid \mathbb{P} \left( \bigvee_{i=1}^T s_i \in \mathcal{S}_{\text{unk}} \mid \pi_{\text{esc}}, s_0 = s \right) \geq E \right\}$  ▷ Known states s.t.  $T$ -step escape is  $E$ -likely
30:         $\mathcal{S}_{\text{inn}} \leftarrow \mathcal{S} \setminus (\mathcal{S}_{\text{unk}} \cup \mathcal{S}_{\text{out}})$  ▷ Known states s.t.  $T$ -step escape is not  $E$ -likely

```

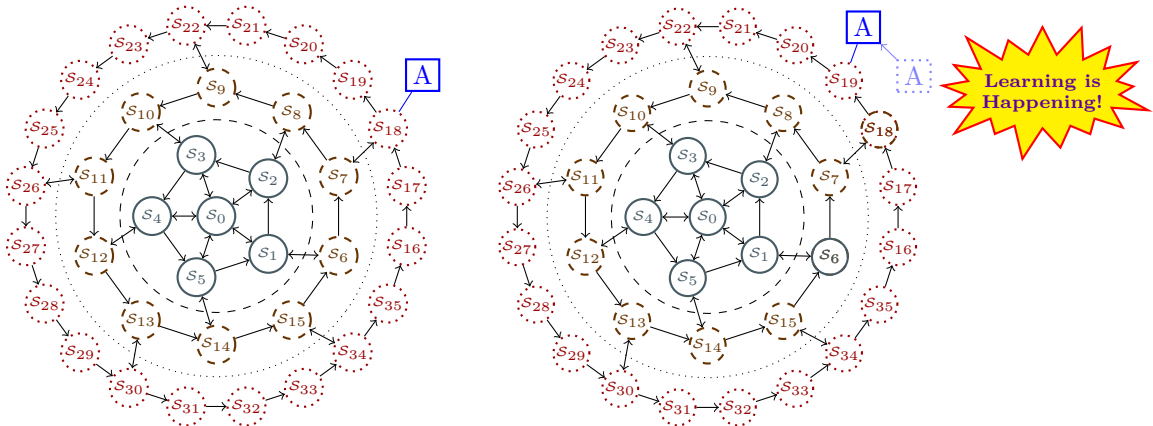


Figure 3: E^4 agent **A** learning on MDP \mathcal{M} . Inner-known set \mathcal{S}_{inn} (solid), outer-known set \mathcal{S}_{out} (dashed), and unknown set \mathcal{S}_{unk} (dotted) illustrated. Assume all actions self-loop with probability $\approx \tau/\sqrt{1-E}$, and let arrows denote 1-step reachability via some action with probability $\approx 1 - \tau/\sqrt{1-E}$. As the agent acts (explores) to reach s_{19} from s_{18} , s_{18} enters \mathcal{S}_{out} , which cascades to s_6 entering \mathcal{S}_{inn} .

states (line 10). Then, if with nonnegligible probability (at least E) in $\hat{\mathcal{M}}$ it is possible to reach \mathcal{S}_{unk} from s within T steps, we place s into \mathcal{S}_{out} (line 29), otherwise we place s into \mathcal{S}_{inn} (line 30). A state s becoming known may also cascade into states $s' \in \mathcal{S}_{\text{out}}$ entering \mathcal{S}_{inn} . This process is graphically illustrated in figure 3. Note that E and T are set so as to ensure E^4 is KWIK-AF (line 5 and 6).

As in the classic E^3 algorithm, within \mathcal{S}_{inn} , if all tail bounds hold simultaneously, the value functions of $\hat{\mathcal{M}}$ approximate the value functions of \mathcal{M} . Furthermore, under Lipschitz continuity of welfare, optimizing welfare in $\hat{\mathcal{M}}$ approximately optimizes welfare in \mathcal{M} . Therefore, at each step, if the agent is in \mathcal{S}_{inn} , it outputs a near-optimal policy (line 16). Otherwise, if the agent is in \mathcal{S}_{out} , it begins an *escape attempt* (line 18), which follows a T -step temporal policy that maximizes the probability of reaching \mathcal{S}_{unk} in $\hat{\mathcal{M}}$ (line 27). The escape attempt either proceeds for T steps (line 14), or until \mathcal{S}_{unk} is reached (line 12) and subsequently explored. The main concrete difference between the classical E^3 and our E^4 is that E^4 has higher sample complexity, due both to vector-valued reward and to the nonlinearity of the welfare function. Furthermore, our analysis is more complex, as we show that E^4 KWIK-AF learns \mathcal{M} , which requires robustness against adversarial state selection even after an infinite number of exploitation steps, whereas the classical E^3 analysis only guarantees a single ε -optimal exploitation policy is output.

4.3 Theoretical Analysis

We are now ready to theoretically analyze E^4 in the context of the KWIK-AF framework. We begin by defining an (α, β) uniform approximation of an MDP \mathcal{M} and compute the per-state sample complexity of the approximation. We then derive the sample complexity of the KWIK-AF framework and show that it is polynomial in all relevant parameters.

Definition 4.2 (Uniform Approximation MDPs). *Let $\text{TVD}(x, y)$ denote the total variation distance between probability distributions x, y . An (α, β) uniform approximation $\mathcal{M}' = \langle \mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{R}', \gamma \rangle$ of a vector-reward MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{R}, \gamma \rangle$ is an MDP that, for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$, satisfies*

1. $\text{TVD}(\mathbf{P}'(\cdot|s, a), \mathbf{P}(\cdot|s, a)) \leq \alpha$; and
2. $\|\mathbf{R}'(s, a) - \mathbf{R}(s, a)\|_\infty \leq \beta$.

Lemma 4.3 (Per-State Sample Complexity). *Suppose MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathbf{R}, \mathbf{P}, \gamma \rangle$, and let*

$$m_{\text{knw}} \doteq \left\lceil \ln \left(\frac{|\mathcal{S}||\mathcal{A}| (2^{|\mathcal{S}|} - 2 + 2g)}{\delta} \right) \max \left(\frac{1}{2\alpha^2}, \frac{R_{\max}}{2\beta^2} \right) \right\rceil, \quad (4)$$

where $R_{\max} \doteq \max_{s \in \mathcal{S}, a \in \mathcal{A}, i \in 1, \dots, g} \mathbf{R}_i(s, a) \in [0, \infty]$. Now if $\hat{\mathcal{M}}$ is estimated from m_{knw} samples of each state-action pair, then, with probability at least $1 - \delta$, $\hat{\mathcal{M}}$ is an α - β approximation of \mathcal{M} .

Theorem 4.4 (E^4 is KWIK-AF). *Algorithm 1 is a KWIK-AF learner that learns \mathcal{M} w.r.t. the class of all $\lambda \|\cdot\|_\infty$ Lipschitz welfare functions, with sample complexity*

$$m(\varepsilon, \delta, |\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max}, g) \in \text{Poly} \left(|\mathcal{S}|, |\mathcal{A}|, \log g, R_{\max}, \lambda, \frac{1}{\varepsilon}, \log \frac{1}{\delta}, \frac{1}{1 - \gamma} \right).$$

5 Conclusion

This work motivates and defines a formal model of welfare-centric fair reinforcement learning. We find that naïve approaches, like planning via policy iteration (example 2.1), and independent per-beneficiary exploration (example 2.2) do not yield fair RL agents. Defining fair RL and quantifying a learner’s efficiency are challenging problems (section 3), as we must consider stochastic policies, and thus can not evaluate learners in terms of the *regret* or *mistakes of individual actions*. We thus define the *Fair Adversarial MDP Policy Agent* (definition 3.5) and the KWIK-AF Learner (definition 3.6) to model fair RL and codify efficient learning in this domain. We then show (section 4) that under mild regularity conditions on the welfare function, it is possible to learn in the KWIK-AF framework while making polynomially many mistakes (algorithm 1, theorem 4.4). Our method adapts the classic E^3 algorithm, which is an appropriate fit, as its exploration strategy is actually independent of the reward function. As a result, the only major change we require is attempting each action more often during exploration to account for the larger number of parameters that must be learned.

References

- A. Agarwal, A. Beygelzimer, M. Dudík, J. Langford, and H. Wallach. A reductions approach to fair classification. In *International Conference on Machine Learning*, pages 60–69. PMLR, 2018.
- A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun. *Reinforcement Learning: Theory and Algorithms*. 2022.
- K. Alikhademi, E. Drobina, D. Prioleau, B. Richardson, D. Purves, and J. E. Gilbert. A review of predictive policing from the perspective of fairness. *Artificial Intelligence and Law*, pages 1–17, 2021.
- R. J. Arneson. Luck egalitarianism and prioritarianism. *Ethics*, 110(2):339–349, 2000.
- P. Auer, T. Jaksch, and R. Ortner. Near-optimal regret bounds for reinforcement learning. *Advances in neural information processing systems*, 21, 2008.
- J. Bentham. An introduction to the principles of morals and legislation. *University of London: the Athlone Press*, 1789.
- J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.
- J. G. Cavazos, P. J. Phillips, C. D. Castillo, and A. J. O’Toole. Accuracy comparison across face recognition algorithms: Where are we on measuring race bias? *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2020.
- Y. Chen, A. Cuellar, H. Luo, J. Modi, H. Nemlekar, and S. Nikolaidis. Fair contextual multi-armed bandits: Theory and experiments. In *Conference on Uncertainty in Artificial Intelligence*, pages 181–190. PMLR, 2020.
- A. Chhabra, K. Masalkovaitė, and P. Mohapatra. An overview of fairness in clustering. *IEEE Access*, 9: 130698–130720, 2021.
- F. Chierichetti, R. Kumar, S. Lattanzi, and S. Vassilvitskii. Fair clustering through fairlets. *Advances in neural information processing systems*, 30, 2017.
- C. M. Cook, J. J. Howard, Y. B. Sirotin, J. L. Tipton, and A. R. Vemury. Demographic effects in facial recognition and their dependence on image acquisition: An evaluation of eleven commercial systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 1(1):32–41, 2019.
- C. Cousins. An axiomatic theory of provably-fair welfare-centric machine learning. In *Advances in Neural Information Processing Systems*, 2021.
- C. Cousins. Uncertainty and the social planner’s problem: Why sample complexity matters. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 2022.
- C. Cousins. Algorithms and analysis for optimizing robust objectives in fair machine learning. In *Columbia Workshop on Fairness in Operations and AI*. Columbia University, 2023a.
- C. Cousins. Revisiting fair-PAC learning and the axioms of cardinal welfare. In *Artificial Intelligence and Statistics (AISTATS)*, 2023b.
- C. Cousins, K. Asadi, and M. L. Littman. Fair E³: Efficient welfare-centric fair reinforcement learning. In *5th Multidisciplinary Conference on Reinforcement Learning and Decision Making (RLDM)*, 2022.
- C. Cousins, I. E. Kumar, and S. Venkatasubramanian. To pool or not to pool: Analyzing the regularizing effects of group-fair training on shared models. In *Artificial Intelligence and Statistics (AISTATS)*, 2024.
- H. Dalton. The measurement of the inequality of incomes. *The Economic Journal*, 30(119):348–361, 1920.
- G. Debreu. Topological methods in cardinal utility theory. *Cowles Foundation Discussion Papers*, 76, 1959.
- D. Ensign, S. A. Friedler, S. Neville, C. Scheidegger, and S. Venkatasubramanian. Runaway feedback loops in predictive policing. In *Conference on Fairness, Accountability and Transparency*, pages 160–171. PMLR, 2018.
- Z. Fan, N. Peng, M. Tian, and B. Fain. Welfare and fairness in multi-objective reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pages 1991–1999, 2023.
- T. Gajdos and J. A. Weymark. Multidimensional generalized Gini indices. *Economic Theory*, 26(3):471–496, 2005.
- W. M. Gorman. The structure of utility functions. *The Review of Economic Studies*, 35(4):367–390, 1968.
- S. Jabbari, M. Joseph, M. Kearns, J. Morgenstern, and A. Roth. Fairness in reinforcement learning. In *International conference on machine learning*, pages 1617–1626. PMLR, 2017.
- S. M. Kakade. *On the sample complexity of reinforcement learning*. University of London, University College London (United Kingdom), 2003.

- M. Kearns and S. Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2): 209–232, 2002.
- J. Kleinberg, J. Ludwig, S. Mullainathan, and A. Rambachan. Algorithmic fairness. In *AEA papers and proceedings*, volume 108, pages 22–27, 2018.
- L. Li, M. L. Littman, T. J. Walsh, and A. L. Strehl. Knows what it knows: A framework for self-aware learning. *Machine learning*, 82(3):399–443, 2011.
- D. J. Lizotte, M. Bowling, and S. A. Murphy. Linear fitted- q iteration with multiple reward functions. *The Journal of Machine Learning Research*, 13(1):3253–3295, 2012.
- B. Metevier, S. Giguere, S. Brockman, A. Kobren, Y. Brun, E. Brunskill, and P. S. Thomas. Offline contextual bandits with high probability fairness guarantees. *Advances in neural information processing systems*, 32, 2019.
- J. S. Mill. *Utilitarianism*. Parker, Son, and Bourn, London, 1863.
- D. Parfit. Equality and priority. *Ratio (Oxford)*, 10(3):202–221, 1997.
- A. C. Pigou. *Wealth and welfare*. Macmillan and Company, limited, 1912.
- M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. USA, 1st edition, 1994. ISBN 0471619779.
- M. Raghavan, S. Barocas, J. Kleinberg, and K. Levy. Mitigating bias in algorithmic hiring: Evaluating claims and practices. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pages 469–481, 2020.
- J. Rawls. *A theory of justice*. Harvard University Press, 1971.
- J. Rawls. *Justice as fairness: A restatement*. Harvard University Press, 2001.
- D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley. A survey of multi-objective sequential decision-making. *Journal of Artificial Intelligence Research*, 48:67–113, 2013.
- H. Satija, P. S. Thomas, J. Pineau, and R. Laroché. Multi-objective SPIBB: Seldonian offline policy improvement with safety constraints in finite MDPs. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- H. Satija, A. Lazaric, M. Pirotta, and J. Pineau. Group fairness in reinforcement learning. *Transactions on Machine Learning Research*, 2022.
- U. Siddique, P. Weng, and M. Zimmer. Learning fair policies in multi-objective (deep) reinforcement learning with average and discounted rewards. In *International Conference on Machine Learning*, pages 8905–8915. PMLR, 2020.
- A. L. Strehl, L. Li, and M. L. Littman. Reinforcement learning in finite MDPs: PAC analysis. *Journal of Machine Learning Research*, 10(11), 2009.
- P. S. Thomas, B. Castro da Silva, A. G. Barto, S. Giguere, Y. Brun, and E. Brunskill. Preventing undesirable behavior of intelligent machines. *Science*, 366(6468):999–1004, 2019.
- K. Van Moffaert and A. Nowé. Multi-objective reinforcement learning using sets of Pareto dominating policies. *The Journal of Machine Learning Research*, 15(1):3483–3512, 2014.
- T. Wang, D. Lizotte, M. Bowling, and D. Schuurmans. Dual representations for dynamic programming. *Journal of Machine Learning Research*, 1:1–29, 01 2008.
- M. Wen, O. Bastani, and U. Topcu. Algorithms for fairness in sequential decision making. In *International Conference on Artificial Intelligence and Statistics*, pages 1144–1152. PMLR, 2021.
- P. Weng. Fairness in reinforcement learning. *arXiv preprint arXiv:1907.10323*, 2019.
- J. A. Weymark. Generalized Gini inequality indices. *Mathematical Social Sciences*, 1(4):409–430, 1981.
- G. Yu, U. Siddique, and P. Weng. Fair deep reinforcement learning with generalized Gini welfare functions. 2023.
- T. Zahavy, B. O’Donoghue, G. Desjardins, and S. Singh. Reward is enough for convex MDPs. *Advances in Neural Information Processing Systems*, 34:25746–25759, 2021.